

ББК 32.973.26-018.2.75

C50

УДК 681.3.07

Издательский дом “Вильямс”

Зав. редакцией *С.Н. Тригуб*

Перевод с английского *Р.Г. Имамутдиновой*, канд. физ.-мат. наук *А.А. Минько*

Под редакцией канд. физ.-мат. наук *А.А. Минько*

По общим вопросам обращайтесь в Издательский дом “Вильямс” по адресу:

info@williamspublishing.com, <http://www.williamspublishing.com>

115419, Москва, а/я 783; 03150, Киев, а/я 152

Смит, Билл.

C50 Методы и алгоритмы вычислений на строках. : Пер. с англ. — М. : ООО “И.Д. Вильямс”, 2006. — 496 с. : ил. — Парал. тит. англ.

ISBN 5–8459–1081–1 (рус.)

Книга представляет собой фундаментальное введение в алгоритмы и методы, эффективно вычисляющие паттерны в строковых последовательностях. Речь идет об общих алгоритмах и методах, которые находят применение во многих областях науки и информационных технологий: сжатие данных, криптография, распознавание речи и компьютерное зрение, вычислительная геометрия и молекулярная биология. Рассмотренные в книге алгоритмы предназначены для нахождения в строковых последовательностях определенных типов паттернов — частных, характеристических и внутренних. Каждому типу паттернов посвящена соответствующая часть книги. Книга отличается последовательным изложением материала, большим количеством иллюстративных примеров, свободным обсуждением текущих исследований в этой области, содержит более 500 упражнений, поясняющих и расширяющих материал, изложенный в тексте книги.

Книга предназначена для тех, кто имеет достаточную подготовку в математике и информатике и хочет познакомиться с этой интересной и важной областью. Материал книги доступен для студентов старших курсов и аспирантов соответствующих специальностей.

ББК 32.973.26–018.2.75

Все названия программных продуктов являются зарегистрированными торговыми марками соответствующих фирм.

Никакая часть настоящего издания ни в каких целях не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами, будь то электронные или механические, включая фотокопирование и запись на магнитный носитель, если на это нет письменного разрешения издательства Addison–Wesley UK.

Authorized translation from the English language edition published by Addison–Wesley UK, Copyright © 2003.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from the Publisher.

Russian language edition published by Williams Publishing House according to the Agreement with R&I Enterprises International, Copyright © 2006

ISBN 5–8459–1081–1 (рус.)
ISBN 0–201–39839–7 (англ.)

© Издательский дом “Вильямс”, 2006
© Addison–Wesley UK, 2003

ОГЛАВЛЕНИЕ

Предисловие	10
Часть I. Строковые последовательности и алгоритмы	17
Глава 1. Свойства строковых последовательностей	19
Глава 2. Паттерны? Что такое паттерны?	56
Глава 3. Такие разные строки	85
Глава 4. Строковые алгоритмы и тестовые данные	113
Часть II. Вычисление внутренних паттернов	137
Глава 5. Деревья для строковых последовательностей	139
Глава 6. Декомпозиция строковых последовательностей	191
Часть III. Вычисление частных паттернов	215
Глава 7. Базовые алгоритмы	217
Глава 8. Наследники Бойера–Мура	245
Глава 9. Алгоритмы вычисления расстояния между строками	278
Глава 10. Приближенное сравнение с паттерном	309
Глава 11. Регулярные выражения и множественные паттерны	343
Часть IV. Вычисление характеристических паттернов	381
Глава 12. Периодичность	383
Глава 13. Обобщение периодичности	417
Литература	467
Предметный указатель	483

СОДЕРЖАНИЕ

Предисловие	10
Часть I. Строковые последовательности и алгоритмы	17
Глава 1. Свойства строковых последовательностей	19
1.1 Последовательность жемчужин	19
1.2 Линейные строковые последовательности	22
1.3 Периодичность	33
1.4 Строковые петли	45
Глава 2. Паттерны? Что такое паттерны?	56
2.1 Внутренние паттерны	56
2.2 Частные паттерны	64
2.3 Характеристические паттерны	75
Глава 3. Такие разные строки	85
3.1 Проблема исключений и морфизмы	85
3.2 Строки Туе, являющиеся (2, 3)-исключениями	89
3.3 Строки Туе, являющиеся (3, 2)-исключениями	97
3.4 Строки Фибоначчи	100
Глава 4. Строковые алгоритмы и тестовые данные	113
4.1 Хорошие строковые алгоритмы	113
4.2 Разные паттерны	120
4.3 Разные грани	127

Часть II. Вычисление внутренних паттернов	137
Глава 5. Деревья для строковых последовательностей	139
5.1 Деревья граней	139
5.2 Деревья суффиксов	141
5.2.1 Предварительные сведения о деревьях суффиксов	142
5.2.2 Алгоритм Мак-Крейта	146
5.2.3 Алгоритм Укконена	150
5.2.4 Алгоритм Фараха	156
5.2.5 Применение и реализация	167
5.3 Альтернативные структуры для представления суффиксов	172
5.3.1 Ориентированные ациклические графы слов	172
5.3.2 Массивы суффиксов	181
Глава 6. Декомпозиция строковых последовательностей	191
6.1 Линдонская декомпозиция: алгоритм Дюваля	192
6.2 Применения линдонской декомпозиции	203
6.3 <i>s</i> -факторизация	211
Часть III. Вычисление частных паттернов	215
Глава 7. Базовые алгоритмы	217
7.1 Алгоритм Кнута–Морриса–Пратта	217
7.2 Алгоритм Бойера–Мура	224
7.3 Алгоритм Карпа–Рабина	236
7.4 Алгоритм Демелки–Бейза–Ятса–Гоннета	240
7.5 Заключение	244
Глава 8. Наследники Бойера–Мура	245
8.1 Алгоритм БМ — цикл с перескоками	246
8.2 Алгоритм Бойера–Мура–Хоспула	248
8.3 Частота встречаемости букв и алгоритмы Бойера–Мура–Санди	250
8.3.1 Сравнение со стражем	250
8.3.2 Наиболее эффективные перескоки	251
8.3.3 Первый алгоритм Бойера–Мура–Санди	252
8.3.4 Второй алгоритм Бойера–Мура–Санди	254
8.4 Алгоритм Бойера–Мура–Гелила	258
8.5 Алгоритм Турбо-БМ	263
8.6 Наследники алгоритма КМП	266
8.7 Создайте собственный алгоритм!	271
8.8 Точные оценки сложности алгоритмов точного сравнения с паттерном	274

Глава 9. Алгоритмы вычисления расстояния между строками	278
9.1 Базисная рекурсия	280
9.2 Алгоритм Вагнера–Фишера	283
9.3 Алгоритмы Хешберга	287
9.4 Алгоритм Ханта–Шиманского	292
9.5 Алгоритм Укконена–Майерса	299
9.6 Заключение	307
Глава 10. Приближенное сравнение с паттерном	309
10.1 Алгоритм для произвольной функции расстояния	310
10.2 Алгоритм, вычисляющий несовпадения	314
10.3 Алгоритмы, вычисляющие разности	320
10.3.1 Алгоритм Укконена	322
10.3.2 Алгоритм Майерса	325
10.4 Быстрый и гибкий алгоритм Ву и Менбера	334
10.5 Сложность алгоритмов приближенного сравнения с паттерном	341
Глава 11. Регулярные выражения и множественные паттерны	343
11.1 Алгоритмы для регулярных выражений	346
11.1.1 Недетерминированные конечные автоматы	346
11.1.2 Детерминированные конечные автоматы	352
11.1.3 Модификация алгоритма ВМ	355
11.2 Алгоритмы сравнения с множественными паттернами	359
11.2.1 Автомат Ахо–Корасика: модификация алгоритма КМП	360
11.2.2 Автомат Комменца–Вальтера: модификация алгоритма БМ	364
11.2.3 Аппроксимирующие паттерны: модификация алгоритма ВМ	366
11.2.4 Аппроксимирующие паттерны: алгоритм Бейза-Ятса–Наварро	368
Часть IV. Вычисление характеристических паттернов	381
Глава 12. Периодичность	383
12.1 Все кратные строки	386
12.1.1 Алгоритм Крочемора	386
12.1.2 Алгоритм Мейна и Лоренца	397
12.2 Серии	406
12.2.1 Крайние левые серии — алгоритм Мейна	407
12.2.2 Все серии — алгоритм Колпакова и Кучерова	413
Глава 13. Обобщение периодичности	417
13.1 Все оболочки — алгоритм Ли–Смита	418

Содержание	9
13.2 Все раппорты — алгоритм Франека–Смита–Танга	430
13.2.1 Вычисление NE-дерева	432
13.2.2 Вычисление NE-массива	435
13.3 k -приближенные раппорты — алгоритм Шмидта	441
13.4 k -приближенные периоды — алгоритмы Сима–Илиопулоса–Парка–Смита	459
Литература	467
Предметный указатель	483

Предисловие

В начале было Слово,
и Слово было у Бога,
и Слово было Бог.

— Иоанн 1.1

Вычисление паттернов в строковых последовательностях — это фундаментальная проблема, которая возникает во многих областях науки и информационных технологий. Манипулирование текстом в текстовых редакторах, лексический анализ компьютерных программ, работа конечных автоматов, извлечение информации из баз данных — это малая часть тех процессов, которые требуют нахождения или вычисления паттернов. Алгоритмы вычисления паттернов находят применение в таких областях, как сжатие данных, криптография, распознавание речи и компьютерное зрение, вычислительная геометрия и молекулярная биология. Тема вычисления паттернов в строковых последовательностях важна не только из-за своего практического применения. Она является частью комбинаторики, где, как известно, существует много просто формулируемых задач, для которых, однако, очень сложно найти решение, и интерпретация таких задач, как вычисление паттернов, часто позволяет найти элегантное и точное их решение.

В этой связи вызывает большое удивление, что академические факультеты математики или компьютерных наук в своем большинстве не включают в магистерские курсы или в курсы для аспирантов эту интересную, важную и сложную для исследований тему. Еще более удивительно, что существует всего несколько книг и обзоров, где собраны вместе и основополагающие теоретические результаты, и практические алгоритмы, которые появились в последнюю четверть прошедше-

го столетия. Я знаю всего пять книг [214, 67, 108, 206, 61] и три объемных обзора научных статей [21, 2, 191], содержимое которых значительно перекрывает материал данной книги. Обзорные статьи и книги [214, 67] написаны в большей мере как обобщение итогов исследования авторов, нежели как книги для студентов. Книги [108, 206] касаются, в основном, применения строковых алгоритмов в молекулярной биологии. Последняя монография [61], написанная легко и элегантно, сочетает в себе как монографию, так и руководство по строковым алгоритмам. К сожалению, в настоящее время она доступна только на французском языке.

Данная книга ставит целью заполнить этот пробел: дать общее введение в алгоритмы вычисления паттернов, которое было бы полезно в качестве отправного пункта для научных исследований и давало бы серьезный фактический материал, доступный для изучения студентам старших курсов и аспирантам. Задержимся на некоторое время на трех словах из предыдущего предложения: “доступный”, “алгоритм” и “паттерн”.

Основное “свойство” этой книги — сделать материал *доступным* для студентов старших курсов и аспирантов, имеющих специализацию по математике или компьютерным наукам и которые знакомы с дискретными структурами и алгоритмами, оперирующими этими структурами.

Первым условием доступности материала книги является достаточная математическая подготовка читателя, а не его знакомство со строковыми последовательностями. Предполагается, что студент прослушал стандартный курс дискретной математики, курсы по структурам данных и анализу алгоритмов. В этом случае он знает, что такое стеки, очереди, связанные списки и массивы, имеет понятие об анализе алгоритмов и о том, как записать “асимптотическую сложность” алгоритма, имеет некоторый опыт работы с математическими утверждениями и методами, используемыми для доказательства корректности алгоритмов, также знаком с алгоритмами, выполняемыми на графах и деревьях. В дополнение к перечисленному, предполагается, что читатель знаком с какими-либо языками программирования и способен читать и понимать алгоритмы, записанные на этих языках.

Второе условие не связано с компетентностью и подготовленностью читателя: моя цель — “соблазнить” студента и читателя интереснейшей и увлекательной областью новых знаний. Я не собирался писать энциклопедию алгоритмов, вычисляющих паттерны в строковых последовательностях. Вместе с тем я хотел представить результаты, которые (я надеюсь) важны и которые можно обобщить и расширить. Но это неизбежно ведет к тому, что некоторые интересные результаты были опущены (нельзя охватить необъятное!). И все-таки я надеюсь, что выбранный подход к изложению материала будет стимулировать читателя и далее к самостоятельному изучению им научной литературы.

Особым “субъектом” материала этой книги является математический объект, который в компьютерных науках называется “строка” (string) или “строковая последовательность” (в Европе, в среде математиков, более распространен термин

“слово”). Но основное внимание в книге уделяется *алгоритмам*, т.е. точным методам и процедурам, предназначенным для выполнения “чего-то”. Исходя из этого данную книгу скорее можно отнести к книгам по компьютерным наукам, чем к математическим книгам. Поэтому она значительно отличается от классической монографии [164], посвященной этой теме, и ее “потомков” [162, 163], ставящих во главу угла математические аспекты данной темы. Нас в первую очередь будут интересовать алгоритмы, находящиеся в строковых последовательностях различного рода паттерны, и только во вторую очередь — математические свойства самих строк. Это, конечно, не означает, что математические результаты не будут представлены строго и последовательно. Это означает, что будут представлены только те математические результаты, которые необходимы для пояснения построения и поведения алгоритмов. И последнее замечание: я сознательно ограничился изложением последовательных алгоритмов для обработки одномерных строк, не делая ссылок на обширную литературу по алгоритмам с распараллеливанием процесса вычисления или на быстро растущую литературу по многомерным (особенно, двухмерным) строкам.

Другой основной “герой” нашей книги — это *паттерн*. Все рассмотренные в книге алгоритмы предназначены для нахождения в строковых последовательностях определенных типов паттернов. Я говорю “определенных типов”, поскольку будем различать три основных типа паттернов — частные, характеристические и внутренние. Каждому типу паттернов посвящена соответствующая часть книги.

Частный паттерн (specific pattern) — это единственный вид паттернов, который можно задать в виде списка символов в нужном порядке. Например, в строке $x = abaababaabaab$ мы можем найти (трижды) паттерн $u = abaab$, но не найдем паттерн $u = ababab$. (Иногда паттерн может содержать специальные “символы замещения”, и в этом случае возможно только “приближенное” (в некотором точно определенном смысле) сравнение паттерна и строки.)

Характеристические паттерны (generic patterns) основаны на специальных представлениях структурной информации о строковых последовательностях. Например, мы можем говорить о “повторениях” в строке x — в этом случае в строке x есть несколько смежных одинаковых подстрок. (Например, в приведенной выше строке x присутствуют повторяющиеся подстроки $(aba)(aba)$, $(abaab)(abaab)$, aa (три отдельные серии) и несколько других, если вы сможете их найти.)

Последний тип паттернов, которые будут рассмотрены в книге, я назвал *внутренними* (intrinsic). Эти паттерны отображают внутреннюю структуру строковых последовательностей. Мы рассмотрим различные паттерны, которые показывают наличие периодических структур в строках, например нормальную форму, дерево суффиксов, линдонскую декомпозицию, s -факторизацию. Эти паттерны вездесущие: они используются почти во всех алгоритмах вычисления частных и характеристических паттернов. Другими словами, они формируют основу для эффективных процедур обработки строковых последовательностей. Раз-

нообразии внутренних паттернов поразительно: для строки x нашего примера нормальная форма имеет вид $(abaababa)(abaab)$, тогда как линдонскую декомпозицию можно записать как $(ab)(aabab)(aab)(aab)$, а s -факторизацию — как $(a)(b)(a)(aba)(baaba)(ab)$, и все эти паттерны полезны и эффективны с вычислительной точки зрения.

Эта книга имеет следующую организацию. В части I приведены основные сведения о строковых последовательностях и алгоритмах обработки строк. Здесь даны терминология, формы записи и основные свойства строковых последовательностей. Глава 2 является ключом к остальной части книги: здесь четко поставлены задачи, алгоритмы для решения которых описаны в последующих частях книги. На основе материала этой главы читатель может выбрать для себя направление дальнейшего чтения в виде тех глав, которые представляют для него наибольший интерес. В этой части также обсуждаются качества “хороших” алгоритмов и вопросы их реализации на практике. В частях II–IV описаны алгоритмы для вычисления внутренних, частных и характеристических паттернов соответственно.

Всего в книге 13 глав, разбитых на четыре части. Каждая глава, как видно из оглавления, разделена на несколько разделов, каждый из которых заканчивается набором упражнений. Там, где это необходимо, главы заканчиваются разделами, обобщающими изложенный материал и рассматривающими смежные вопросы, дополнительные темы и нерешенные проблемы.

Отметим, что в книге приведено примерно 500 упражнений, которые составляют неотъемлемую часть книги и могут использоваться для следующих целей.

- Для проверки степени усвоения читателем прочитанного материала.
- Сделать более ясными или показать в другом контексте понятия или результаты, приведенные в тексте.
- Показать расширения или модификации алгоритмов и математических результатов, которые приведены в тексте без подробного обсуждения.
- Показать детали алгоритмов и доказательств, которые не включены в основной текст из-за их громоздкости или их “ухода” от основной темы.

С помощью упражнений я также пытался вовлечь читателя в процесс разработки и анализа представленных алгоритмов. Этим я хотел показать, что в основе большинства алгоритмов и их модификаций лежат простые идеи, проникнуть в суть которых не составляет труда, но которые, может быть, “затемнены” или отброшены предыдущими исследователями. Если идея понята и принята, остается только ее техническая реализация. Это общее замечание относится и к строковым алгоритмам.

Признаюсь, что непосредственно в процессе написания книги я мог допускать ошибки. Поэтому при ее вычитке я старался исправить все замеченные ошибки и оплошности и сгладить шероховатости изложения материала. Но, конечно, я не

могу гарантировать, что внимательный читатель не найдет “дефектов” в моей книге. Я поддерживаю Web-узел

<http://www.cs.curtin.edu.au/~smyth/patterns.shtml>,

где вы можете оставить свои сообщения о замеченных ошибках и предложения по улучшению книги. Я также буду благодарен, если читатели по этой же причине свяжутся со мной по электронной почте

smyth@computing.edu.au или smyth@mcmaster.ca

Материал книги можно использовать, как минимум, для чтения двухсеместрового курса (каждый семестр по 12–14 недель) для студентов старших курсов и аспирантов. Материал первых глав в разное время уже читался аспирантам факультета компьютерных наук и систем и факультета вычислительной техники и программного обеспечения университета Мак-Мастера, г. Гамильтон, Онтарио, Канада, и аспирантам факультета компьютерных наук университета Дебрецена, Венгрия. Аспиранты, прослушавшие эти курсы, внесли свой вклад в создание этой книги.

Я хотел бы выразить глубочайшую благодарность школе вычислительной техники университета Кортин, Перт, Западная Австралия, и ее бывшим и настоящему руководителям Деннису Муру, Терри Силли, Свете Венкатеш и Джеффу Уесту (Dennis Moore, Terry Cealli, Svetha Venkatesh, Geoff West) за щедрую поддержку и за содействие, как интеллектуальное, так и практическое, на протяжении нескольких лет. Большая часть книги написана во время моих продолжительных визитов в Кортин. Я также благодарен профессору Петью Аттиле (Pethő Attila), декану факультета компьютерных наук университета Дебрецена, за его интерес к моей работе и поддержку, особенно на последнем этапе создания электронного варианта книги. Хотел бы также выразить свою глубокую благодарность моим друзьям и коллегам Лейле Багдади, Джерри Чепплу, Франью Франеку, Костасу Илиопулосу, Терри Лекроку, Ян Ли, Деннису Муру, Пату Риану, Джеми Симпсону и Ксиангдонгу Ксиао (Leila Baghdadi, Jerry Chapple, Franya Franěk, Costas Iliopoulos, Thierry Lecroq, Yin Li, Dennis Moore, Pat Ryan, Jaime Simpson, Xiangdong Xiao) за их дружеское и полезное содействие. Большая благодарность двум анонимным рецензентам, которые сделали конструктивные замечания и предложения по книге. Наконец, возношу хвалу своей дочери Жаклин за ее восхитительный выбор и попытки попробовать на вкус слова “строка” и “слово”.

W.F.S.

Моим родителям

ЧАСТЬ I

Строковые последовательности и алгоритмы

Слово спасет мир (если не будет поздно)

— Аноним

ГЛАВА 1

Свойства строковых последовательностей

Истинные слова не могут быть красивыми;
Красивые слова не могут быть истинными.

— Лао-цзы (604–531 гг. до н.э.). Путь Лао-цзы

1.1 Последовательность жемчужин

Предположим, что у нас есть кучка жемчужин. Для наглядности расположим жемчужины на столе слева направо в одну линию рядом друг с другом. Пусть всего есть n жемчужин, на каждую жемчужину приклеим маленькую бирку с меткой. Предположим, что метками являются целые числа от 1 до n и метки присваиваются жемчужинам по следующим правилам.

- Метка самой левой жемчужины равна 1.
- Для жемчужины с меткой i ($i = 1, 2, \dots, n - 1$) соседняя справа жемчужина имеет метку $i + 1$.

Эти правила удовлетворяют нашему интуитивному представлению о том, как сделать из кучки жемчужин “строковую последовательность”: жемчужины расположены в одномерную цепочку (ряд), по которой можно перемещаться из одного конца в другой через соседние жемчужины. Но опираясь на приведенные выше правила, можно сформулировать более общее понятие последовательности. Во-первых, можно отказаться от “жемчужин”, если ввести понятие *элемент* последовательности. Во-вторых, метки не обязаны принимать значения только целых

положительных чисел. Метки могут быть, например, цветными или буквами какого-нибудь алфавита. Что действительно важно, так это следующие правила, определяющие строковую последовательность.

0. Каждый элемент последовательности имеет уникальную метку.
1. Каждый элемент с некоторой меткой x (за исключением не более одного элемента, который называется *самый левый элемент*) имеет единственный *предшествующий элемент* с меткой $p(x)$.
2. Каждый элемент с некоторой меткой x (за исключением не более одного элемента, который называется *самый правый элемент*) имеет единственный *последующий элемент* с меткой $s(x)$.
3. Для любого элемента с меткой x , который не является самым левым, выполняется равенство $x = s(p(x))$.
4. Для любого элемента с меткой x , который не является самым правым, выполняется равенство $x = p(s(x))$.
5. Для двух различных элементов с метками x и y существует такое целое положительное число k , что $x = s^k(y)$ или $x = p^k(y)$.

Эти правила охватывают сущность понятия *сочленения* и операции *конкатенации*, выполняющей это сочленение: каждый элемент последовательности, за исключением самого левого и самого крайнего элементов, имеет единственный предшествующий и единственный последующий элементы. Самый левый и самый крайний элементы имеют или единственного последующего, или единственного предшествующего элемента. Более того, начиная с любого элемента с меткой x мы, перебирая конечную последовательность предшествующих и последующих элементов, можем достичь любого другого элемента с меткой y . Отсюда вытекает следующее утверждение.

Определение 1.1.1. Строковая последовательность — это набор элементов, которые удовлетворяют правилам 0–5. ■

Критическим (но не очевидным) фактором этого определения являются условия в правилах 1 и 2, что может быть *не более одного* самого левого и *не более одного* самого правого элементов. Для примера рассмотрим, что случится, если из цепочки жемчужин сформировать петлю, т.е. расположить их по кругу. Теперь в этой последовательности жемчужин нет “самого левого” и “самого правого” элементов. Однако отметим, что это не порождает проблем, поскольку все правила 0–5 выполняются.

Теперь предположим, что количество жемчужин бесконечно — есть левый конец цепочки жемчужин, а правый конец “скрыт за горизонтом”. Для этой бесконечной последовательности также выполняется определение 1.1.1, при этом существует самый левый элемент, но не существует самого правого. Нетрудно

проверить, что правила 0–5 будут выполняться, если цепочку жемчужин продолжить до бесконечности и в левом направлении; в этом случае нет ни самого левого, ни самого правого элементов.

В этой книге в свое время будут рассмотрены все возможные виды строковых последовательностей. Примем следующее соглашение по терминологии. Строковую последовательность (string) для краткости будем называть *строкой*.¹ Строку с конечным числом элементов, содержащую как самый левый, так и самый правый элементы, назовем *линейной строковой последовательностью* или *линейной строкой*. Строковую последовательность с конечным (не нулевым) числом элементов, не имеющую ни самого левого, ни самого правого элементов, назовем *строковой петлей* (в литературе такой вид последовательности также имеет название *циклическая* (или *круговая*) *строковая последовательность*). Строковую последовательность с бесконечным количеством элементов, содержащую самый левый элемент, назовем *бесконечной строковой последовательностью* или *бесконечной строкой*; бесконечную строковую последовательность, не имеющую ни самого левого, ни самого правого элементов, назовем *бесконечной петлей*. Далее термин “строковая последовательность” (или просто “строка”) будет означать любой объект, который удовлетворяет определению 1.1.1. Тип строки обычно легко определить из контекста. Кроме того, на практике нетрудно отличить линейную строку от петли, поскольку, как правило, петля определяется через соответствующую линейную строку x и записывается как $C(x)$, т.е. петля $C(x)$ формируется из линейной строки x путем назначения самого левого элемента этой последовательности последующим элементом для самого правого элемента.

Упражнения 1.1

1. Объясните, почему для сочленения строковой последовательности необходимо правило 3? Приведите пример математического объекта, для которого выполняются правила 0–2, но не выполняется правило 3.
2. Можно ли правило 4 вывести из правил 0–3? Объясните свой ответ.
3. Объясните, почему для сочленения строковой последовательности необходимо правило 5? Приведите пример математического объекта, для которого выполняются правила 0–4, но не выполняется правило 5.
4. Является ли бесконечное множество $\{a, a^2, \dots\}$ бесконечной строковой последовательностью?
5. Может ли в соответствии с определением 1.1.1 строковая последовательность состоять из одного элемента a ? Может ли такая последовательность быть петлей?

¹В русской математической литературе строковые последовательности часто называют “словами”. — *Примеч. ред.*

6. Может ли в соответствии с определением 1.1.1 строковая последовательность не иметь элементов (такая последовательность или строка называется *пустой*)? Может ли пустая строка быть линейной строковой последовательностью? А петель?
7. С учетом предыдущих упражнений определите, сколько различных типов строк нужно дополнительно включить в классификацию строковых последовательностей, данную в этом разделе?
8. Классификация строковых последовательностей, приведенная в этом разделе, не включает следующие типы строк.
 - а) строковая последовательность с бесконечным числом элементов, имеющая самый правый элемент, но не имеющая самого левого;
 - б) строковая последовательность с конечным числом элементов, имеющая или самый правый элемент, или самый левый, но не оба сразу.
 Объясните эти исключения.
9. Существует ли какой-нибудь способ *доказать*, что определение 1.1.1 определяет строковую последовательность?

1.2 Линейные строковые последовательности

Определение строковых последовательностей, данное в предыдущем разделе, очень общее и охватывает многие типы последовательностей, встречающиеся на практике. Так, строками будут

- слова в английском языке, здесь элементами последовательностей будут прописные и строчные буквы английского алфавита вместе с апострофом и дефисом;
- текстовые файлы, здесь элементами будут символы в кодировке ASCII;
- книга, написанная на китайском языке, где элементами последовательности являются китайские иероглифы;
- компьютерная программа, где элементами будут определенные (пробел, двоеточие, точка с запятой и т.п.) вместе со “словами”, со “словами”, заключенными между этими разделителями;
- последовательность молекул ДНК (иногда длиной до трех миллиардов элементов-молекул), состоящая только из букв *C*, *G*, *A* и *T*, которые соответствуют нуклеотидам цитозину, гуанину, аденину и тимину;
- поток частиц, бомбардирующих космический аппарат;
- список длин сторон выпуклого многоугольника.

Это примеры того, что в предыдущем разделе мы назвали “линейной строковой последовательностью”. В книге, в основном, рассматриваются именно такие строки, поэтому в данном разделе даны некоторые определения и соответствующая терминология, необходимые для работы с линейными строковыми последовательностями. Большинство (но не все) из представленного ниже применимо также к петлям, бесконечным и пустым строкам.

В приведенных выше примерах видно, что важным свойством любой строки является природа ее элементов, будь-то действительные числа, слова какого-либо языка, иероглифы и т.п. На практике удобно представлять элементы строк как члены некоторого множества. Такое множество называется *алфавитом*, поэтому члены этого множества естественно назвать *буквами* (далее мы увидим, что термин “буква” можно интерпретировать значительно шире, чем просто буква какого-нибудь естественного языка).² Будем говорить, что строковая последовательность *определена на этом алфавите*. Конечно, если строка x определена на алфавите A , то она определена и на любом множестве, подмножеством которого является алфавит A . Другими словами, алфавит для данной строки x определяется неоднозначно. Можно определить минимальный алфавит для строки x как множество различных элементов, каждый из которых обязательно входит в x . Иногда принимается соглашение, что алфавитом может быть только минимальный алфавит (например, “двоичный” алфавит), иногда это условие снимается (пример — алфавит “действительных чисел”).

В книге алфавит будет обозначаться как A , а $\alpha = |A|$ будет обозначать его мощность. В случае α , равного 2, 3 и 4, алфавит будет называться *бинарным*, *тернарным* и *кватернарным* соответственно. Далее мы увидим, сколько интересных строковых последовательностей можно построить на бинарном алфавите. Последовательности на основе кватернарного алфавита находят применение при анализе последовательностей ДНК.

В общем случае, кроме требования “различности” элементов множества A , на само это множество особых условий не накладывается. Алфавит может быть конечным (даже нулевым!) множеством, счетным множеством (например, множество целых чисел) и даже несчетным (как континуальное множество вещественных чисел). Кроме того, элементы алфавита могут быть упорядоченными (тогда для любой пары различных элементов определен результат их сравнения “больше” или “меньше”), неупорядоченными и частично упорядоченными. Для многих алгоритмов, описанных в книге, используется неупорядоченный алфавит. В главе 4 показано, что для многих алгоритмов свойство упорядоченности алфавита значительно влияет на их вычислительную эффективность и на выбор тестовых данных для проверки алгоритмов.

²Из этих определений вытекает обоснованность терминов “строковая последовательность” и “строка”. — *Примеч. ред.*

Для данного алфавита A обозначим через A^+ множество всех возможных непустых конечных сочетаний букв алфавита A . Например, если $A = \{a\}$, тогда $A^+ = \{a, a^2, a^3, \dots\}$, где запись a^2 обозначает строку aa , запись a^3 — строку aaa и т.д. Если $A = \{0, 1\}$, тогда A^+ состоит из всех различных непустых конечных двоичных (битовых) последовательностей или, что то же самое, A^+ представляет множество всех неотрицательных целых чисел. Как предлагалось в упражнении 1.1.6, имеет смысл ввести *пустую строковую последовательность* (или *пустую строку*), которая обычно обозначается символом ε . Если определена пустая строка ε , тогда можно определить еще два множества: $A' = A \cup \{\varepsilon\}$ и $A^* = A^+ \cup \{\varepsilon\}$. В этих обозначениях можно дать другое определение линейной строковой последовательности, эквивалентное данному в предыдущем разделе.

Определение 1.2.1. **Линейными строковыми последовательностями (линейными строками)** на алфавите A называются элементы множества A^+ . Элементы множества A^* называются **конечными линейными строковыми последовательностями (конечными линейными строками)** на алфавите A . ■

Таким образом, множество A^+ является множеством всех линейных строк на алфавите A . Отметим, что пустая строка *не является* линейной, поскольку не имеет ни самого левого, ни самого правого элементов.

На протяжении всей книги строковые последовательности будут обозначаться строчными жирными буквами: p , t , x и т.п. Мы будем трактовать строковую последовательность, удовлетворяющую правилам 0–5 из раздела 1.1, как одномерный массив. Конечно, можно использовать и другие представления строковых последовательностей, например в виде связанных списков, но одномерный массив является наиболее естественной моделью для таких строк. Исходя из этой модели произвольную строку x на алфавите A в виде массива можно задать следующим образом:³

$$x : \text{array}[1..n] \text{ of } A.$$

В этом случае будем говорить, что строка x имеет *длину* $n = |x|$ и *позиции* $1, 2, \dots, n$. Для любого целого $i \in 1..n$ буква в позиции i обозначается как $x[i]$, поэтому можно записать

$$x = x[1]x[2] \cdots x[n],$$

что, в соответствии с определением 1.1.1, является сочленением n строк длины 1. Фактически, здесь позиция i играет роль метки из правил 0–5.

Отметим, что модель массива применима также к пустым последовательностям $x = \varepsilon$, которым соответствует пустой массив и которые имеют нулевую длину.

³ Автор по всей книге, не оговаривая это специально, широко использует нотацию, применяемую для записи алгоритмов. Считаем, что у читателя такая нотация не вызовет недоразумений или трудностей. — *Примеч. ред.*

Обсуждение 1.2.1. Выше было сказано, что массивы являются “простыми и естественными” представлениями строковых последовательностей. Это утверждение в значительной степени отражает “вычислительную” точку зрения на данную проблему. Действительно, массив — это простая структура данных, но будет ли массив “естественен”, зависит от наших представлений о механизмах получения доступа к элементам последовательностей. Поскольку строковые последовательности определяются через понятие сочленения и операцию конкатенации, то, очевидно, “естественно” осуществить поиск нужного элемента последовательности (получить доступ к нему) через предшествующие или последующие элементы. Такой механизм получения доступа лучше реализуется на представлении строк с помощью связанных списков,⁴ где для получения элемента в позиции i из “текущей” позиции j требуется время, пропорциональное $|j - i|$. С другой стороны, для получения доступа к элементу массива путем простого указания его позиции требуется фиксированное (постоянное) время, не зависящее от положения искомого элемента и “текущего” элемента. С такой точки зрения представление строковых последовательностей в виде массивов кажется более предпочтительным, чем представление в виде связанных списков, поскольку получаем оценки времени выполнения алгоритмов, построенных на использовании массивов, меньшие, чем аналогичные оценки для алгоритмов, использующих списки.

На практике алгоритмы, выполняемые над строками, почти всегда начинаются с самого левого элемента (с позиции 1) или самого правого элемента (с позиции n) и далее проверяются соседние (предшествующие или последующие) элементы (или позиции) один за другим. С другой стороны, выходом таких алгоритмов часто являются позиции определенных элементов в строке, и здесь неявно предполагается, что получить доступ к этим позициям можно за фиксированное время. Один из способов “примерить” эти противоречивые представления о доступе к элементам строки заключается в том, что первоначально строка предоставляет доступ к своим элементам как связанный список — элементы доступны только один за другим либо слева направо или справа налево — и в таком порядке элементы копируются в массив. Для такого копирования требуется времени порядка $O(n)$ и дополнительной памяти такого же порядка $O(n)$. Такое преобразование не повлияет на асимптотические показатели любых алгоритмов, рассмотренных в книге. Поэтому примем в некоторой степени “лишнее” соглашение, что на входе алгоритмов строковые последовательности обрабатываются как связанные списки, а на выхо-

⁴Связанным списком называется структура для хранения данных, например элементов a_1, a_2, \dots, a_n , где каждая “ячейка”, в которой хранится элемент a_i , имеет также указатель на “ячейку”, в которой хранится элемент a_{i+1} . Как правило, связанный список имеет заголовок, где хранится указатель на “ячейку” с элементом a_1 . Таким образом, чтобы получить доступ к некоторому элементу a_i , надо последовательно пройти по указателям от “ячейки” с элементом a_1 до “ячейки”, содержащей элемент a_i . — *Примеч. ред.*

де алгоритмов, если необходимо, они представимы в виде массивов. Мы обещаем предупредить читателя, когда будем отступать от этого соглашения. ■

Равенство строковых последовательностей определяется обычным способом. Строка x длины n и строка y длины m *равны* (записывается как $x = y$) тогда и только тогда, когда $n = m$ и выполняется равенство $x[i] = y[i]$ для всех $i = 1, \dots, n$. Из этого определения следует, что пустая строка ε равна только самой себе. Отметим на будущее, что присоединение пустой строки слева или справа к любой строке x не изменяет эту строку, т.е. $x = \varepsilon x \varepsilon$.

Для любой пары целых чисел i и j , таких, что $1 \leq i \leq j \leq n$, определим *подстроку* $x[i..j]$ строки x следующим образом:

$$x[i..j] = x[i]x[i+1] \dots x[j].$$

Будем говорить, что подстрока $x[i..j]$ *начинается* с позиции i (или, что тоже самое, что ее *начальная позиция* i) и что ее *длина* равна $j - i + 1$. Если $j - i + 1 < n$, то подстрока $x[i..j]$ называется *собственной подстрокой* строки x . Отметим два крайних случая подстрок: подстрока $x = x[1..n]$ длины n и подстрока $x[i] = x[i..i]$ длины 1. Для пары целых чисел i и j , таких, что $i > j$, примем соглашение, что $x[i..j] = \varepsilon$, т.е. является пустой подстрокой нулевой длины. Поскольку $x = \varepsilon x \varepsilon$, то можно утверждать, что ε является подстрокой любой строки из множества A^* .

Пусть k — целое число, $k \in 1..n$, и рассмотрим позиции $i_k = 1, 2, \dots, n - k + 1$ в строке x . Каждая из этих $n - k + 1$ позиций является начальной позицией подстроки $x[i_k..i_k + k - 1]$ длины k . Таким образом, любая строка длины n имеет $n - k + 1$ подстрок (не обязательно различных) длины k .

Если u является подстрокой (соответственно, собственной подстрокой) строки x , тогда x называется *надстрокой* (соответственно, *собственной надстрокой*) строки u . Конечно, подстроки и надстроки являются строковыми последовательностями.

Если множество A упорядочено, то с помощью подстрок можно индуцировать упорядоченность элементов множества A^* . Такой порядок называется *лексикографическим*. Дадим более точное определение: пусть имеется две строки $x = x[1..n]$ и $y = y[1..m]$, где $n \geq 0$ и $m \geq 0$. Будем говорить, что x лексикографически меньше y (записывать это будем как $x < y$), если выполняется одно из следующих (взаимно исключающих) условий:

- $n < m$ и $x[1..n] = y[1..n]$ (говоря кратко, здесь строка x является “собственным префиксом” строки y);
- существует целое число i ($i \in 1.. \min\{n, m\}$), такое, что $x[1..i-1] = y[1..i-1]$ и $x[i] < y[i]$ (здесь i — первая позиция, в которой элементы строк x и y различны).

Например, на основе порядка английского алфавита можно сделать такие сравнения строк.

- $ab < abc$ (поскольку $i = 2 = n < 3 = m$).
- $\varepsilon < a$ (поскольку $i = 0 = n < 1 = m$).
- $ab < aca$ (поскольку $i = 2$ и $b < c$).

Отметим, что данное выше определение работает и в случае, когда одна или обе сравниваемые строки бесконечны. Основываясь на этом определении, нетрудно ввести другие отношения порядка: $x \leq y$ тогда и только тогда, когда $x = y$ или $x < y$; неравенство $x > y$ выполняется только тогда, когда справедливо $y < x$; наконец, $x \geq y$, если $y \leq x$.

Запись $x = u_1 u_2 \dots u_k$, где u_i — непустые подстроки, называется *декомпозицией* или *факторизацией* строки x по *факторам* u_i (применение этого определения см. в разделе 1.4 и главе 6).

Есть два типа подстрок $x[i..j]$, которые очень важны и поэтому имеют специальные названия. Для произвольного целого $j \in 0..n$ подстрока $x[1..j]$ называется *префиксом* строки x , иногда такая подстрока обозначается как $\text{pref}(x)$; если $j < n$, то подстрока $x[1..j]$ называется *собственным префиксом* строки x . Аналогично для произвольного целого $i \in 1..n + 1$ подстрока $x[i..n]$ называется *суффиксом* строки x (записывается как $\text{suff}(x)$) и *собственным суффиксом*, если $i > 1$. Отметим, что из данного определения и равенства $x = \varepsilon x \varepsilon$ следует, что пустая строка ε является собственным префиксом и собственным суффиксом любой строковой последовательности x . Для примера, строка

$$f = abaababaabaab$$

имеет префиксы $\varepsilon, a, ab, aba, \dots, f = abaababaabaab$ и суффиксы $\varepsilon, b, ab, aab, \dots, f$. Собственные префиксы и собственные суффиксы получаем из приведенных путем исключения из них самой строки f .

Важным понятием, которое будет полезным во многих задачах, является понятие *грань* (border).

Определение 1.2.2. *Гранью b строки x называется любой собственный префикс этой строки, равный суффиксу x .* ■

Из этого определения видно, что любая строка x всегда имеет пустую грань $b = \varepsilon$ длины $\beta = 0$, но сама строка x не является собственной гранью. Будем использовать символ β для обозначения длины $|b|$ грани b . Часто представляет интерес *наибольшая грань b^** длины $\beta^* = |b^*|$, где $0 \leq \beta^* \leq n - 1$. Строковая последовательность f , приведенная выше, имеет две непустые грани: ab и $abaab$. Строка $g = abaabaab$ имеет в точности такие же две грани, но в этом случае наибольшая грань $abaab$ частично перекрывает сама себя (точнее, перекрываются

грань и суффикс, с которым совпадает эта грань). Аналогично строка a^n имеет грани a, a^2, \dots, a^{n-1} , среди которых грани $a^i, a^{i+1}, \dots, a^{n-1}$, где $i = \lceil 9(n + 1)/2 \rceil$, перекрываются.⁵ Далее будет показано, что перекрывающиеся грани являются характеристикой строковых последовательностей, содержащих повторяющиеся подстроки. Например, строку g можно записать в следующем виде $(aba)^2ab = (aba)(aba)ab$, где видны повторяющиеся подстроки.

Применим понятие грани для получения так называемой “нормальной формы” для данной непустой строки x длины n . Предположим, что для этой строки x найдена грань b и вычислена ее длина β . (В разделе 1.3 показано, как найти любую грань строки x за время порядка $\Theta(n)$.⁶) Из определения 1.2.2 имеем

$$x[1..\beta] = x[n - \beta + 1..n]. \quad (1.1)$$

Отсюда видно, что величина $p = n - \beta \geq 1$ определяет сдвиг позиций равных элементов этих подстрок (отметим, чем больше значение β , тем меньше значение p), т.е. для любых целых $i \in 1..\beta$ должно выполняться равенство

$$x[i] = x[i + p]. \quad (1.2)$$

Если $\beta = 0$ ($p = n$), то равенства (1.1) и (1.2) тривиальны. Если же $2\beta \geq n$ ($2p \leq n$), тогда строка x содержит, по крайней мере, две одинаковые смежные подстроки $x[1..p]$ и $x[p+1..2p]$. Более точно, строка x содержит $\lfloor n/p \rfloor$ одинаковых подстрок, каждая из которых имеет длину p , за которыми следует (возможно, пустой) суффикс длины $n \bmod p$.⁷ Положив $r = n/p$ и введя обозначение $u = x[1..p]$, получаем следующее выражение для любой строки x :

$$x = u^{\lfloor r \rfloor} u', \quad (1.3)$$

где $u' = x[1..n - \lfloor r \rfloor p]$ — собственный префикс (возможно, пустой) подстроки u . Если представить число r как сумму целой и дробной частей $r = \lfloor r \rfloor + k/p$ ($k \in \{0, \dots, p-1\}$) и интерпретировать запись $u^{k/p} = u[1..p]^{k/p}$ как подстроку $u[1..k]$, тогда равенство (1.3) можно переписать в более компактной форме

$$x = u^r. \quad (1.4)$$

Число p называется *периодом*, а число r — *порядком (показателем степени)* строковой последовательности x . Префикс $u = x[1..p]$ назовем *образующей (образующей подстрокой)* строки x . Отметим, что поскольку любая строка x

⁵Здесь и далее запись $\lceil a \rceil$ обозначает наименьшее целое число, не меньшее числа a . — *Примеч. ред.*

⁶Объяснение обозначению $\Theta(n)$ дано в разделе 1.3. — *Примеч. ред.*

⁷Здесь и далее запись $\lfloor a \rfloor$ обозначает наибольшее целое число, не превышающее числа a . Для положительных чисел $\lfloor a \rfloor$ совпадает с целой частью этих чисел. — *Примеч. ред.*

имеет пустую грань $b = \varepsilon$, то она всегда имеет *тривиальный период* $p = n$, *тривиальный порядок* $r = 1$ и *тривиальную образующую* x .

Представляет интерес вычисление периода $p = p(\beta)$ и порядка $r = r(\beta)$ как функций длины β грани b . Очевидно, что p является монотонно убывающей, а r монотонно возрастающей функциями от β . Поэтому если $\beta = \beta^*$ (β^* — длина наибольшей грани), то p принимает минимальное значение p^* , а r — максимальное значение r^* , которые называются *минимальным периодом* и *максимальным порядком* соответственно. В общем случае именно значения p^* и r^* представляют наибольший интерес, поэтому, если это не приводит к неопределенности, значение p^* будем называть просто *периодом*, а значение r^* — *порядком*. Аналогично в таком случае *образующей* будем называть подстроку $u = x[1..p^*]$.

Определение 1.2.3. Пусть строковая последовательность $x = x[1..n]$ имеет минимальный период p^* и, соответственно, $r^* = n/p^*$ и $u = x[1..p^*]$. Тогда декомпозиция

$$x = u^{r^*} \tag{1.5}$$

называется **нормальной формой** строковой последовательности x . ■

На основе нормальной формы (1.5) можно ввести следующую полезную классификацию строковых последовательностей.

- Строковая последовательность x называется *примитивной*, если $r^* = 1$. Если $r^* \neq 1$, то строковая последовательность называется *периодической*.
- В случае $r^* \geq 2$ строковая последовательность x называется *сильно периодической*, а при $1 < r^* < 2$ — *слабо периодической*.
- Если r^* целое число и $r^* \geq 2$, то строковая последовательность x называется *кратной* (repetition).⁸ В частных случаях $r^* = 2$ и $r^* = 3$ строку x будем называть *квадратом* и *кубом* соответственно.

Итак, строковая последовательность x может быть только или примитивной или периодической. Если строка x периодическая, то она может быть сильно или слабо периодической. Если x кратная, то она обязательно сильно периодическая. Отметим, что неравенство $r^* \geq 2$ возможно только тогда, когда x имеет грань длиной $\beta \geq n/2$.

Приведем несколько примеров применения введенных определений.

- Строка $x = aaabaabab$ примитивна ($p^* = n$).
- Строка $f = abaababaabaab = (abaababa)(abaab)$ слабо периодическая с периодом $p^* = 8$, порядком $r^* = 13/8$ и образующей $abaababa$.

⁸В русской литературе такие строки называются *строго периодическими* или даже просто *периодическими*. — Примеч. ред.

- Строка $g = abaabaab = (aba)^2ab$ строго периодическая с периодом $p^* = 3$, порядком $r^* = 8/3$ и образующей aba .
- Строка $x = (ab)^4$ кратная с периодом $p^* = 2$, порядком $r^* = 4$ и образующей ab .
- Строка $x = (abcabd)^2$ является квадратом с периодом $p^* = 6$ и образующей $abcabd$.

Отметим, что нормальная форма (1.5) фактически является “внутренним” паттерном в том смысле, каком это понятие применяется в части II этой книги: каждая строка x имеет паттерн, называемый “нормальной формой”, с помощью которого можно охарактеризовать периодические свойства этой строки. В будущем те простые понятия и определения (грань, период, порядок и т.д.), которые введены в данном разделе, будут постоянно возникать при исследовании различных алгоритмов, рассмотренных в книге.

Упражнения 1.2

1. Попробуйте “примирить” определения линейных строк, данные в разделах 1.1 и 1.2. Другими словами, докажите, что строка, являющаяся линейной в соответствии с одним определением, будет удовлетворять условиям линейности строк другого определения.
2. Предположим, что алфавит A состоит из α букв. Для произвольного неотрицательного целого n определите, сколько элементов множества A^* имеют длину n' . Сколько элементов этого множества имеют длину, не превышающую n ?
3. Пусть в условиях предыдущего упражнения $A = \{0, 1\}$. Каждому элементу множества A^+ можно поставить в соответствие неотрицательное целое число. Определите, какому количеству элементов этого множества соответствует каждое целое число?
4. Основываясь на определении равенства строк, докажите, что $\varepsilon = \varepsilon^2$.
5. На основе определений равенства строк и лексикографического упорядочения докажите, что для произвольных строк x и y на упорядоченном алфавите равенство $x = y$ выполняется тогда и только тогда, когда $x \not\prec y$ и $y \not\prec x$. В частности, покажите, что этот результат справедлив в случае $x = y = \varepsilon$. Также покажите, что ε является лексикографически наименьшим элементом множества A^* .
6. Основываясь на обычном порядке строчных букв английского алфавита, расположите следующие строки в возрастающем лексикографическом порядке.

$abbac, abbba, abb, abc, a, \varepsilon^2, ab, aba, \varepsilon b.$

7. Докажите, что оператор $<$ сравнения строк, определенный в этом разделе, обладает свойством транзитивности, т.е. докажите, что

$$x < y \text{ и } y < z \Rightarrow x < z.$$

8. Дайте независимое определение отношения порядка $>$, затем на основе этого определения совместно с определением отношения $<$, данного в этом разделе, покажите, что отношение $x > y$ имеет место только тогда, когда выполняется отношение $y < x$.
9. Пусть дана некоторая строка x длины n . Определите длины следующих подстрок строки x и наложите условия на i и k , при которых будут справедливы ваши ответы.
- $x[i..i + k - 1]$;
 - $x[i - k + 1..i]$;
 - $x[i + 1..k - 1]$;
 - $\varepsilon x[1..k]$.
10. Какое максимальное число различных подстрок можно образовать из строки длины n ? Приведите пример строки, на которой достигается этот максимум. Как можно охарактеризовать множество таких строк?
11. В предыдущем упражнении не накладывалось неявного ограничения на размер алфавита. Эта задача становится более интересной (и значительно более трудной), если наложить условие, что размер алфавита конечен и фиксирован. Такая задача на сегодняшний день пока не решена, что вы можете предложить для ее решения?
- Совет.** По-видимому, решение этой задачи следует начать с изучения следующего вопроса: для данных положительных чисел $\alpha = |A|$ и k надо найти наибольшую строку на алфавите A , которая содержит не более одной заданной подстроки длины k .
12. Опишите алгоритм, который находил бы все различные подстроки строки x . Докажите корректность вашего алгоритма и оцените его асимптотическую сложность (постарайтесь добиться порядка $O(n^2)$).
13. Пусть y — непустая строка длины m , а k — положительное целое число. Определите $|x|$ как функцию от m и k в каждом из следующих случаев.
- $x = y^k$;
 - $x = y^{|y|}$
 - $x = y^{|y|^k}$.
14. Какова длина строки

$$x = (ab)^n a (ab)^{n-1} a \cdots (ab)^2 a (ab) a?$$

15. Как определить префикс, суффикс и грань для пустой строки ε ?
16. Определите наибольшую грань, период и порядок для каждой из следующих строк и затем классифицируйте их как примитивную, сильно периодическую, слабо периодическую или кратную.
- $abaababaabaababaababa$;
 - $abcabacabcbacbcacb$;
 - $abcabdabcabdabcabd$;
 - $a(ab)^3a(ab)^3aa$.
17. Строка x называется **палиндромной последовательностью** (или просто **палиндромом**), если она читается справа налево так же, как слева направо. Если говорить более точно, то строка x называется **палиндромом**, если для каждого $i = 1, 2, \dots, \lfloor n/2 \rfloor$ выполняется равенство $x[i] = x[n - i + 1]$. Докажите или опровергните утверждение, что каждая грань палиндрома сама является палиндромом.
- Примечание.** Некоторые нетривиальные примеры палиндромов можно найти по электронным адресам
www.growndodo.com/wordplay/palindromes/dogseesada.html
complex.gmu.edu/neural/personnel/ernie/witty/palindromes.html
18. Покажите, что период и порядок “хорошо определены”; другими словами, если имеет место равенство $x = u^k u'$ для некоторой строки u и целого числа k , где u' является собственным префиксом строки u , тогда существует единственная соответствующая грань строки x .
19. Докажите, что если u^{r^*} является нормальной формой строки x , тогда подстрока u не может быть кратной. (Важность и полезность этого утверждения будут показаны в разделе 2.3.)
20. Пусть $(ab)^{3/2} = aba$. Какое из следующих равенств верно: $((ab)^{3/2})^2 = (aba)^2$ или $((ab)^{3/2})^2 = (ab)^{(3/2)^2} = (ab)^3$?
21. Докажите, что не существует строк с двумя различными примитивными образующими.
22. Предложите способ вычисления количества строк на алфавите $\{a, b\}$, которые имеют заданную длину n и являются примитивными.
- Совет.** Заметьте, что при любом нечетном n произвольную строку $x[1..n]$ можно сформировать из строк $x[1..(n-1)/2]$ и $x[(n+1)/2..n]$, а при четном n — из строк $x[1..n/2]$ и $x[n/2+1..n]$. Если вы встречаете затруднения при решении этого упражнения, обратитесь к работам [105, 106].
23. Приведенная в этом разделе классификация строковых последовательностей основана на значениях порядка (показателя степени) r^* . Повторите эту же классификацию на основе значений длины β^* наибольшей грани.

24. Нормальную форму строки x , введенную в этом разделе, более правильно назвать *левой* нормальной формой, поскольку вместо формулы $x = u^r u'$, где u' — префикс u , можно также ввести запись $x = v' v^s$, где v' — суффикс v .
- получите формулу для *правой* нормальной формы;
 - покажите, что классификация строковых последовательностей, построенная на основе левых нормальных форм, совпадает с такой же классификацией, построенной на основе правых нормальных форм.
25. Может ли петля быть подстрокой линейной строки? Может ли линейная строка быть подстрокой петли?
26. Нарисуйте дерево классификации строковых последовательностей, введенной в этом разделе.

1.3 Периодичность

Как будет показано в последующих главах, наибольшая грань строковой последовательности x порождает ряд числовых характеристик, которые полезны не только для классификации строк. Поэтому в данном разделе мы подробно рассмотрим вопрос о нахождении наибольших граней, в частности, покажем алгоритм нахождения наибольших граней, выполняющий вычисления за время порядка $\Theta(n)$.

Существует очевидный способ вычисления длины β^* наибольшей грани b^* строки x . Сначала положим $\beta^* \leftarrow 0$.⁹ Затем сравниваем подстроки $x[1]$ и $x[n]$; если они равны, то полагаем $\beta^* \leftarrow 1$. Далее сравниваем $x[1..2]$ и $x[n-1..n]$; если они равны, полагаем $\beta^* \leftarrow 2$. Продолжаем этот процесс, каждый раз сравнивая префиксы и суффиксы, длины которых на 1 больше, чем на предыдущем шаге, и увеличивая на 1 значение β^* , если они равны. На последнем шаге сравниваются подстроки $x[1..n-1]$ и $x[2..n]$, и если они равны, то β^* получает значение $n-1$. Конечное значение β^* равно длине наибольшей грани b^* строки x . Мы не даем этот алгоритм в формальной записи, поскольку он не эффективен — в упражнении 1.3.1 определяется его асимптотическая сложность.

Предположим, что в описанном алгоритме последовательные значения β^* сохраняются в массиве (или в строке!) $\beta[1..n]$. Тогда для любого $i = 1, 2, \dots, n$ $\beta[i]$ дает длину наибольшей грани подстроки $x[1..i]$. Назовем $\beta[1..n]$ **массивом граней** строки $x[1..n]$. Здесь $\beta[n]$ представляет длину наибольшей грани строки x , т.е. это то значение, которое необходимо для построения нормальной формы строки x .

Сделаем следующие заключения относительно массива граней и самих граней.

⁹Здесь и далее автор использует одну из форм нотаций для записи вычислительных алгоритмов. В этой нотации символ “ \leftarrow ” означает оператор присваивания, который в других нотациях обозначается как “ $=$ ” или “ $:=$ ”. — *Примеч. ред.*

- $\beta[1] = 0$ (поскольку ε является наибольшей гранью строки $x[1..1]$).
- Если $x[1..i]$ имеет грань длины $k > 0$ (здесь $2 \leq i \leq n$), тогда подстрока $x[1..i-1]$ имеет грань длины $k-1$. Отсюда следует, что для всех i , $1 \leq i \leq n-1$, справедливо неравенство $\beta[i+1] \leq \beta[i] + 1$.
- Для всех i , $1 \leq i \leq n-1$, равенство $\beta[i+1] = \beta[i] + 1$ выполняется только тогда, когда $x[i+1] = x[\beta[i] + 1]$ (поскольку $\beta[i] + 1$ — позиция в строке x элемента, расположенного справа от префикса $x[1..\beta[i]]$, который является наибольшей гранью подстроки $x[1..i]$).
- Если b — грань строки x , а b' — грань подстроки b , тогда b' является гранью строки x .

Эти заключения, в частности, третье и четвертое, показывают, что имеется возможность вычислить $\beta[i+1]$ на основании значений $\beta[1], \beta[2], \dots, \beta[i]$. Здесь основная трудность возникает тогда, когда $\beta[i] > 0$ и одновременно $x[i+1] \neq x[\beta[i] + 1]$. В этом случае надо обратиться ко второй по величине грани подстроки $x[1..i]$. Если эта грань пуста, тогда $\beta[i+1] = 1$, если $x[1] = x[i+1]$, или $\beta[i+1] = 0$, если $x[1] \neq x[i+1]$. Если же эта грань не пуста (обозначим ее длину как $\beta^2[i] > 0$), тогда необходимо сравнить элементы $x[i+1]$ и $x[\beta^2[i] + 1]$. В случае их равенства полагаем $\beta[i+1] \leftarrow \beta^2[i] + 1$. Если же эти элементы не равны, то в этом случае надо рассмотреть третью по величине грань подстроки $x[1..i]$ и повторить описанные выше действия. При необходимости рассматриваются следующие по величине грани подстроки $x[1..i]$.

Рассмотрим следующий пример. Пусть $x = abaababa*$, где символ $*$ показывает, что 9-я буква в строке x пока не определена. Массив граней строки x имеет вид 00112323?, где символ $?$ указывает на то, что остался не вычисленным только последний элемент массива. Отметим, что для $i = 8$ $\beta[i] = 3$. Если выполняются равенства

$$x[9] = x[i+1] = x[\beta[i] + 1] = x[4] = a,$$

тогда сразу получаем $\beta[9] = 4$. Если эти равенства не выполняются, тогда надо рассмотреть вторую по величине грань подстроки $x[1..8]$. Имеем

$$\beta^2[8] = \beta[\beta[8]] = \beta[3] = 1.$$

Если выполняются равенства

$$x[8+1] = x[\beta^2[8] + 1] = x[2] = b,$$

тогда получим $\beta[9] = 2$. Если же элемент $x[9]$ не совпадает ни с a ни с b , а принимает значение новой буквы (например, c), то в этом случае $\beta[9] = 0$.

Рассмотрим величины $\beta^j[i]$, $j = 1, 2, \dots, k$, представляющие значения длин j -х по величине граней подстроки $x[1..i]$. (Здесь мы полагаем, что $\beta^1[i] = \beta[i]$)

и $\beta^k[i] = 0$.) Поскольку для каждого j ($j = 2, 3, \dots, k$) $\beta^j[i] < \beta^{j-1}[i]$, то отсюда следует, что $x[1..\beta^j[i]]$ является гранью подстроки $x[1..\beta^{j-1}[i]]$. Таким образом, j -я по величине грань подстроки $x[1..i]$ будет гранью $(j - 1)$ -й по величине грани подстроки $x[1..i]$. Символически это запишется как

$$\beta^j[i] = \beta[\beta^{j-1}[i]]. \quad (1.6)$$

Если $x[i + 1] \neq x[\beta^{j-1}[i] + 1]$, тогда с помощью формулы (1.6) можно определить следующую позицию элемента, с которым будет сравниваться элемент $x[i + 1]$, — это будет элемент $x[\beta^j[i] + 1]$.

Приведенные рассуждения можно обобщить в виде следующей леммы.

Лемма 1.3.1. Пусть для любого целого $n \geq 1$ $x = x[1..n]$ — строковая последовательность с массивом граней $\beta = \beta[1..n]$. Пусть k — наименьшее целое, такое, что $\beta^k[n] = 0$. Тогда

- а) для каждого целого $j \in 1..k$ строка $x[1..n]$ имеет грань $x[1..\beta^j[i]]$;
- б) для произвольной буквы λ возможны только такие грани строки $x[1..n + 1] = x[1..n]\lambda$, длины которых являются членами убывающей последовательности

$$\langle \beta[n] + 1, \beta^2[n] + 1, \dots, \beta^k[n] + 1, 0 \rangle. \quad (1.7)$$

■

На основе равенства (1.6) можно построить эффективный алгоритм вычисления массива $\beta[1..n]$, который содержал бы значения длин наибольших граней всех префиксов $x[1..i]$ данной строки x . Этот алгоритм представлен ниже как алгоритм 1.3.1. Как будет показано в главе 7, все эти значения длин (не только $\beta[n]$) находят применение при сравнении различных строковых последовательностей, когда x является паттерном.

Алгоритм 1.3.1. Вычисление массива граней

```

▷ Вычисление массива граней  $\beta$  строки  $x[1..n]$ 
 $\beta[1] \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n - 1$  do
   $b \leftarrow \beta[i]$ 
  while  $b > 0$  and  $x[i + 1] \neq x[b + 1]$  do
     $b \leftarrow \beta[b]$ 
  if  $x[i + 1] = x[b + 1]$  then
     $\beta[i + 1] \leftarrow b + 1$ 
  else
     $\beta[i + 1] \leftarrow 0$ 

```

Теорема 1.3.2. Алгоритм 1.3.1 корректно вычисляет массив $\beta[1..n]$.

Доказательство. Сначала рассмотрим цикл **while**: в этом цикле на основе равенства (1.6) $b = \beta^{j-1}[i]$ заменяется на $b = \beta^j[i]$. Поскольку b монотонно убывает, этот цикл должен завершиться. Цикл **for** завершится, когда будет достигнуто равенство $i = n$. Следовательно, алгоритм 1.3.1 завершится за $n - 1$ шагов.

Теперь рассмотрим, что происходит после завершения цикла **while**. Выход из этого цикла осуществится тогда, когда выполнится хотя бы одно из равенств $b = 0$ или $x[i + 1] = x[b + 1]$. Если выполняется равенство $x[i + 1] = x[b + 1]$, то должно произойти присваивание $\beta[i + 1] \leftarrow b + 1$ независимо от того, будет ли b равно 0 или нет. С другой стороны, если $b = 0$, тогда должно произойти присваивание $\beta[i + 1] \leftarrow 0$. Эти варианты корректно обрабатываются условным оператором **if**, поскольку после выполнения цикла **while** имеем наибольшую возможную длину из последовательности (1.7). Таким образом, алгоритм 1.3.1 правильно вычисляет $\beta[i + 1]$ для любого $i \in 0..n - 1$. ■

Теперь вычислим время, необходимое для выполнения алгоритма. Все шаги цикла **for**, за исключением, возможно, цикла **while**, выполняются за константное время. Поэтому алгоритм 1.3.1 выполняется за время $\Theta(n)$ плюс общее время выполнения циклов **while**. Для получения оценки этого общего времени рассмотрим те значения, которые может принимать b во время выполнения алгоритма: b первоначально равно нулю и возрастает на 1 при выполнении оператора $b \leftarrow \beta[b]$ в начале $(i + 1)$ -й итерации цикла **for** только в том случае, если выполнялось присвоение $\beta[i + 1] \leftarrow b + 1$ в конце i -й итерации. Таким образом, b возрастает на единицу не более $n - 2$ раза, и каждое такое возрастание приходится на одну итерацию цикла **for**. В то же время уменьшение значения b может происходить только внутри цикла **while**. Поскольку количество уменьшений значения b (b всегда уменьшается не менее чем на 1) не может превосходить количества увеличений этого значения (увеличение всегда происходит в точности на 1), поэтому оператор присвоения $b \leftarrow \beta[b]$ внутри цикла **while** выполняется не более $n - 2$ раза. Таким образом, мы доказали следующую теорему.

Теорема 1.3.3. Алгоритм 1.3.1 для выполнения требует $\Theta(n)$ и фиксированного дополнительного и фиксированного дополнительного пространства памяти. ■

Обсуждение 1.3.1. В этом разделе мы впервые рассмотрели вопрос об асимптотической сложности алгоритма и об объеме памяти, необходимой для его выполнения. Эти вопросы будут многократно исследоваться на протяжении всей книги. Поэтому сейчас уделим достаточно внимания тем предположениям и рассуждениям, которые лежат в основе утверждений, подобных теореме 1.3.3.

- В алгоритме 1.3.1 необходимо хранить целые числа i и b и элементы массива β , которые также являются неотрицательными целыми числами. Каждое из

этих чисел не превышает $n - 1$, и поэтому для хранения одного числа надо зарезервировать $\lceil \log_2 n \rceil$ битов. Более точно, для каждого элемента массива β необходима память объемом $\log_2 n/w = O(\log n)$, где w — длина машинного слова компьютера.

Отсюда следует, что для хранения значений i и b также необходима память объемом $O(\log n)$, а для хранения всего массива β — память объема $O(n \log n)$. Таким образом, общее время выполнения алгоритма становится равным $O(n \log n)$, а не $\Theta(n)$, как указано в теореме 1.3.3. Мы разрешим это противоречие здесь и далее по всей книге путем дополнительного предположения, что размер задачи не может быть каким угодно большим, а всегда имеет некоторую, вполне реалистическую, границу. Так, если n — размер задачи, предположим существование такой константы k (по возможности, небольшой по величине), что выполняется неравенство $\log_2 n/w \leq k$. С практической точки зрения это не слишком ограничительное предположение, даже если мы примем, что $k = 1$. Действительно, поскольку современные компьютеры 32- или 64-разрядные, т.е. $w \geq 32$, то из неравенства $\log_2 n/w \leq 1$ получаем границу для числа n : $n \leq 2^{32} \approx 4,3 \times 10^9$. Если эта граница кажется слишком ограничительной, то можно принять $k = 2$ — тогда для n получаем границу $n \leq 2^{64} \approx 1,8 \times 10^{19}$.

Большое или маленькое последнее число? Предположим, что надо прочитать $1,8 \times 10^{19}$ машинных слов (строковых элементов) на компьютере со скоростью один миллиард слов в секунду. Тогда на выполнение всей операции чтения потребуется примерно 570 лет. Таким образом, только наши далекие потомки смогут увидеть завершение этой операции!

Я не являюсь оппонентом “блюстителей чистоты математических нравов”, которые требуют при вычислении асимптотических оценок времени выполнения или сложности алгоритмов обязательно умножать размер задачи n на “вычислительный фактор” $\log n$. В определенном смысле такой фактор действительно существует и его следует учитывать “в теории”, но в обычной вычислительной практике он излишен.

- Остановимся также на обозначениях O , Ω и Θ , которые используются на протяжении всей книги для характеристики асимптотических оценок времени и объема памяти, необходимых для выполнения алгоритмов. Мы предполагаем, что читатель знаком с этими обозначениями; если это не так, советуем обратиться к работам [134, 198].¹⁰ Здесь поясним различия между ними.

¹⁰Хотя автор отсылает читателя к соответствующим источникам, приведем краткое определение этих понятий. Говорят, что величина $T(n)$ имеет порядок $O(f(n))$, если существуют положительные константы c и n_0 , такие, что для всех n , больших или равных n_0 , выполняется неравенство $|T(n)| \leq cf(n)$. Величина $T(n)$ имеет порядок $\Omega(f(n))$, если существуют положительные константы c и n_0 , такие, что для бесконечного числа значений n (но не обязательно для всех значений n),

Требует некоторых пояснений применение обозначения $O(n)$. Если мы говорим “время выполнения (алгоритма) имеет порядок $O(n)$ ”, это означает, что мы рассматриваем время выполнения как функцию от размера n некоторой конкретной задачи, являющейся представителем *класса задач*. Другими словами, эта функция представляет характеристику целого класса задач, но “вычисленная” на некоторой конкретной (экстремальной) задаче. Как правило, утверждение “время выполнения (алгоритма) имеет порядок $O(n)$ ” выполняется только для очень больших значений n (фактически, для бесконечных значений n) и для любой задачи размера n из данного класса задач.

В книге мы трактуем время выполнения и требуемый объем памяти как функцию размера для всего данного класса задач. Тогда утверждение

$$\text{“время выполнения имеет порядок } \Theta(n)\text{”} \quad (1.8)$$

мы будем понимать так, что для всех задач данного класса, размер которых превышает некоторое фиксированное значение n_0 , алгоритм требует времени выполнения не меньше k_1n и не больше k_2n , где k_1 и k_2 также фиксированные (не зависящие от n) числа. Утверждение (1.8) более слабое, чем аналогичное утверждение с использованием $O(n)$. Фактически, здесь утверждается, что для всех задач данного класса, размер которых превышает значение n_0 , алгоритм требует времени выполнения не больше k_2n . В этом случае возможны ситуации, когда в данном классе задач есть такие задачи, на которых для выполнения алгоритма требуется объем памяти, пропорциональный, например $\log n$, $n^{1/2}$, $n/\log n$ или даже 1, т.е. асимптотически требуемый объем памяти может быть меньше n . Причем таким свойством могут обладать *все* задачи данного класса.

В книге будем использовать обозначение Θ для описания свойств алгоритмов только тогда, когда будут строго выполняться все условия утверждения (1.8). При использовании обозначения $O(f(n))$ (соответственно, обозначения $\Omega(f(n))$) будем понимать, что в данном классе задач существует ряд задач, на которых данное свойство имеет в точности порядок $f(n)$; но в этом классе задач также существует другой набор задач, на которых данное свойство имеет меньший (соответственно, больший) порядок, чем $f(n)$. ■

Итак, алгоритм 1.3.1 требует времени выполнения $\Theta(n)$ на всех задачах размера n . Поскольку любой алгоритм, вычисляющий массив граней, обязательно хотя бы один раз обратится к каждой из n позиций в исходной строке, то отсюда следует вывод, что алгоритм 1.3.1 **асимптотически оптимален**: не возможен

больших или равных n_0 , выполняется неравенство $|T(n)| \geq cf(n)$. Определение обозначения $\Theta(f(n))$ приведено в последующих абзацах. — *Примеч. ред.*

алгоритм, вычисляющий массив $\beta[1..n]$, с временем выполнения, меньшим $\Theta(n)$, при этом требующий фиксированного объема памяти для любой задачи размера n .

Как показано в разделе 4.1, алгоритм 1.3.1 также можно назвать, в определенном смысле, *онлайновым*¹¹, поскольку на каждом шаге алгоритма результат вычислений для текущей позиции i получается путем просмотра исходной строки слева направо. Поэтому без повторных пересчетов алгоритм позволяет получить результат не только для исходной строки, но и для любого ее префикса. Однако если говорить более строго, алгоритм 1.3.1 нельзя считать в полном смысле онлайновым, поскольку на любом шаге алгоритма может потребоваться вся исходная строка.

В качестве примера для вычисления граней, еще раз рассмотрим строковую последовательность из раздела 1.2:

$$f = \begin{array}{cccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ a & b & a & a & b & a & b & a & a & b & a & a & b \end{array}$$

Для этой строки массив $\beta[1..n]$ содержит значения 0011232345645. Отметим, что если в строке есть кратные подстроки (как, например, в строке f), то значения в массиве граней могут не составлять монотонно возрастающую последовательность, т.е. последовательные значения в массиве могут как убывать, так и возрастать. Например, в данном случае $\beta[6..8] = 323$ и $\beta[11..13] = 645$.

Способ вычисления массива β показывает, что *все* грани строки x (в действительности, грани всех префиксов строки x) можно вычислить на основе значений массива β . Это вытекает из равенства (1.6), которое показывает, что вторая по величине грань подстроки $x[1..i]$ является наибольшей гранью подстроки $x[1..\beta[i]]$. Сформулируем этот важный результат в виде теоремы.

Теорема 1.3.4. Массив граней β любой строковой последовательности x содержит всю информацию, для вычисления всех граней, периодов и порядков любого непустого префикса и порядков любого непустого префикса строки x . ■

Обсуждение 1.3.2. Массив граней — это одна из возможных структур строковых данных, представляющих информацию в линейной информации в линейной форме. Чтобы показать это, вернемся к примеру со строкой f . В данном случае с помощью значений массива β можно получить следующую информацию о гранях строки f .

¹¹Онлайновым считается алгоритм, который может работать в режиме реального времени или в режиме оперативной обработки данных (т.е. по мере поступления информации). — *Примеч. ред.*

i	длина граней
13	5, 2, 0
12	4, 1, 0
11	6, 3, 1, 0
10	5, 2, 0
\vdots	\vdots

Таким образом, информацию, содержащуюся в массиве граней, можно представить в явном виде в форме более удобной таблицы, подобной приведенной выше. Однако отметим, что для вычисления такой таблицы требуется уже не линейное (по n) время, а время порядка $\Theta(n \log n)$. ■

В заключение данного раздела приведем важный результат, который устанавливает числовую зависимость между различными периодами одной строки.

Теорема 1.3.5 (Лемма периодичности [87, 167]). Пусть p и q — два периода строки $x = x[1..n]$ и пусть $d = \gcd(p, q)$.¹² Если выполняется неравенство $p+q \leq n+d$, тогда d также является периодом строки x .

Доказательство этой теоремы интересно. Оно построено на сведениях условий теоремы, базирующихся на параметрах (d, p, q, n) к эквивалентному множеству условий на основе параметров $(d, p, q - n, n - p)$. Эта замена условий позволяет использовать алгоритм Евклида [135] для вычисления наибольшего общего делителя $d = \gcd(p, q)$, и, фактически, корректность доказательства теоремы зависит от корректности алгоритма Евклида. Напомним, что при $q > p$ этот алгоритм вычисляет значение d , основываясь на равенстве

$$d = \gcd(p, q) = \gcd(p, q - p), \tag{1.9}$$

и останавливается после конечного числа шагов при достижении $d = \gcd(d, d)$.

Обозначим через $H(d, p, q, n)$ условия теоремы так, как они сформулированы: p и q являются периодами строки $x[1..n]$, $d = \gcd(p, q)$ и выполняется неравенство $p + q \leq n + d$. Без потери общности можем положить, что $q > p$. (Если $q = p$, то теорема тривиальна.) В этих обозначениях запись $H(d, p, q - p, n - p)$ определяет условия, что p и $q - p$ являются периодами строки $x[1..n - p]$, $d = \gcd(p, q - p)$ и выполняется неравенство $q \leq n - p + d$. Покажем, что из условий $H(d, p, q, n)$ следуют условия $H(d, p, q - p, n - p)$. Это позволит вместо строки $x[1..n]$ рассматривать строку $x[1..n - p]$. Попытаемся сводить одни условия к другим последовательно и за конечное число шагов, пока в результате не будет достигнуто условие $H(d, d, d, n')$, которое выполняется тривиально, где d — период

¹²Здесь и далее $\gcd(p, q)$ обозначает наибольший общий делитель чисел p и q . — *Примеч. ред.*

строки $x[1..n']$. Далее покажем, что d является периодом $x[1..n-p]$ только тогда, когда d будет периодом x ; следовательно, если d будет периодом строки $x[1..n']$, то d будет и периодом строки $x[1..n]$.

Итак, пусть выполняются условия $H(d, p, q, n)$. Обозначим через $\beta_1 = n - p$ и $\beta_2 = n - q$ длины граней \mathbf{b}_1 и \mathbf{b}_2 , соответствующих периодам p и q . Тогда, как следует из леммы 1.3.1, \mathbf{b}_2 является гранью \mathbf{b}_1 . Другими словами, подстрока $x[1..\beta_1] = x[1..n-p]$ имеет грань $x[1..\beta_2]$ с периодом

$$\beta_1 - \beta_2 = n - p - n + q = q - p.$$

Поскольку $d = \gcd(p, q)$, то должно выполняться неравенство $d \leq q - p$. Из последнего неравенства и неравенства $p + q \leq n + d$ получаем

$$p \leq q - d < n - p. \quad (1.10)$$

Следовательно, поскольку строка x имеет период p , то и подстрока $x[1..n-p]$ должна иметь такой же период p . Далее надо показать, что $x[1..n-p]$ имеет как период p , так и период q . Из (1.9) имеем $d = \gcd(p, q - p)$, а неравенство $p + q \leq n + d$ можно переписать как

$$p + (q - p) \leq (n - p) + d.$$

Это доказывает, что из условий $H(d, p, q, n)$ вытекают условия $H(d, p, q - p, n - p)$.

Осталось показать, что d является периодом строки x тогда и только тогда, когда он является периодом подстроки $x[1..n-p]$. Очевидно, что если d является периодом x , то он также является периодом $x[1..n-p]$. Чтобы доказать утверждение в другую сторону, заметим, что если d является периодом $x[1..n-p]$, то, согласно (1.10), он также является периодом строки $x[1..p]$. Поскольку d должен быть делителем p , то справедливо равенство

$$x[1..p] = x[1..d]^{p/d}.$$

Так как подстрока $x[1..n-p]$ является гранью строки x , тогда d должен быть периодом строки

$$x[1..d]^{p/d} x[1..n-p] = x[1..p] x[p+1..n] = x[1..n],$$

что и требовалось. Таким образом, мы доказали, что подстрока $x[1..n-p]$ имеет период d только в том случае, если d является периодом исходной строки x .

Итак, после конечного числа шагов применения алгоритма Евклида для некоторого $n' \in d..n-p$ мы пришли к условиям $H(d, d, d, n')$, где d является периодом строки $x[1..n']$ и, как доказано, периодом исходной строки x . ■

В разделе 1.4 покажем одно из многих возможных применений этой теоремы. Здесь рассмотрим применение леммы периодичности к строке

$$x = abaaabaabaabaabaab,$$

которая имеет длину $n = 18$ и периоды $q = 12, p = 8$. Поскольку $d = \gcd(p, q) = 4$ и $p + q = 20 \leq n + d = 22$, то отсюда следует, что $d = 4$ также является периодом строки x .

Иногда лемма периодичности выполняется даже тогда, когда не все ее условия реализованы. Например, строка

$$y = abaabaabaabaab$$

имеет длину $n = 18$ и периоды $q = 12, p = 8$. Но в этом случае $p + q = 20 > n + d = 18$. Однако все равно $d = 4$ является периодом строки x .

Необходимость условия $p + q \leq n + d$ в теореме 1.3.5 показывает пример строки $z = abaaba$ длиной $n = 6$ и с периодами $q = 5, p = 3$. В этом примере $p + q = 8 = n + d + 1$ и $d = \gcd(3, 5) = 1$ не является периодом строки z .

Недавно лемма периодичности была обобщена на случай строк с тремя периодами [43].

Упражнения 1.3

1. Определите асимптотическую сложность “очевидного” алгоритма нахождения наибольших граней данной строки x . Приведите пример строки произвольной длины n , на которой этот алгоритм будет выполняться наиболее долго. Подсчитайте, сколько необходимо произвести сравнений между буквами строки в самом худшем случае.
2. Поскольку необходимо найти наибольшие грани строки x , то более эффективным кажется алгоритм, в котором возможные грани проверяются в *нисходящем* порядке относительно длин этих граней. Разработайте такой алгоритм и определите его асимптотическую сложность. В каком случае время выполнения такого алгоритма будет одинаковым с временем выполнения “обычного” алгоритма? В каком случае время его выполнения будет линейным по длине исходных строк?
3. Если массив β рассматривать как строку, то на каком алфавите определена эта строка?
4. Докажите утверждение, сделанное в этом разделе, что если b является гранью строки x , а b' — гранью подстроки b , то b' является гранью строки x .
5. Докажите следующую лемму (доказанную Яндонгом Янгом (Jiandong Jiang)).

Если $x = bx'b$, где b — примитивная, а x' — пустая или примитивная строка и $x' \neq b$, тогда b может быть только непустой гранью x .

6. При доказательстве леммы 1.3.1 мы не показали, почему не существует грани строки $x[1..n]$ длины, отличной от $\beta^j[n]$, $j = 1, 2, \dots, k$. Завершите доказательство этой леммы.
7. Разработчики программного обеспечения располагают богатым арсеналом формальных методов доказательства корректности программ. В частности, при исследовании циклических структур различают такие их составляющие.
- **Инвариант цикла** — условное выражение, принимающее значение “истина” в начале или в конце выполнения цикла.
 - **Вариант цикла** — выражение, которое принимает целые положительные значения и уменьшается при каждом выполнении цикла.
- Найдите в двух циклах алгоритма 1.3.1 инвариант и вариант цикла.
8. Примените алгоритм 1.3.1 к каждой из следующих строк и вычислите для них массивы граней $\beta[1..n]$.
- а) $abcdefg$;
 - б) a^n ;
 - в) $(ab)^n$;
 - г) $abaababaabaab$;
 - д) $(aba)^3 a(aba)^2 a(aba)^3$;
 - е) $a^4 b a^5 b$.
9. Оператор **if** в алгоритме 1.3.1 немного раздражает, поскольку в нем дублируется сравнение букв, которое уже выполнялось в операторе **while**. Путем введения дополнительного элемента $\beta[0]$ в массиве граней найдите способ исключить этот оператор.
10. В упражнении 1.2.24 утверждалось, что более правильно “нормальную форму” называть “левой нормальной формой”, там же введено понятие “правой нормальной формы”. Аналогично, поскольку массив граней содержит значения длин наибольших граней всех *префиксов* строки x , его также имеет смысл назвать “левый массив граней”. Тогда **правый массив граней** будет содержать значения длин наибольших граней всех *суффиксов* строки x . Напишите алгоритм вычисления правого массива граней с линейным временем выполнения.
- Замечание.** Это упражнение представляет не только академический интерес. Правый массив граней находит широкое применение, например, в известном строковом алгоритме Бойера–Мура, описанном в разделе 7.2.
11. Обозначим через ρ правый массив граней строки x длины n и через β_T — массив граней транспонированной строки

$$x_T = x[n]x[n-1] \dots x[1].$$

Докажите, что для каждого $i \in 1..n$ выполняется $\rho[i] = \beta_T[n - i + 1]$.

12. На основе леммы 1.3.1 разработайте алгоритм, который определяет, является ли данная строка $y = y[1..n]$ целых неотрицательных чисел массивом граней какой-нибудь строки на некотором алфавите. Ваш алгоритм должен выполняться за время порядка $\Theta(n)$, докажите, что это так и есть. И, конечно, докажите корректность алгоритма.

Предупреждение. Это трудное упражнение, которое лучше рассматривать как исследовательскую задачу. Она тесно связана с материалом разделов 4.2 и 4.3 и ведет к двум опубликованным работам [90, 79]. Несмотря на то что уже сделано по этой теме, здесь есть еще обширное поле для исследований.

13. Обозначим через x' строку, полученную из заданной строки x путем замены букв, входящих в эту строку, другими буквами, также входящими в эту строку. (Например, если $x = abcab$, то после замены $a \rightarrow c$, $b \rightarrow a$, $c \rightarrow b$ получим строку $x' = cabca$.) Покажите, что массив граней строки x' совпадает с массивом граней строки x . Этот результат можно интерпретировать следующим образом: в любом алфавите, содержащем более одной буквы, не существует строк, которые однозначно определялись бы своим массивом граней. (Более подробно этот вопрос рассмотрим в разделе 4.3.)
14. В этом разделе было высказано утверждение, что алгоритм 1.3.1 асимптотически оптимален. Сохранится ли справедливость этого утверждения, если предположить, что алгоритм 1.3.1 будет вычислять только $\beta[n]$?
15. В обсуждении, ведущем к теореме 1.3.3, было сделано несколько туманное утверждение, что “количество уменьшений значения $b \dots$ не может превосходить количества увеличений этого значения”. Чтобы прояснить этот вопрос, рассмотрим функции $x(v) = v + 1$ и $y(v) = v - 1$, определенные на всем множестве действительных чисел. Обозначим через s строку, состоящую из j и k ($j, k \geq 0$) последовательных результатов вычисления функций x и y соответственно. Тогда $|s| = j + k$ и $s(v)$ равно результирующему значению этих вычислений. Например, для $s = xyx|s| = 3$ и $s(0) = 1$.
Докажите, что $s(v) - v = j - k$, и покажите, как этот результат связан с теоремой 1.3.3.
16. Докажите более общий результат, подобный предыдущему упражнению, для случая функций $x(v) = v + r$ и $y(v) = v - s$, где r и s — положительные действительные числа.
17. Покажите, что если выполняются условия леммы периодичности и неравенство $p < q$, тогда $p \leq n/2$.

1.4 Строковые петли

Строковые петли часто преподносят неожиданности строковым алгоритмам, даже если они открыто не связаны с решаемой задачей. В этом разделе рассмотрены некоторые свойства строковых петель.

Как упоминалось в разделе 1.1, строковые петли обычно определяются через линейные строки. Будем использовать обозначение $C(x)$ для строковой петли, сформированной из линейной строки x путем сочленения ее самого левого элемента $x[1]$ с самым правым элементом $x[n]$. Таким образом строка x преобразуется в строковую последовательность, не имеющую ни самого левого, ни самого правого элементов. Например, если такое преобразование применить к линейной строке $f = abaababa$, тогда получим петлю $C(f)$ вида

$$\begin{array}{ccccc} & & a & & \\ & a & & b & \\ & & & & \\ b & & & & a \\ & & & & \\ a & & & a & \\ & & b & & \end{array}$$

В алгоритмах, формирующих петли, часто $C(x)$ выражают через циклические сдвиги строки x . Для данной строки $x = x[1..n]$ и произвольного целого $j \in 0..n - 1$ j -й **циклический сдвиг** (rotation или cyclic shift) строки x определяет строку¹³

$$R_j(x) = x[j + 1..n]x[1..j].$$

Из этого определения следует, что $R_0(x) = x$.

Приведем восемь сдвигов строки $f = abaababa$:

$$\begin{aligned} R_0(f) &= abaababa, \\ R_1(f) &= baababaa, \\ R_2(f) &= aababaaab, \\ &\vdots \\ R_7(f) &= aabaabab. \end{aligned}$$

Отметим, что сдвиги имеют подстроки, которые не являются подстроками исходной строки. Например, здесь строка $aabaa$ не является подстрокой строки f , однако является подстрокой сдвига $R_7(f)$ и, следовательно, петли $C(f)$. Поэтому в процессе формирования петли важно сохранить все возможные циклические

¹³Здесь и далее термином “циклический сдвиг” обозначается как операция циклического сдвига, так и строка, полученная в результате применения этой операции. Как правило, это не приводит к недоразумениям, так как из контекста всегда ясно, о чем идет речь. — *Примеч. ред.*

сдвиги исходной строки. На практике для этого достаточно сохранить строку x^2 , которая содержит все сдвиги $R_j(x)$ в виде своих подстрок.

Теперь можно более четко определить, что мы будем понимать под термином *подстрока петли* $C(x)$ — это подстрока любого сдвига строки x или, что равнозначно, подстрока (длины не более n) строки x^2 . Но поскольку петля не имеет ни самого левого, ни самого правого элементов, особые подстроки, введенные в разделе 1.2, — префикс, суффикс и грань — можно определить только по отношению к конкретным сдвигам, но не к петле в целом. Такое понятие, как период, которое определялось через грани, также нельзя применить к петле. Вместо этого следующие теоремы показывают, как посредством циклических сдвигов можно определить различные свойства петель $C(x)$, если порождающая линейная строка x не является кратной. Сначала докажем такую теорему.

Теорема 1.4.1. Пусть x — строковая последовательность длиной $n \geq 3$ на двоичном алфавите ($\alpha = 2$).

- а) если строка x является кратной, то таким же будет каждый циклический сдвиг этой строки;
- б) если строка x не является кратной, тогда существует по крайней мере один циклический сдвиг этой строки, который будет периодическим.

Доказательство. Сначала докажем утверждение а). Из условия кратности строки x следует, что $x = u^r$, где $r > 1$ и существует такое целое $p < n$, что $n = pr$ и $u = x[1..p]$. Тогда

$$R_1(x) = x[2..n]x[1] = (u[2..p]u[1])^r.$$

Отсюда по индукции получаем требуемое утверждение.

Для доказательства утверждения б) попробуем построить строку x , такую, чтобы любой ее циклический сдвиг был примитивной строкой. В такой строке $x[1] \neq x[n]$, иначе строка x будет периодической. Без потери общности положим, что $x[1] = a$ и $x[n] = b$. Тогда должно выполняться равенство $x[2] = b$, в противном случае $R_1(x)[1] = R_1(x)[n] = a$, и поэтому $R_1(x)$ будет периодической строкой. Подобные рассуждения ведут к заключению, что $x[n-1] = a$. Но тогда уже при $n = 3$ получаем противоречие: $x[2] \neq x[3-1]$. Отсюда следует, что при $n > 3$ строка x имеет грань ab , которая является периодической. ■

Как указано в разделе 3.3, можно построить некратную строку произвольной длины на алфавите из трех букв. Если мы возьмем одну из таких строк длиной $n-1$ и затем присоединим к ней одну букву, которой пока нет в этой строке, тогда, как показано в упражнении 1.4.3, будем иметь строку на алфавите из четырех букв, каждый сдвиг которой является примитивной строкой. Поэтому утверждение б) теоремы 1.4.1 нельзя обобщить на алфавиты размером 4 и выше. Но утверждение а) теоремы 1.4.1 можно обобщить, что сделано в упражнении 1.4.2.

Другой полезный результат связан с циклическими сдвигами, кратностью и периодичностью строк. Согласно этому результату, если строка x не кратная, то все ее циклические сдвиги различны. Доказательство этого результата основано на лемме периодичности.

Теорема 1.4.2. Пусть x — произвольная строковая последовательность длины n с минимальным периодом p^* . Тогда для любого целого числа $j \in 1..n - 1$ равенство $R_j(x) = x$ будет выполняться только в том случае, когда x является кратной строкой и p^* является делителем j .

Доказательство. Для доказательства достаточности утверждения теоремы нетрудно показать, что если выполняется равенство $x = x[1..p^*]^{r^*}$, где $r^* = n/p^*$ — целое число, тогда для $j = p^*, 2p^*, \dots, (r^* - 1)p^*$ справедливо $R_j(x) = x$.

Для доказательства необходимости предположим, что $R_j(x) = x$. Тогда $x[i + j] = x[i]$ для всех целых $i \in 1..n - j$, поэтому x имеет грань $x[1..n - j]$ и период j . Теперь применим теорему 1.3.5 с параметризованными условиями $H(d, j, n - j, n)$, где $d = \gcd(j, n - j)$. Поскольку выполняется $j + (n - j) < n + d$, то на основании теоремы делаем заключение, что d является периодом строки x , т.е.

$$x = x[1..n - j]x[1..j] = x[1..d]^{(n-j)/d}x[1..d]^{j/d} = x[1..d]^{n/d}.$$

Таким образом, строка x является кратной.

Чтобы показать, что p^* делит d и, следовательно, делит j , снова применим теорему 1.3.5 с условиями $H(d', p^*, d, n)$, где $d' = \gcd(p^*, d)$. Поскольку $p^* \leq d \leq n/2$, из теоремы следует, что $p^* + d < n + d'$, тогда d' является периодом строки x . Это означает, что $d' = p^*$ и, следовательно, p^* делит d , что и требовалось доказать. ■

Теорема 1.4.2 удручает, поскольку из нее следует, что, за исключением случая кратных строк, для проверки равенства $C(x) = C(y)$ необходимо проверить все циклические сдвиги исходных строк x и y . Такая проблема возникает в компьютерной графике и вычислительной геометрии, где выпуклый многоугольник можно представить посредством двух различных строк: одна содержит последовательность длин сторон многоугольника при его обходе по часовой стрелке, другая — последовательность углов между смежными сторонами, также полученную при обходе многоугольника по часовой стрелке. Для решения задачи, будут ли два представленных многоугольника конгруэнтными, надо проверить, будут ли равны строки, представляющие эти многоугольники (будь то строки длин сторон или строки углов). Пусть строка x представляет один многоугольник, а строка y — другой. Поскольку многоугольники могут иметь различную ориентацию, первое значение $x[1]$ строки x может не совпадать с первым значением $y[1]$ строки y .

Другими словами, даже если в действительности $C(x) = C(y)$, то совсем не обязательно, что $x = y$. Пример такой ситуации показан на рис. 1.1.

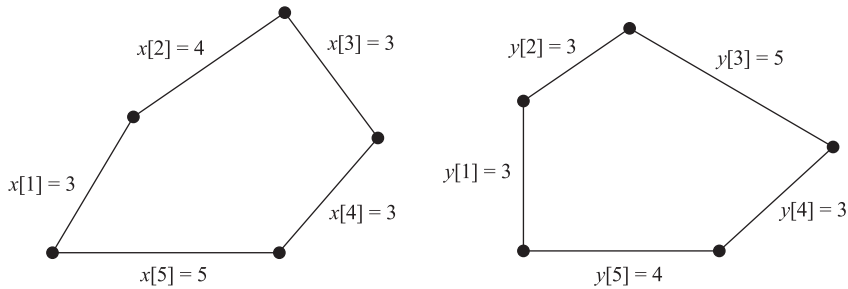


Рис. 1.1. Два конгруэнтных многоугольника

Естественным подходом к решению этой проблемы было бы введение “канонической формы” для строковых петель такой, чтобы она позволяла “напрямую” сравнивать различные петли. На практике такую роль могут выполнять циклические сдвиги исходной строки, которые можно вычислить достаточно эффективно. Если строки определены на упорядоченном алфавите, то решением проблемы может быть следующее определение.

Определение 1.4.3. Канонической формой строковой петли $C(x)$ на упорядоченном алфавите называется лексикографически наименьший циклический сдвиг $R_j(x)$ строки x , т.е. $R_j(x) \leq R_i(x)$ для всех $i \in 0..n - 1$. ■

Каноническая форма строковых петель очень тесно связано еще с одним понятием.

Определение 1.4.4. Непустая строка w на упорядоченном алфавите называется линдонским словом, если $w < R_j(w)$ для всех целых $j \in 1..n - 1$. ■

Строка, состоящая из одной буквы, будет линдонским словом¹⁴, поскольку, согласно теореме 1.4.2, такая строка не может быть кратной. Далее мы увидим, что строка ab является линдонским словом, строка ba — нет. Аналогично строка aab является линдонским словом, а строки aba и baa — нет. В общем случае, основываясь на следствии теоремы 1.4.2, которое заключается в том, что если строка x некратная, то все ее циклические сдвиги различны, можно сделать вывод: у некратной строки существует только один сдвиг, который будет линдонским словом. Таким образом, если строка x некратная, тогда существует линдонское

¹⁴Такое название строкам специального вида дано в честь Роджера Линдона (Roger C. Lyndon), который предложил особую факторизацию строк, которая позже получила название линдонской декомпозиции. Эта факторизация описана дальше в этом разделе. — *Примеч. ред.*

слово $w = R_j(x)$, которое в точности совпадает с канонической формой петли $C(x)$. Это означает, что множество всех линдонских слов на данном упорядоченном алфавите A совпадает с множеством всех канонических форм, которые не являются кратными.

Далее докажем три леммы (взяты из работы [80]), которые описывают основные свойства линдонских слов. Эти леммы необходимы для теоремы о линдонской декомпозиции.

Лемма 1.4.5. Каждое линдонское слово примитивно.

Доказательство. Необходимо показать, что не существует линдонских слов, которые имели бы непустые грани. Предположим противное — существует линдонское слово w , такое, что $w = uv' = v''u$, где u , v' и v'' — непустые строки. Очевидно, что $|v'| = |v''|$. Отсюда и из определения линдонского слова следует, что

$$uv' = w < uv'' \text{ и } v''u = w < v'u.$$

Итак, получаем два неравенства: $v' < v''$ и $v'' < v'$ — т.е. пришли к противоречию. ■

Лемма 1.4.6. Строковая последовательность w будет линдонским словом только тогда, когда для любого непустого префикса v строки w выполняется неравенство $w < v$.

Доказательство. Для доказательства достаточности условий леммы заметим, что если $w = uv < v$ для непустых строк u и v , тогда $uv < vu$.

Для доказательства необходимости предположим, что w — линдонское слово и $w = uv$, где u и v — непустые строки. Тогда $uv < vu$. Очевидно, что $uv \neq v$, поскольку $|uv| > |v|$. Предположим, что $uv > v$. Возможны два взаимоисключающих случая: v является префиксом строки uv и v не является префиксом этой же строки. Но первый случай невозможен, поскольку тогда v является гранью строки w , что противоречит лемме 1.4.5. Второй случай тоже невозможен, так как тогда $w[1..|v|] > v$, откуда следует неравенство $w > uv$, противоречащее свойствам линдонских слов. Таким образом, мы пришли к противоречию с предположением $uv > v$. Следовательно, $uv < v$, что и требовалось доказать. ■

Лемма 1.4.7. Пусть w_1 и w_2 являются линдонскими словами. Тогда строка w_1w_2 будет линдонским словом только в том случае, когда $w_1 < w_2$.

Доказательство. Необходимость условий леммы следует непосредственно из леммы 1.4.6: если w_1w_2 является линдонским словом, тогда $w_1 < w_1w_2 < w_2$.

Для доказательства достаточности обозначим через v непустой собственный префикс строки w_1w_2 и предположим, что $w_1 < w_2$. Возможна одна из следующих ситуаций.

- $v = w'_1 w_2$, где w'_1 — непустой собственный префикс строки w_1 . Поскольку w_1 линдонским словом, то на основании леммы 1.4.6 имеем $w_1 < w'_1$. Так как $|w_1| > |w'_1|$, то отсюда получаем $w_1 w_2 < w'_1 w_2 = v$.
- $v = w_2$. Если w_1 не является префиксом w_2 , тогда $w_1 w_2 < w_2 = v$, поскольку $w_1 < w_2$. Если же w_1 является префиксом w_2 , тогда можно записать $w_2 = w_1 w''_2$, где w''_2 — собственный суффикс строки w_2 . Поскольку w_2 является линдонским словом, можно применить лемму 1.4.6, вследствие которой имеем $w_2 < w''_2$. Следовательно, $w_1 w_2 < w_1 w''_2 = w_2 = v$.
- $v = w'_2$, где w'_2 — собственный суффикс строки w_2 . Поскольку w_2 является линдонским словом, из леммы 1.4.6 следует, что $w_1 < w'_2$. Так как w_2 не является префиксом w'_2 и имеется условие $w_1 < w_2$, отсюда следует, что w_1 не может быть префиксом w'_2 . Поэтому $w_1 w_2 < w_2 < w'_2 = v$, что и требовалось.

В каждом из рассмотренных случаев получаем, что $w_1 w_2 < v$, где v — суффикс строки $w_1 w_2$. На основании леммы 1.4.6 заключаем, что $w_1 w_2$ является линдонским словом. На этом завершается доказательство леммы. ■

Чтобы использовать полученные результаты, нам необходимо еще одно определение.

Определение 1.4.8. Пусть x — непустая строковая последовательность на упорядоченном алфавите A . Декомпозиция $x = w_1 w_2 \cdots w_k$ называется **линдонской**, если все w_i ($i = 1, 2, \dots, k$) являются линдонскими словами, которые удовлетворяют неравенствам $w_1 \geq w_2 \geq \cdots \geq w_k$. ■

Теперь можно сформулировать и доказать одну замечательную теорему [48].

Теорема 1.4.9. Для любой непустой строковой последовательности x на упорядоченном алфавите существует линдонская декомпозиция, причем единственная.

Доказательство. Сначала покажем, что существует по крайней мере одна линдонская декомпозиция строки x . Начнем с первоначальной декомпозиции

$$x = w_1^{(0)} w_2^{(0)} \cdots w_n^{(0)},$$

где факторы $w_i^{(0)} = x[i]$, $i \in 1..n$ представляют отдельные буквы. Здесь все факторы $w_i^{(0)}$ являются линдонскими словами, но сама декомпозиция не обязательно будет линдонской, поскольку возможны неравенства $w_j^{(0)} < w_{j+1}^{(0)}$ для некоторых номеров j . Для таких номеров сочленение строк $w_j^{(0)} w_{j+1}^{(0)}$ будет линдонским словом, что следует из леммы 1.4.7. Объединяя такие факторы, получим новую декомпозицию

$$x = w_1^{(1)} w_2^{(1)} \cdots w_{n_1}^{(1)},$$

где каждый фактор $w_i^{(1)}$ является линдонским словом и $n_1 \leq n$. Если $n_1 < n$ и есть последовательные факторы, для которых выполняется неравенство $w_j^{(1)} < w_{j+1}^{(1)}$, продолжаем объединение таких факторов и получаем новую декомпозицию. Через конечное число таких шагов получим линдонскую декомпозицию

$$x = w_1^{(s)} w_2^{(s)} \cdots w_{n_s}^{(s)},$$

где будут выполнены условия $w_1^{(s)} \geq w_2^{(s)} \geq \cdots \geq w_{n_s}^{(s)}$.

Теперь предположим, что линдонская декомпозиция не единственная, т.е. существуют две различные декомпозиции

$$x = w_1 w_2 \cdots w_k \text{ и } x = w'_1 w'_2 \cdots w'_k.$$

Без потери общности можно предположить, что существует минимальное положительное целое число $i < \min\{k, k'\}$, такое, что для некоторой непустой подстроки v и для всех $j \in 1..i-1$ выполняется равенство $w'_j = w_j$. Тогда существует единственное целое число $r \geq 1$ и непустая строка v^* , такая, что

$$w_i w_{i+1} \cdots w_{i+r-1} v^* = w_i v \text{ и } |w_i w_{i+1} \cdots w_{i+r}| \geq |w_i v|.$$

Обозначим $w^* = w_i w_{i+1} \cdots w_{i+r-1}$, так что $w'_i = w^* v^*$. Поскольку w'_i является линдонским словом, из леммы 1.4.6 следует, что $w'_i < v^*$. Строка w^* — собственный префикс w'_i , поэтому $w^* < v^*$.

Далее, поскольку v^* является префиксом w_{i+r} , можно записать $v^* \leq w_{i+r}$. Так как w_i — префикс w^* , то $w_i \leq w^*$. Поэтому

$$w_i \leq w^* < v^* \leq w_{i+r}.$$

Последние неравенства противоречат предположению, что $x = w_1 w_2 \cdots w_k$ является линдонской декомпозицией. ■

В конце этого раздела вернемся к сделанному выше замечанию о том, что строка x не должна быть кратной и что каноническая форма петли $C(x)$ совпадает с циклическим сдвигом строки x и этот сдвиг является линдонским словом. Предположим, что имеется алгоритм (например, один из описанных в разделе 6.1), вычисляющий линдонскую декомпозицию для произвольной строки x . Запишем линдонскую декомпозицию в более сжатой форме

$$x = w_1^{q_1} w_2^{q_2} \cdots w_k^{q_k}, \tag{1.11}$$

где $w_1 > w_2 > \cdots > w_k$ и q_i — целые положительные числа, $i = 1, 2, \dots, k$. Как показано в упражнении 1.4.14, с помощью нашего алгоритма можно сразу получить все значения q_i .

Кратко рассмотрим задачу определения в строке x позиции $\text{MSP}(x) = i^*$ (эта позиция называется **минимальной начальной точкой**¹⁵), такой, что $R_{i^*-1}(x)$

¹⁵Отсюда название функции MSP — это сокращение от англ. Minimum Starting Point, что и переводится как *минимальная начальная точка*. — Примеч. ред.

будет минимумом по всем циклическим сдвигам строки x . Как мы увидим дальше, позиция i^* определяется однозначно и $R_{i^*-1}(x)$ является линдонским словом только в том случае, когда строка x не будет кратной. Сначала рассмотрим несколько специальных случаев.

- Если $k = 1$, тогда $\text{MSP}(x) = 1$.
- Если $k = 2$, тогда $\text{MSP}(x) = q_1|w_1| + 1$.
- Если $w_k^{q_k}$ не является префиксом $w_{k-1}^{q_{k-1}}$, тогда $\text{MSP}(x) = n - q_k|w_k| + 1$.
- Если строка x кратная, т.е. $x = u^m$, $m > 1$, тогда $\text{MSP}(x) = \text{MSP}(u)$.

Теперь рассмотрим применение алгоритма к строке x^2 . Если строка x не кратная, алгоритм создаст декомпозицию $v_1 R_{i^*-1}(x) v_2$, где v_1 и v_2 — декомпозиции строк $x[1..i^* - 1]$ и $x[i^*..n]$ соответственно и $R_{i^*-1}(x)$ является линдонским словом. В этом случае для того, чтобы удостовериться в том, что $\text{MSP}(x) = |v_1| + 1$, необходимо просмотреть только линдонское слово длины n . В случае, когда $x = u^m$ (m целое и больше 1), необходимо просмотреть блок одинаковых линдонских слов $w_i^{q_i}$, таких, что $q_i|w_i| = n$, и снова получим результат $\text{MSP}(x) = |v_1| + 1$. Таким образом, любой алгоритм для вычисления линдонской декомпозиции можно (с минимальными изменениями) использовать для вычисления канонических форм. Однако, как показано в упражнении 1.4.16, основываясь только на линдонской декомпозиции, нельзя непосредственно определить, какой фактор в декомпозиции определяет $\text{MSP}(x)$. Мы вернемся к этому вопросу в разделе 6.2.

В разделе 6.2 также показано, что линдонская декомпозиция является основой для решения многих задач, в частности, эффективного нахождения лексикографически наименьшего и лексикографически наибольшего суффиксов любого префикса строки.

Упражнения 1.4

1. Покажите, что для сохранения всей информации о петле $C(x)$ достаточно сохранить строку $x^{(2n-1)/n}$.
2. По индукции докажите, что для произвольной непустой строки $x = u^k$ и любого целого $j \in 0..|u| - 1$ выполняется равенство $R_j(x) = (R_j(u))^k$. (Запишите в явном виде “индукцию”, упоминающуюся при доказательстве утверждения *a*) теоремы 1.4.1. Этим вы докажете, что утверждение *a*) этой теоремы справедливо для любого алфавита.)
3. Пусть x — некратная строка длиной $n - 1$ на алфавите $A = \{a, b, c\}$. Докажите утверждение, сделанное в этом разделе, что любой циклический сдвиг строки x будет примитивной строкой.
4. Внимательный читатель, конечно, обратил внимание на то, что при рассмотрении вопроса о существовании примитивных циклических сдвигов

у кратных и некратных строк мы не исследовали случай, когда исходная строка примитивная на алфавите, состоящем из трех букв. Причина такого упущения станет ясной, если обратиться к следующему результату, полученному путем машинного моделирования. Оказывается, не существует примитивных строк длиной $n = 5, 7, 9, 10, 14$ и 17 , у которых все циклические сдвиги примитивны. Для всех других значений $n \in 1..19$ такие строки существуют. Таким образом, имеем интересную исследовательскую задачу, поставленную Джейми Симпсоном (Jamie Simpson): определить значения n , для которых существуют примитивные строки на алфавите из трех букв, у которых все циклические сдвиги примитивны.

Дополнительная информация. После длительных исследований этой задачи, ее автор пришел к выводу, что для всех n , больших 17 , такие примитивные строки существуют, но доказать это очень трудно. Однако недавно Джеймс Карри (James Currie) [74] нашел другой подход к решению этой задачи.

5. Без привлечения леммы периодичности покажите, что если для двух непустых строк x и y выполняется равенство $xy = yx$, тогда существуют положительные целые числа k_1 и k_2 и строка u , такие, что и $y = u^{k_2}$. На этом строится альтернативное доказательство утверждения (см. теорему 1.4.2), что если $x = R_j(x)$ для некоторого $j \in 1..n - 1$, тогда строка x кратная.
6. Вдумчивый читатель заметит, что в доказательстве леммы 1.4.6 присутствует “потерянная лемма”, т.е. доказательство этой леммы можно сделать более простым и понятным, если ранее доказать следующее утверждение.

Пусть x и y — строки на упорядоченном алфавите A . Если $x > y$ и y не является префиксом строки x , тогда для любых строк $z, z' \in A^*$ справедливо неравенство $xz < yz'$.

Докажите эту “потерянную лемму”.

7. Подобно предыдущему упражнению, в доказательстве леммы 1.4.7 также можно выделить два утверждения, использование которых упростит доказательство леммы. Вот эти утверждения.

■ Если $x < y$ и x не является префиксом строки y , тогда для любых строк $z, z' \in A^*$ справедливо неравенство $xz < yz'$.

■ Если $x < y$ и $|x| \geq |y|$, тогда для любых строк $z, z' \in A^*$ справедливо неравенство $xz < yz'$.

Докажите сформулированные утверждения.

8. В продолжение предыдущего упражнения докажите, что если строки $x, y, x', y' \in A^*$ такие, что $xy < x'y'$ и строки x и x' не являются префиксами друг друга, тогда выполняется неравенство $x < x'$.

9. Первая часть доказательства теоремы 1.4.9 предлагает схему алгоритма для определения линдонской декомпозиции. Разработайте такой алгоритм и определите верхнюю границу его асимптотической сложности.
10. Найдите линдонскую декомпозицию для следующих строк.
- а) $xyzabc$;
 - б) $abcxyz$;
 - в) $abaababa$;
 - г) $(baa)^n$;
 - д) $abacabadabacabac$;
 - е) $abacabadabacabad$;
 - ж) $abacabadabacabab$;
 - з) $abacabadabacabaa$.
11. Для каждой строки, перечисленной в предыдущем упражнении, найдите каноническую форму для соответствующих петель.
12. Приведите примеры
- а) примитивной строки, которая содержит кратную составляющую;
 - б) некратной строки, которая не является примитивной.
13. Назовем строковую петлю $C(x)$ **некратной**, если любой циклический сдвиг строки x является некратной строкой. Докажите следующие два утверждения.
- Каноническая форма петли $C(x)$ будет примитивной строкой только тогда, когда строка x не является кратной.
 - Петля $C(x)$ будет некратной только тогда, когда любой циклический сдвиг строки x будет примитивной строкой.
14. Предположим, что алгоритм построения линдонской декомпозиции выполняется над строкой x слева направо, и пусть на некотором этапе алгоритма получено следующее представление строки x :

$$x = w_1^{q_1} w_2^{q_2} \dots w_j^{q_j} u^m \lambda \dots$$

для некоторых целых $j \geq 0$, $m \geq 2$ и некоторой буквы λ . Опишите дальнейший процесс построения линдонской декомпозиции для следующих случаев.

- а) $\lambda < u[1]$.
- б) $\lambda = u[1]$.
- в) $\lambda > u[1]$.

15. Докажите правильность выводов для $\text{MSP}(x)$ в четырех “специальных случаях”, приведенных в конце этого раздела.
16. “Соблазнительно” предположить (автор однажды “купился” на это), что минимальную начальную точку можно вычислить как $\text{MSP}(x) = i_k$, где i_j — начальная точка строки $w_j^{q_j}$ в линдонской декомпозиции (1.11), $j = 1, 2, \dots, k$. Приведите пример строки, где $\text{MSP}(x) = i_{k-1}$ или даже $\text{MSP}(x) = i_{k-2}$.

ГЛАВА 2

Паттерны? Что такое паттерны?

Я еще не настолько растерял свои познания в лексикографии, чтобы забыть, что слова — это дочери земли, а вещи — сыновья Небес.

— Сэмюэль Джонсон (1709–1784).
Предисловие к словарю

В главе 1, посвященной свойствам строковых последовательностей, только раз упоминалось о паттернах. Пришло время разобраться, что же это такое. В этой главе сделан обзор паттернов, которые будут вычисляться алгоритмами, описанными в частях II–IV книги, и представлено компактное определение основных задач, решаемых этими алгоритмами. Соответственно, паттерны раздела 2.1 рассмотрены в части II, паттерны раздела 2.2 — в части III и паттерны раздела 2.3 — в части IV.

В этой главе, как и во всей книге, $x = x[1..n]$ обозначает строковую последовательность длиной n .

2.1 Внутренние паттерны

Как упоминалось в предисловии, внутренние паттерны (intrinsic patterns) в определенном смысле являются “собственным” или “природным” свойством строковых последовательностей — не каждая строка содержит частный паттерн (specific pattern), например подстроку *bb*, или характеристический паттерн (generic pattern), такой как кратные строки. С другой стороны, внутренние паттерны можно найти

во всех строках, и они также в определенной степени характеризуют строковые последовательности. Так, “нормальная форма” из раздела 1.2 — это паттерн, который можно рассматривать как стандартное описание любых строк, которое не зависит от других возможных описаний строк. Интересно отметить, что существует много полезных и полностью различных таких “стандартных описаний”.

Нормальная форма определена в разделе 1.2, она вычисляется путем двойной обработки массива граней, как показано в разделе 1.3. Сформулируем задачу вычисления нормальной формы как задачу нахождения периодических структур в строках.

Задача 2.1 (Вычисление нормальной формы). Выразить строковую последовательность x в форме u^r , где u — префикс строки x , выбранный так, чтобы максимизировать значение r (раздел 1.3). ■

Хотя массивы граней не рассматриваются как паттерны, по сути они являются своеобразными внутренними паттернами соответствующих строковых последовательностей. Основой для такой интерпретации массивов граней служит теорема 1.3.4, которая показывает роль этих массивов в вычислении граней, периодов, порядков и нормальных форм.

Задача 2.2 (Вычисление массива граней). Вычислить наибольшие грани каждого префикса строки x и представить их в виде массива (строки) $\beta = \beta[1..n]$ (раздел 1.3). ■

В главе 5 мы изучим два вида деревьев: дерево граней и дерево суффиксов.¹ Дерево граней — это структура, которая уже содержится в скрытом виде в массиве граней, как показано ниже. Смысл сделать эту структуру явной заключается в том, чтобы помочь решить проблему “оболочек”, рассмотренную в разделе 2.3. В главе 13 показано, что каждая строка имеет внутренний шаблон, так называемый *массив оболочек*, который является структурным “образом” массива граней. Дерево граней имеет в точности $n + 1$ узлов, помеченных различными целыми числами от 0 до n . Отношения наследования (родители-сыновья) на таком дереве определяются просто: узел i является сыном узла-родителя $\beta[i]$. В упражнении 1.4.3 предложено доказать, что структура, получаемая на основе такого отношения наследования, действительно является деревом. В разделе 5.1 дан пример построения дерева граней.

Задача 2.3 (Вычисление дерева граней). Представить массив граней строковой последовательности x в виде корневого дерева (раздел 5.1). ■

¹Деревом называется связный неориентированный граф без циклов и без петель, имеющий один узел (называется корнем дерева), из которого выходят ребра к другим узлам. Дерево хотя и является неориентированным графом, но его узлы подчиняются отношениям наследования, и поэтому дерево имеет упорядоченную структуру. Конечные узлы дерева называются листьями. — *Примеч. ред.*

Три сформулированные выше задачи рассматриваются уже в этом разделе, поскольку они тесно связаны между собой. Но следующая задача значительно отличается от них и дает новый взгляд на внутреннюю структуру строковых последовательностей. Чтобы подойти к понятию дерева суффиксов, рассмотрим сначала общую структуру данных, которая называется *синтаксическое дерево* [136].²

Пусть $X = \{x_1, x_2, \dots, x_m\}$ — множество попарно различных строк. Тогда *синтаксическое дерево* на множестве X — это дерево поиска, имеющее в точности $m+1$ конечных узлов: по одному для каждой строки x_i ($i = 1, 2, \dots, m$) плюс один для пустой строки ε . Ребра синтаксического дерева помечены буквами, которые содержатся в строках из множества X , и специальным “сигнальным” символом $\$$, обозначающим конец строки. Строка x_i раскладывается на отдельные буквы — от корня дерева до конечного узла, которым оканчивается ребро, помеченное символом $\$$. В общем случае с *каждым* нисходящим путем от узла N_1 до узла N_2 , например, ассоциируется некая строка u . Чтобы избежать излишней сложности выражений и вместе с тем возможных терминологических осложнений, будем говорить, что узел N_2 — это строка u , и если N_1 является корнем дерева, тогда строка u состоит из букв, которыми помечен путь от узла N_1 до узла N_2 . Заметим, что в этом случае строка u является префиксом хотя бы одной какой-нибудь строки $x_i\$$. Если же $N_1 = N_2$, тогда $u = \varepsilon$. В частности, корень дерева является пустой строкой.

Использование символа $\$$ обусловлено тем, что $m+1$ конечных узлов являются листьями дерева и, как показано на рис. 2.1, *a* для множества $X = \{ab, abc\}$, без символа $\$$ невозможно было бы выделить узел 1, отображающий тот факт, что строка ab является префиксом строки abc . На этом рисунке можно заметить еще одно характерное свойство такой структуры данных — общие префиксы (такие, как ab или ε) на синтаксическом дереве отображаются только один раз. Еще раз подчеркнем, что каждый конечный узел соответствует строке, отличной от других.

На рис. 2.1, *b* показано *компактное* синтаксическое дерево (или синтаксическое дерево *Patricia*³) [185], которое получено из синтаксического дерева путем исключения узлов степени 2, которые имеют родителя и только одного сына, т.е. корень дерева в число исключаемых узлов не входит.⁴ Таким образом,

²В оригинале такая структура называется *trie*, это слово получено из букв, стоящих в середине слова *retrieval* (поиск, выборка, возврат). Устоявшегося термина для этой структуры в русской литературе пока нет. (О терминологическом разбросе можно судить по двум русским переводам книги [136], на которую ссылается автор.) На наш взгляд термин *синтаксическое дерево* наилучшим образом определяет это понятие. — *Примеч. ред.*

³Термин *Patricia trie* происходит от названия алгоритма Patricia — Practical Algorithm To Retrieve Information Coded In Alphanumeric, применяемый для построения специальных деревьев поиска и извлечения с его помощью информации на основе ключей [185]. — *Примеч. ред.*

⁴Степенью узла называется количество ребер, инцидентных с данным узлом, т.е. количество ребер, входящих и исходящих из данного узла. — *Примеч. ред.*

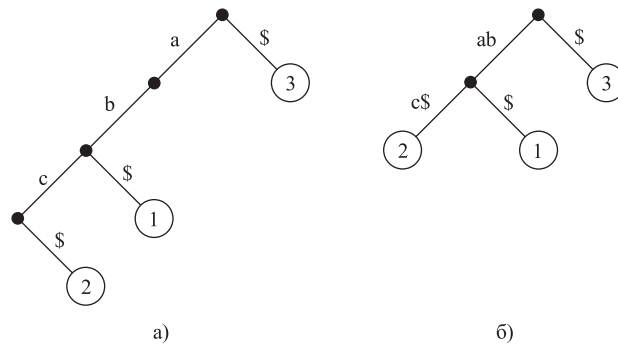


Рис. 2.1. Синтаксическое дерево для множества $X = \{ab, abc\}$

в компактном синтаксическом дереве будем иметь ребра, помеченные не отдельными буквами, а подстроками, как показано на рис. 2.1, б).

На каждом рис. 2.1 целыми числами помечены конечные узлы, соответствующие строкам $x_i\$$, $i \leq n$. “Обычное” дерево, соответствующее данной строке u , определяется, начиная от корня синтаксического дерева и далее через ребра, метки которых начинаются с очередной буквы строки u . Если возникнет одна из следующих ситуаций: из данного узла не выходят ребра, метки которых начинаются с очередной буквы строки u , либо достигнут неконечный узел, который не является префиксом u , или конечный узел, который не является строкой $u\$$, то во всех этих случаях $u \notin X$. Во всех других возможных случаях $u \in X$. В упражнении 2.1.3 предлагается доказать, что для любого узла T_x у всех ребер, выходящих из него, метки начинаются с различных букв.

Отметим, что синтаксическое дерево можно использовать для проверки не только того, что данная строка принадлежит множеству X , но и того, будет ли какой-нибудь ее префикс элементом этого множества. Однако следует помнить, что проверка префикса по компактному синтаксическому дереву несколько сложнее, чем по “обычному” синтаксическому дереву, где префикс u обязательно будет узлом дерева, тогда как в компактном дереве это не обязательно (например, префиксы a и abc на рис. 2.1, б). На компактном синтаксическом дереве конец префикса может быть *на ребре*, а не только в узле. Поэтому для решения подобных задач на компактном синтаксическом дереве (а также, как увидим в разделе 5.2.1, на дереве суффиксов) требуются дополнительные возможности для просмотра меток на ребрах, а не только для просмотра узлов.

Дерево суффиксов (suffix tree) строки x длиной n (иногда называемое **деревом подслов** (subword tree) [12]), если говорить кратко, — это компактное синтаксическое дерево для множества X всех суффиксов строки x (включая пустой суф-

фикс ε). Будем обозначать дерево суффиксов строки x как T_x . Отметим, что T_x имеет в точности $n + 1$ конечных узлов и, как показано в упражнении 2.1.4, не более n внутренних узлов, которые называются **узлами ветвления** (branch nodes). Таким образом, дерево суффиксов T_x имеет не более $2n + 1$ узлов и не более $2n$ ребер. Память, необходимую для хранения метки каждого узла, можно уменьшить, если вместо самой метки хранить два целых числа, указывающих позиции (начало и конец) этой метки-подстроки в строке x . (Но как указывалось в обсуждении 1.2, такая замена (подстроки на ее позиции в строке) предполагает, что получить доступ к подстроке на основе значений ее позиций можно за фиксированное время. Это, в свою очередь, предполагает, что исходная строка хранится в виде массива.) Таким образом, дерево суффиксов является “подходящей” структурой данных в том смысле, что для хранения требует память объемом порядка $\Theta(n)$ (см. также далее обсуждение 2.1). На рис. 2.2 показано дерево суффиксов T_g для строки g

$$g = a \overset{1}{b} \overset{2}{a} \overset{3}{a} \overset{4}{b} \overset{5}{a} \overset{6}{a} \overset{7}{b} \overset{8}{}$$

Отметим, что номера конечных узлов определяют позиции в исходной строке, где начинается данный суффикс. Узел с номером $n + 1$ соответствует пустому суффиксу.

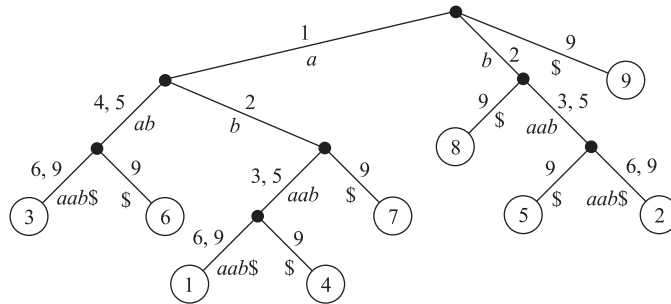


Рис. 2.2. Дерево суффиксов

Дерево суффиксов, введенное Винером [230], является очень важным внутренним паттерном. Для данной строки x дерево суффиксов T_x можно использовать для непосредственного определения того, является ли данная строка u суффиксом строки x , но, конечно, это можно сделать более эффективно напрямую в строке x . Поэтому более существенно то, что поскольку каждая подстрока строки x является префиксом некоторого суффикса этой же строки, дерево суффиксов можно использовать для определения того, является ли данная строка u подстрокой строки x . Более того, на основе дерева суффиксов можно определить позицию первого вхождения подстроки u в строку x , а также позиции *всех* возможных

вхождений подстроки u в строку x . Можно также применить дерево суффиксов для нахождения кратных строк в данной строке [18, 107, 216]. Поскольку дерево суффиксов можно сформировать для набора строк путем их сочленения, все вычислительные возможности дерева суффиксов распространяются на *множества* строковых последовательностей. Эффективные алгоритмы построения дерева суффиксов описаны в разделе 5.2.

В главе 11 мы покажем, что дерево суффиксов можно естественным образом интерпретировать как конечные автоматы. В этой интерпретации поиск по дереву суффиксов строки p эквивалентен доказательству выводимости выходного слова p конечным автоматом.

Задача 2.4 (Вычисление дерева суффиксов). Построить дерево суффиксов T_x для данной строки x (раздел 5.2).

Обсуждение 2.1.1. Сделаем замечание об эффективности построения дерева суффиксов и выполнения поиска по нему.

В разделе 5.2 мы изучим три алгоритма построения дерева суффиксов, два из которых требуют упорядоченности алфавита (раздел 4.1). Напомним, что в упорядоченном алфавите для каждой пары различных букв λ и μ за фиксированное время можно определить, выполняется ли неравенство $\lambda > \mu$. Эти два алгоритма выполняются за время порядка $O(n \log \alpha)$, где α — размер алфавита. Поскольку возможно, что $\alpha \in \Theta(n)$, то в самом худшем случае время выполнения алгоритма имеет порядок $\Omega(n \log n)$. Это нижняя граница временной сложности алгоритмов построения деревьев суффиксов для упорядоченных алфавитов.

Третий алгоритм из раздела 5.2 — это недавнее “открытие”; время выполнения этого алгоритма имеет порядок $\Theta(n)$ для индексированных алфавитов (раздел 4.1). Индексированный алфавит можно трактовать как множество целых чисел $\{1, 2, \dots, \alpha\}$ для некоторого $\alpha \in O(n)$. Индексированные алфавиты очень часто встречаются на практике, например двоичный алфавит можно рассматривать как индексированный, точно так же как любое подмножество символов ASCII (в том числе английский алфавит), алфавит ДНК и РНК и многие другие.

Фактически, если алфавит упорядочен, но не индексирован, в большинстве случаев построение дерева суффиксов выполняется эффективно. В разделе 5.2 представлен *онлайновый* алгоритм построения дерева суффиксов, когда дерево $T_{x\lambda}$ (для некоторой буквы λ) строится на основании дерева T_x за время порядка $O(\log n)$. Это означает, что для новых строк, которые получены путем добавления *справа* новых букв (т.е. путем добавления новых суффиксов), построение нового дерева суффиксов (на основе уже существующего дерева) требует только незначительных временных ресурсов. Подобным образом этот алгоритм можно эффективно использовать для построения дерева суффиксов при расширении набора строк.

Дополнительной к задаче построения дерева суффиксов является задача поиска по этому дереву. В этой задаче в каждом неконечном узле дерева T_x необходимо определить (основываясь на текущем значении буквы $u[i]$ строки u) нижележащие узлы, которые соответствуют этой строке. Это требует сравнения буквы $u[i]$ с буквой по крайней мере одного из нижележащих узлов, а возможно и всех.

Если алфавит упорядочен, то такое сравнение для каждой буквы $u[i]$ можно выполнить за время порядка $O(\log \alpha)$ путем использования подходящих структур данных (например, структуры дерева поиска или упорядоченного массива) для каждого узла. Для алфавитов малого размера эффективна реализация дерева T_x в виде бинарного дерева, как показано в упражнении 2.1.7. В общем случае для поиска строки u в строке x может потребоваться время порядка $\Omega(|u| \log \alpha)$. Если $|u|$ мало по сравнению с $n = |x|$ или размер алфавита также относительно небольшой, то для поиска всей строки x обычно требуется время порядка $O(n)$.

Поиск по дереву суффиксов остается наиболее эффективным в случае индексированных алфавитов, поскольку здесь оценочный фактор $\log \alpha$ можно проигнорировать на основании того, что при использовании процедур просмотра таблиц определить необходимый узел можно за конечное время. С другой стороны, если α велико, эффективный по времени алгоритм может потребовать большого объема памяти (порядка $\Theta(\alpha)$) для каждого узла дерева T_x . Когда алфавит упорядочен, иногда дает хороший результат применение техники хеш-таблиц, однако здесь снова для получения “почти постоянного” времени выполнения поиска приносится в жертву необходимый объем памяти.

В целом, дерево суффиксов — это чрезвычайно эффективное и широко используемое средство для поиска и сравнения строк. Его ахиллесовой пятой является необходимость использования больших объемов памяти, особенно в случае больших упорядоченных алфавитов (но и для малых алфавитов это может стать проблемой — см. обсуждение в подразделе 5.2.5). В общем случае применение дерева суффиксов для поиска и сравнения строк наиболее привлекательно, когда строка x неизменна, алфавит небольшой по размеру и фиксированный, а строка u относительно короткая. Этим условиям удовлетворяют, например, последовательности ДНК: алфавит фиксирован (C, G, A, T), как и x , при этом, как правило, $|u|$ значительно меньше n . ■

Кроме деревьев суффиксов, предложено большое количество разнообразных структур суффиксов. В разделе 5.3 мы рассмотрим две из них: ориентированные ациклические графы слов и массивы суффиксов.

Задача 2.5 (Вычисление структур суффиксов). Построить различные структуры суффиксов (раздел 5.3). ■

Внутренние паттерны представляют большой интерес для декомпозиции (факторизации) строковых последовательностей. В главе 6 рассмотрим линдонскую

декомпозицию (кратко описанную в разделе 1.4), где покажем ее связь с каноническими формами для петель. Решение задач декомпозиции основано на вычислении лексикографически наименьшего и наибольшего суффиксов для каждого префикса строки x . Эффективные алгоритмы линдонской декомпозиции описаны в разделах 6.1 и 6.2.

Задача 2.6 (Вычисление линдонской декомпозиции). Вычислить линдонскую декомпозицию строковой последовательности x (раздел 6.1). ■

Еще одной важной формой декомпозиции, весьма отличной от линдонской декомпозиции, является так называемая s -факторизация [156], которая впервые была применена как средство сжатия строковых последовательностей [235] и затем послужила основой для эффективных алгоритмов вычисления периодических составляющих строковых последовательностей. Если говорить кратко, то *s -факторизацией* строки x называется декомпозиция $x = w_1 w_2 \cdots w_k$, где каждая подстрока w_j ($j = 1, 2, \dots, k$) является или одиночной буквой, которая не встречается в подстроке $w_1 w_2 \cdots w_{j-1}$, или наибольшей подстрокой, которая одновременно будет суффиксом и собственной подстрокой строки $w_1 w_2 \cdots w_j$. Мы определим s -факторизацию более формально и приведем алгоритмы для ее вычисления в разделе 6.3, а ее применение покажем в разделе 12.2.

Задача 2.7 (Вычисление s -факторизации). Вычислить s -факторизацию данной строковой последовательности x (раздел 6.3). ■

Упражнения 2.1

1. Покажите, что “дерево граней” является ациклическим и связным деревом с корнем в узле 0.⁵
2. Приведите примеры компактного синтаксического дерева и дерева суффиксов, у которых из корневого узла исходит только одно ребро.
3. Покажите, что метки ребер, исходящих из какого-нибудь узла компактного синтаксического дерева, не имеют общего непустого префикса.
4. Пусть в данном дереве с k конечными узлами каждый внутренний узел имеет в точности два исходящих ребра. Покажите, что это дерево содержит ровно $k - 1$ внутренних узлов. На основании этого утверждения докажите, что соответствующее дерево суффиксов содержит не более n внутренних узлов.

⁵В принципе, если структура является деревом, то дерево ациклично (т.е. не имеет циклов) по определению. Здесь предлагается доказать, что структура “дерево граней” не имеет циклов. — *Примеч. ред.*

5. Обобщите предыдущий результат: докажите, что для компактного синтаксического дерева, построенного для множества из k строк, требуется хранить информацию не более чем о $O(k)$ узлах и ребрах. Приведите пример обычного синтаксического дерева, для которого это утверждение не выполняется.
6. Нарисуйте дерево суффиксов для строки $x = ababbabbaba$.
7. Предположим, что построено дерево суффиксов T_x для некоторой строки x на упорядоченном алфавите A , где $\alpha = |A|$ — известное фиксированное целое число. Для каждого из приведенных ниже алфавитов предложите структуру данных для узлов, которая позволяла бы эффективно находить правильное исходящее ребро при поиске по дереву. (Таких структур данных можно предложить несколько, в зависимости от того, считать ли α “маленьким” или “большим”.)
 - а) $A = \{1, 2, \dots, \alpha\}$.
 - б) Алфавит A содержит упорядоченные, но не обязательно последовательные буквы: $\lambda_1 < \lambda_2 < \dots < \lambda_\alpha$.

Разработайте собственный алгоритм построения дерева суффиксов для строк на фиксированном алфавите. Попробуйте разработать такой алгоритм, который был бы максимально эффективным, — время его выполнения не должно быть больше $O(n^2)$!

8. В упражнении 1.2.12 предлагалось разработать алгоритм для вычисления всех различных подстрок данной строки x за время порядка $O(n^2)$. Покажите, что если уже построено дерево суффиксов для этой строки x , тогда найти все различные подстроки этой строки можно за время порядка $\Theta(n)$.

Совет. Заметьте, что различные подстроки, которые начинаются в позиции i (т.е. строки вида $x[i..j]$, $i \leq j$), являются префиксами различных подстрок минимальной длины, которые также начинаются в позиции i . Таким образом, для нахождения искомым подстрок надо просто выписать все различные префиксы минимальной длины.

2.2 Частные паттерны

Паттерны, обсуждаемые в этом разделе, относятся к тому типу, которые вычисляются “классическими” алгоритмами обработки паттернов: имеется строка (иногда называемая *текстом*) и требуется найти одно или все вхождения в нее другой заданной строки (которая называется *паттерном*). Такую операцию, например, многократно (за один сеанс работы) выполняют текстовые редакторы. Другой пример этой задачи не такой очевидный: молекулярные биологи получают из международной базы данных генетической информации сегменты ДНК длиной

5 млн комплиментарных пар оснований нуклеиновых кислот (в нашей терминологии — это строка из 5 млн букв) для определения местоположения в этом большом сегменте отдельных коротких сегментов ДНК, состоящих из 300 комплиментарных пар оснований нуклеиновых кислот. Здесь опять применяются алгоритмы обработки паттернов.

Алгоритм 2.2.1

```

▷ Поиск всех вхождений строки  $p$  в строку  $x$ 
for  $i \leftarrow 1$  to  $n - m + 1$  do
   $j \leftarrow \text{сравнить}(i, m)$ 
  if  $j = m + 1$  then
    output  $i$ 

```

Удивительно простой алгоритм! Но с другой стороны, ни для кого не составит труда написать подобный алгоритм, поскольку он очевиден: для каждой позиции i в строке x выполняется тривиальная процедура $\text{сравнить}(i, m)$, которая сравнивает по одной букве слева направо строку $p[1..m]$ с подстрокой $x[i..i + m - 1]$. Эта процедура возвращает значение $m + 1$, если $p = x[i..i + m - 1]$, либо наименьшее целое $j \in 1..m$, такое, что $p[1..j] \neq x[i..i + j - 1]$. Таким образом, если $j = m + 1$, то это означает, что строка p входит в строку x начиная с позиции i . Все очень просто! Алгоритм 2.2.1 выполняется корректно, что доказывается в упражнении 2.2.2.

Недостатком такого алгоритма является его неэффективность. Например, в случае $x = a^n$, $p = a^{m-1}b$ процедура сравнить должна выполнить m побуквенных сравнений для каждого $i \in 1..n - m + 1$, т.е. всего $m(n - m + 1)$ побуквенных сравнений. Таким образом, алгоритм 2.2.1 требует $O(mn)$ времени выполнения, что чрезвычайно много, когда, например, $n = 5\,000\,000$ и $m = 300$ (как в исследованиях ДНК). Хотя алгоритм 2.2.1 в среднем не требует столь большого времени выполнения, все же стараются избегать его применения из-за временной затратности и иногда непредсказуемого поведения. Как мы уже видели, использование дерева суффиксов позволяет решить ту же задачу за время порядка $O(m \log \alpha)$, по крайней мере тогда, когда уже построено дерево суффиксов. Как будет показано в главах 7 и 8, существует много алгоритмов, которые гарантируют, что все вхождения строки p в строку x можно найти за время порядка $O(n + m)$. В главе 7 рассмотрены четыре “базовых” алгоритма обработки паттернов, а в главе 8 — еще несколько дюжин (не сотен) более сложных алгоритмов, которые являются оптимальными или в теории или на практике (но редко будут оптимальными одновременно и на практике и в теории).

Задача 2.8 (Вычисление всех паттернов). С гарантированной эффективностью вычислить все вхождения непустой строки (паттерна) p в заданную строку x (главы 7 и 8). ■

Во многих практических ситуациях сформулированные выше задачи поиска и сравнения строковых последовательностей не всегда адекватно отображают действительность. Например, при исследовании ДНК частичные паттерны часто повторяются в различных позициях, но эти повторяющиеся подстроки не совпадают абсолютно точно. Другими словами, в практических ситуациях представленный паттерн p может распознаваться на каких-либо подстроках исследуемой строки x , даже если между ними (между паттерном и сравниваемой подстрокой) нет полного совпадения. Например, если $p = CGAT$ и допустима перестановка двух первых или двух последних букв, тогда $p \approx GCAT$ и $p \approx CGTA$. Поэтому в строке

$$x = TCTAGG\underline{CGAT}TC\underline{GCAT}ATTC\underline{GCGT}AGCTCTA \quad (2.1)$$

подчеркнутые подстроки будут считаться совпадающими с паттерном p .

Основная идея в использовании аппроксимирующих (приближенных) паттернов состоит в определении “расстояния” между двумя строками. В общем случае *расстояние* между строками p_1 и p_2 рассчитывается как взвешенное количество стандартных операций редактирования, необходимых для выполнения преобразования $p_1 \rightarrow p_2$. Обычно рассматривают следующие операции редактирования.

- **Вставка** — например, вставка символа G в строку $GCAT$ сформирует строку $GCGAT$.
- **Удаление** — например, удаление символа G из строки $GCGAT$ сформирует строку $GCAT$.
- **Подстановка** — например, подстановка символа G вместо символа C из строки $GCAT$ сформирует строку $GGAT$.

Конечно, перечисленные операции определены только для непустых символов. На основе этих операций можно определить несколько типов расстояния между строками, некоторые из которых мы сейчас рассмотрим.

1. Если две строки x_1 и x_2 имеют одинаковую длину n , *расстояние Хемминга* $d_H(x_1, x_2)$ [109] определяется как минимальное количество *подстановок*, необходимых для преобразования строки x_1 в строку x_2 . Так, $d_H(GCAT, CGAT) = 2$. Отметим, что для замены некоего символа $x_1[i]$ на какой-нибудь другой может потребоваться одна операция удаления этого символа и одна операция вставки нового символа после символа $x_1[i - 1]$.
2. Для произвольных строк x_1 и x_2 *расстояние Левенштейна* $d_L(x_1, x_2)$ [158] определяется как минимальное количество операций *удаления* и *вставки*, необходимых для преобразования строки x_1 в строку x_2 . Так, $d_L(GCAT, CGAT) = 2$, но

$$d_L(GCGAT, CGAT) = d_L(CAT, CGAT) = 1.$$

Если строки x_1 и x_2 имеют одинаковую длину, то совсем не обязательно, что $d_H(x_1, x_2) = d_L(x_1, x_2)$. Например, если $x_1 = CGA$ и $x_2 = AGT$, то $d_H(CGA, AGT) = 2$, тогда как $d_L(CGA, AGT) = 4$. Отметим, что последовательность операций удаления и вставки, преобразующих строку x_1 в строку x_2 , может определяться неоднозначно, даже если общее количество таких операций одинаково. Например, строку $CGATA$ можно преобразовать в строку $ATACG$ путем применения четырех операций удаления символов CG и AT и последующих четырех вставок символов AT и CG либо путем применения двух операций удаления префикса CG и последующих двух вставок суффикса CG .

3. Для произвольных строк x_1 и x_2 **расстояние преобразования** (edit distance) $d_E(x_1, x_2)$ [158] определяется как минимальное количество операций *удаления*, *вставки* и *подстановки*, необходимых для преобразования строки x_1 в строку x_2 . Здесь предполагается, что подстановка одной буквы вместо другой выполняется за одну операцию, тогда как при вычислении расстояния d_L подстановка выполняется за две операции: удаление и последующая вставка. Чтобы лучше понять, как d_E соотносится с d_H и d_L , рассмотрим строки $x_1 = CGACG$ и $x_2 = GTCGA$. Для этих строк $d_H(x_1, x_2) = 5$, поскольку ни на одном месте в этих строках нет попарно совпадающих букв, и поэтому для получения из одной строки другой требуется ровно пять подстановок. Но $d_L(x_1, x_2) = 4$, так как необходимо удалить префикс C , вставить суффикс A и заменить букву A на букву T , для чего потребуется одна операция удаления и одна операция вставки. Однако $d_E(x_1, x_2) = 3$, поскольку здесь для замены буквы A на букву T требуется только одна операция.

Читатель может заметить, что в литературе термины “расстояние Левенштейна” и “расстояние преобразования” используются не последовательно: иногда расстояние преобразования называется расстоянием Левенштейна и *наоборот*. Я обещаю, что в книге эти термины будут применяться именно так, как они определены выше.

4. Предыдущие определения расстояния основаны на том, что все операции редактирования подсчитываются с одинаковыми значениями. Однако по крайней мере для операции подстановки одинаковые значения не всегда верно отражают условия решаемой задачи. Например, если некоторые подстановки (скажем, $C \rightarrow G$) выполняются чаще или с большей вероятностью, чем другие подстановки (например, $C \rightarrow T$). Чтобы отобразить такие различия, разным подстановкам назначаются разные веса, которые представляются в виде **матрицы весов** W . (Так, например, делается в молекулярной биологии.) Матрица весов имеет размерность $\alpha \times \alpha$ и состоит из неотрицательных действительных чисел, где число $W[i, j]$ является весом замены

(подстановки) j -й буквы i -й буквой алфавита A . Для произвольных строк x_1 и x_2 **взвешенное расстояние** $d_W(x_1, x_2)$ определяется как минимальная сумма количества операций удаления и вставки и весов операций подстановок, необходимых для преобразования строки x_1 в строку x_2 , при этом веса задаются матрицей весов W . Например, предположим, что $A = \{C, G, A, T\}$ и матрица W имеет вид

	C	G	A	T
C	0,0	0,4	0,9	1,0
G	0,4	0,0	1,0	0,9
A	0,9	1,0	0,0	0,5
T	1,0	0,9	0,5	0,0

Тогда, используя пример из определения 3, получим

$$d_W(CGACG, GTCGA) = 2,5,$$

поскольку теперь подстановка $A \rightarrow T$ “стоит” только 0,5. Более сложные варианты матрицы весов рассматриваются в разделе 13.3.

Обычно требуется, чтобы функция расстояния d , определенная на какой-либо области D , удовлетворяла условиям **метрики**: для любых $u, v, w \in D$

$$d(u, v) \geq 0, \quad (\text{условие неотрицательности}) \quad (2.2)$$

$$d(u, v) = 0 \Leftrightarrow u = v, \quad (\text{условие единственности}) \quad (2.3)$$

$$d(u, v) = d(v, u), \quad (\text{условие симметричности}) \quad (2.4)$$

$$d(u, w) \leq d(u, v) + d(v, w), \quad (\text{неравенство треугольника}) \quad (2.5)$$

В упражнении 2.2.12 предложено доказать, что расстояния d_H , d_L и d_E в действительности являются метриками, а в упражнении 2.2.13 — что расстояние d_W будет метрикой в том случае, если матрица W симметричная. Однако во многих практических приложениях матрица W не симметричная или же необходимо, чтобы функция расстояния учитывала не только веса операции подстановки, но и возможные веса операций вставки и удаления. Кроме того, пример строки (2.1) показывает, что все приведенные функции расстояния не рассматривают возможность операции **взаимообмена** соседних (или даже не соседних) букв в строке вместо или в дополнение операций подстановки, удаления и вставки. Учет такой операции в функции расстояния — это не просто “корректировка” определения расстояния между строками, от этого очень сильно зависит эффективность алгоритмов сравнения и поиска строк.

Первой и основной задачей, связанной с расстоянием между строками, конечно, является следующая задача.

Задача 2.9 (Вычисление расстояния). Для произвольных строк x_1 и x_2 на основании определения данной функции расстояния d вычислить $d(x_1, x_2)$ (глава 9). ■

В этом контексте представляет интерес еще одна проблема. Для строки x построим строку

$$x' = x[i_1]x[i_2] \cdots x[i_k],$$

где $1 \leq i_1 < i_2 < \cdots < i_k \leq n$ и $1 \leq k \leq n$. Строка x' называется **подпоследовательностью** строки x . Для полноты положим, что пустая строка ε является подпоследовательностью любой строки. Очевидно, что любая подстрока строки x является также ее подпоследовательностью. С другой стороны, строка $x' = cold$ будет подпоследовательностью, но не подстрокой строки $x = scrolled$. Имея две строки x_1 и x_2 , мы можем поставить задачу нахождения **наибольшей общей подпоследовательности** $LCS(x_1, x_2)$ ⁶. Тогда, например,

$$cold = LCS(scrolled, could).$$

Эта задача возникает в молекулярной биологии и тесно связана с задачей 2.9 и особенно с расстоянием Левенштейна. Фактически, как показано в упражнении 2.2.9, любой алгоритм, определяющий минимальную последовательность операций удаления и вставки для преобразования строки x_1 в строку x_2 , также эффективно находит $LCS(x_1, x_2)$. Для этого сначала надо выполнить в обеих строках необходимые операции удаления, и результирующая строка будет LCS . Таким образом, алгоритм нахождения наибольшей общей подпоследовательности можно построить на основе алгоритма вычисления расстояния Левенштейна. Однако, как показано в упражнении 2.2.15, $LCS(x_1, x_2)$ не обязательно определяется однозначно.

Задача 2.10 (Вычисление LCS). Вычислить наибольшую строку x' , которая будет подпоследовательностью двух заданных строк x_1 и x_2 (глава 9). ■

Хотя задача вычисления $LCS(x_1, x_2, \dots, x_k)$, $k > 2$, является NP-полной⁷ [130], для случая $k = 2$ и строк длиной n показано, что эту задачу можно решить за время порядка $\Omega(n^2)$, если сравнения выполняются на буквах некоего алфавита [4], и за время порядка $\Omega(n \log n)$, если алфавит упорядочен [116]. Различные типы алфавитов рассматриваются в разделе 4.1.

При решении задачи вычисления аппроксимирующих паттернов используется решение задачи 2.9 для нахождения максимального расстояния между паттернами.

⁶Здесь LCS (от англ. Longest Common Subsequence) переводится как “наибольшая общая подпоследовательность”. — Примеч. ред.

⁷NP-полной, если говорить кратко, называется такая задача, которая эквивалентна задаче полного перебора. — Примеч. ред.

Задача 2.11 (Вычисление всех аппроксимирующих паттернов). Для данных строки x , непустого паттерна p , функции расстояния d и числа $k \geq 0$ вычислить все подстроки x' строки x , такие, что $d(p, x') \leq k$ (глава 10). ■

Еще один важный способ реализовать идею приближенного равенства строк заключается в том, чтобы сгруппировать строки, соответствующие некоторому “паттерну”, основанному на определенных регулярных выражениях. Прежде чем описывать такие паттерны в полном объеме, рассмотрим применение в паттернах *символов замещения*, которые замещают в строке одну или несколько букв данного алфавита. Такие символы используются для поиска определенных слов в тексте на каком-нибудь “стандартном” языке (например, английском), находят применение в вычислительной технике (например, шаблоны (маски) имен файлов в системе DOS или функция `grep` в системе Unix). Символы замещения являются “метасимволами”, которые могут присутствовать в обычных выражениях. Существуют два общепринятых символа замещения:

символ \bullet замещает одну любую букву из алфавита A ;

символ $*$ замещает любую строку из множества A^* .

Например, в строке $x = aaabaabbaabbbba$ паттерну $p_1 = a \bullet a$ соответствуют подстроки aaa и aba , тогда как паттерну $p_2 = a * a$ — подстроки aa , aaa , aba , $abba$, $abbbba$ и многие другие, включая саму строку x . Данный пример показывает основную проблему, которая возникает при использовании паттернов с символами замещения, — это потеря свойства транзитивности при сравнении строк, т.е. строки, совпадающие с паттерном (например, с паттерном $a \bullet a$), могут не совпадать друг с другом. Это значительно затрудняет использование “обычных” алгоритмов сравнения строк с паттернами.

Задача 2.12 (Вычисление совпадений с паттернами, содержащими символы замещения). Вычислить все подстроки данной строки x , совпадающие с непустым паттерном p , который, возможно, содержит символы замещения \bullet и $*$ (раздел 10.4).

Вернемся к паттернам, являющимся регулярными выражениями. Такие паттерны сочетают в себе свойства аппроксимирующих паттернов и паттернов, по которым можно сделать множественный выбор. Как увидим далее, многие паттерны можно записать с помощью регулярных выражений. Построение регулярных выражений начнем с введения *метасимволов* — специальных символов, присутствующих в паттерне, которые *не являются* буквами алфавита A и имеют определенное значение. Неявно мы уже использовали метасимволы, например, когда строку $abababaabaabacccc$ записывали в более краткой форме, как $(ab)^2(aba)^3c^4$. И, конечно, метасимволами являются символы замещения, рассмотренные выше.

С помощью метасимволов можно представить целые классы строковых последовательностей, которые невозможно указать путем задания “простых” отдельных паттернов. Дадим определение нескольких метасимволов, которые значительно расширяют “представительские возможности” регулярных выражений.

М1. Для любой пары паттернов p_1 и p_2 запись $p_1 | p_2$ означает паттерн $\{p_1, p_2\}$, т.е. в сравнении участвуют и паттерн p_1 и паттерн p_2 . Например,

- паттерну $a^3 | b^3$ соответствуют и строка a^3 и строка b^3 ,
- паттерну $C | G | A$ соответствует любой элемент множества $\{C, G, A\}$,
- паттерну $\varepsilon | a | b | ab$ соответствует любой элемент множества $\{\varepsilon, a, b, ab\}$.

На конечном алфавите A метасимвол $|$ не только обобщает символ замещения \bullet , но и расширяет его возможности.

В операционной системе Unix и во многих языках формирования запросов к базам данных выражение $a | b$ обозначает запись $[ab]$. Если алфавит A упорядочен, то в дальнейшем мы, возможно, будем использовать выражение типа $[a-zA-Z]$, которое будет обозначать $a | b | \dots | z | A | B | \dots | Z$. Это только синтаксическое сокращение записи и никак не влияет на создание паттернов.

М2. Для любого паттерна p запись p^* означает нулевую или любую (конечную) конкатенацию (сочленение) строк p . Например,

- паттерну a^* соответствует любой элемент множества $\{a^i\}$, где i — произвольное неотрицательное целое число и считается, что $a^0 \equiv \varepsilon$;
- паттерну $(a | b)^*$ соответствует любой элемент множества A^* , где $A = \{a, b\}$;
- паттерну $(a | b)^*(a | b)(a | b)^*$ соответствует любой элемент множества A^+ , где $A = \{a, b\}$;
- паттерну $a^* | b^*$ соответствует любой элемент множества $\{a^i, b^i\}$, где i — произвольное неотрицательное целое число;
- паттерну a^*b^* соответствует любая строка вида a^ib^j , где i и j — неотрицательные целые числа;
- паттерну $(ab)^*$ соответствует любая строка вида $(ab)^i$, где i — неотрицательное целое число.

Из этих примеров видно, что комбинация метасимволов $|$ и $*$ может соответствовать символу замещения $*$. Но с помощью метасимвола $*$ можно представить больше паттернов, чем с помощью любой комбинации символов замещения, например паттерны a^* и $(ab)^*$ нельзя записать с помощью символов замещения.

М3. Для любого паттерна p запись $\sim p$ означает любую строку, не соответствующую паттерну p . Например, на алфавите $A = \{a, b\}$

- паттерну $\sim(a^*)$ соответствует паттерн $(a | b)^*b(a | b)^*$, т.е. любая строка, содержащая букву b ;
- паттерну $\sim(a | b)^*$ не соответствует ни одна строка, даже пустая;
- паттерну $\sim((a | b)^*(a | b)(a | b)^*)$ соответствует только пустая строка ε ;
- паттерну $\sim(a^*b^*)$ соответствует любая строка, в которой буквы b предшествуют буквам a , т.е. этот паттерн эквивалентен паттерну $a^*bb^*a(a | b)^*$;
- паттерну $\sim(a^* | b^*)$ соответствует любая строка, содержащая как букву a , так и букву b , т.е. этот паттерн эквивалентен паттерну

$$((a | b)^*a(a | b)^*b(a | b)^*) | ((a | b)^*b(a | b)^*a(a | b)^*) \quad (2.6)$$

- паттерну $\sim(ab)^*$ соответствует любая непустая строка, в которой нет степеней строки ab , т.е. этот паттерн эквивалентен паттерну $(ab)^*(a | (aa | b)(a | b)^*)$.

Приведенные примеры показывают, что любые паттерны, которые включают метасимвол \sim , можно переписать без его использования. Такая ситуация имеет место и в общем случае — не существует классов строковых последовательностей, которые нельзя было бы описать с помощью лишь метасимволов $|$ и $*$. Таким образом, метасимвол \sim является только “метасимволом для удобства”, поскольку он позволяет сократить запись многих регулярных выражений.

Сделаем замечание: в приведенных выше выражениях с метасимволами была допущена небольшая вольность — символ p использовался как для обозначения отдельных строк, так и для паттерна, который представляет *множество* строк. Я надеюсь, что такая вольность не вызовет проблем у читателя. В оправдание могу сказать, что это сделано для сокращения текста и во избежание излишнего формализма.

Теперь можно сформулировать следующее определение.

Определение 2.2.1. Регулярным выражением называется паттерн, который может содержать метасимволы $|$ и $*$. **Регулярным множеством** называется множество строковых последовательностей, представленных регулярным выражением. ■

В связи с этим определением возникает следующая задача.

Задача 2.13 (Вычисление всех совпадений с регулярным выражением). Для данной строки x вычислить все ее подстроки, которые соответствуют данному регулярному выражению p (раздел 11.1). ■

Использование метасимволов позволяет реализовать идею сравнения строк не с одним паттерном, а с множеством паттернов. В разделе 11.2 соответствующая

задача будет представлена в двух формах: точной и приближенной. Задача нахождения точного совпадения множества паттернов $P = \{p_1, p_2, \dots, p_r\}$ с подстроками строки $x = x[1..n]$ фактически является частным случаем задачи сравнения с регулярным выражением $p = p_1 | p_2 | \dots | p_r$. Такое сравнение можно выполнить за время порядка $\Theta(m + n)$, где $m = |p_1| + |p_2| + \dots + |p_r|$. В разделе 11.2 также описаны два алгоритма для решения более сложной задачи нахождения приближенного совпадения подстрок строки x с множеством паттернов.

Задача 2.14 (Вычисление множественных паттернов). Вычислить все (точные или приближенные) совпадения всех паттернов $P = \{p_1, p_2, \dots, p_r\}$ с подстроками заданной строки x . ■

Упражнения 2.2

1. Напишите алгоритм “тривиальной” процедуры *сравнить* и докажите ее корректность.
2. Докажите, что алгоритм 2.2.1 корректно выполняется для любого паттерна p . Удостоверьтесь, что вы рассмотрели все возможные случаи, включая случаи, когда $n = 0$, $m = 0$, $m = n$, $m > n$. Имеет ли смысл рассматривать случаи пустой строки и пустого паттерна?
3. Предположим, что алгоритм 2.2.1 применяется к строке x длиной n и паттернам p фиксированной длины m , которые генерируются “случайным образом” на алфавите размером $\alpha \geq 2$. Покажите, что ожидаемое количество сравнений можно вычислить по приближенной формуле

$$(n - m + 1) \left(\frac{\alpha}{\alpha - 1} \right) \left(1 - \left(\frac{1}{\alpha} \right)^m \right).$$

4. Для строк x_1 и x_2 длиной n покажите, что
 - а) $d_L(x_1, x_2) \leq 2d_H(x_1, x_2) \leq 2n$;
 - б) $d_H(x_1, x_2) - d_L(x_1, x_2) \leq n - 2$, $n \geq 2$;
 - в) $d_E(x_1, x_2) \leq \min\{d_H(x_1, x_2), d_L(x_1, x_2)\}$.

Приведите примеры строк, для которых приведенные неравенства будут выполняться в виде равенств.

5. Вычислите следующие расстояния для функций расстояния d_H , d_L и d_E соответственно.
 - а) $d(x, \varepsilon)$;
 - б) $d((ab)^3, (ba)^3)$;
 - в) $d(ababbba, aabbaab)$.

6. При обсуждении расстояния d_L указывалось, что $d_L(CGACG, GTCGA) = 4$, поскольку для преобразования $CGACG \rightarrow GTCGA$ необходимы две операции удаления и две операции вставки. Предложите другие четыре операции удаления и вставки для выполнения этого преобразования.
7. Для произвольной строки $x[1..n]$ покажите, что для любого $j \in 0..n - 1$ $d_H(x, R_j(x)) \neq 1$, где $R_j(x)$ — j -й циклический сдвиг строки x (см. раздел 1.4).
8. Для расстояний преобразования и Левенштейна найдите такие непустые строки x_1, x_2 и x_3 , что $|x_1| = |x_2|$ и $d(x_1, x_2) > d(x_1, x_2x_3)$.
9. Пусть имеется три строки x_1, x_2 и x_3 , такие, что x_1 является подпоследовательностью строки x_2x_3 , и выполняется неравенство

$$|x_2x_3| - |x_1| < d(x_1, x_2),$$

где d — расстояние преобразования или Левенштейна. Вейлин Лу (Weilin Lu) предположил, что в этом случае должно выполняться неравенство $d(x_1, x_2) > d(x_1, x_2x_3)$. Прав ли Вейлин Лу?

10. Определим **расстояние удаления** (deletion distance) $d_D(x_1, x_2)$ между двумя строками x_1 и x_2 как минимальное количество операций удаления, необходимых для преобразования строк x_1 и x_2 в одинаковые строки. (Например, $d_D(ab, bac) = 3$, поскольку необходимо три операции удаления для преобразования строк ab и bac в строку a либо столько же операций удаления для преобразования в строку b .) Докажите, что

$$d_D(x_1, x_2) = |x_1| + |x_2| - 2|\text{LCS}(x_1, x_2)| = d_L(x_1, x_2).$$

11. Определим **асимметричное расстояние преобразования** $d_A(x_1, x_2)$ между двумя строками x_1 и x_2 как минимальное количество операций удаления в строке x_2 плюс минимальное количество операций подстановки/удаления в строке x_1 , необходимых для преобразования строк x_1 и x_2 в одинаковые строки. Покажите, что $d_A(x_1, x_2) = d_E(x_1, x_2)$.
12. Покажите, что функции расстояния d_H, d_L и d_E являются метриками в произвольной области A^* , т.е. для них выполняются условия (2.2)–(2.5).
13. Покажите, что функция расстояния d_W будет метрикой тогда и только тогда, когда матрица W будет симметричной с нулевой диагональю и положительными внедиагональными элементами.
14. Покажите, что расстояние Левенштейна не зависит от порядка, в котором выполняются операции удаления и вставки. Интерпретируя операцию подстановки как последовательность операций удаления и вставки, докажите аналогичное утверждение для расстояния преобразования и взвешенного

расстояния. Какие метрические условия используются при доказательстве этих утверждений?

15. Пусть строки x_1 и x_2 такие, что ни одна из них не является подпоследовательностью другой. Покажите, что в этом случае $x = \text{LCS}(x_1, x_2)$ только тогда, когда $d_L(x_1, x) + d_L(x_2, x) = d_L(x_1, x_2)$.
16. Покажите на примерах, что результат операции $\text{LCS}(x_1, x_2)$ не единственный. Далее покажите, что для произвольного положительного целого числа k можно найти две строки x_1 и x_2 длиной $2k$, такие, что для них различные наибольшие общие подпоследовательности будут различаться в k позициях.
17. Для произвольных двух строк x_1 и x_2 подобно $\text{LCS}(x_1, x_2)$ можно также вычислить *наибольший общий фактор* (подстроку), который обозначается как $\text{LCF}(x_1, x_2)$. Разработайте алгоритм вычисления $\text{LCF}(x_1, x_2)$ на основе деревьев суффиксов и найдите его алгоритмическую сложность.
18. На алфавите $A = \{a, b\}$ охарактеризуйте строки, которые соответствуют паттерну $p = a \bullet b * a$.
19. Упростите следующие выражения для паттернов.
 - а) $(a | b)^*(a | b)(a | b)^*$;
 - б) $(a | b)^*b(a | b)^*$;
 - в) $(a^* | b^*)^*$.
20. Покажите, что выражение (2.6) эквивалентно выражению $((a^*ab) | (b^*ba))(a | b)^*$. Можно ли еще упростить это выражение?

2.3 Характеристические паттерны

В этом разделе мы рассмотрим задачи, которые требуют вычисления в строках характеристических паттернов. Эти паттерны отображают внутреннюю структуру строк и не выражаются в виде набора определенных подстрок или классов подстрок, которые также со своей стороны характеризуют строковые последовательности. В большинстве случаев структура строк основана на идее периодичности, рассмотренной в разделе 1.2. Грубо говоря, характеристические паттерны в некотором смысле отображают “приближенную периодичность”.

Начнем с паттернов, которые действительно отображают *точную* периодичность строк — наличие в строке кратных подстрок, определенных в разделе 1.2. Кратные строки привлекли внимание математиков еще в начале 20 столетия, когда Аксель Туе (Axel Thue) [220] рассмотрел конструкцию бесконечных некрatных строковых последовательностей на алфавите из трех букв. В разделах 3.2 и 3.3 мы рассмотрим некоторые “строки Туе”. Проблема возродилась в компьютерных

науках как задача вычисления всех кратных подстрок данной строки, которую поставили в 70-х годах (по-видимому, первыми) Мейн и Лоренц (Main and Lorentz) [169]. Эта задача и ее расширения нашли применение в молекулярной биологии и в алгоритмах сжатия данных.

В алгоритмах решения подобных задач критическим вопросом является определение максимального количества возможных кратных подстрок, находящихся в данной строке. Это количество определяет объем выходных результатов выполнения алгоритмов и устанавливает нижнюю границу алгоритмической сложности алгоритмов. Например, рассмотрим строку $x = aaaaaa$. Эта строка содержит пять вхождений квадрата a^2 , три вхождения квадрата $(aa)^2$ и одно вхождение квадрата $(aaa)^2$. Таким образом, данная строка имеет девять различных вхождений квадратов. В общем случае нетрудно показать, что строка a^n содержит $\lfloor n^2/4 \rfloor$ вхождений различных квадратов. Отсюда следует естественное заключение, что для вычисления всех квадратов строки длиной n необходимо время порядка $\Omega(n^2)$.

Но это “естественное” заключение ошибочно! Например, если мы действительно захотим “сделать перепись” всех квадратов строки a^6 , то без труда заметим, что их легко получить на основании того факта, что подстрока $x[1] = a$ повторяется шесть раз. Эту информацию можно записать с помощью “тройки Крочемора” [69]: $(i, p^*, r^*) = (1, 1, 6)$. Эта запись обозначает, что подстрока длиной $p^* = 1$, которая начинается в позиции $i = 1$, повторяется $r^* = 6$ раз. Таким образом, здесь p^* является периодом, а r^* — степенью строки x , как они определены в разделе 1.2. Другими словами, тройка (i, p^*, r^*) предлагает разложение в виде нормальной формы любой максимальной кратной подстроки, присутствующей в данной строке x . Например, строка $x = a^2b^3a^2b^3a$ имеет кратные подстроки

$$(1, 1, 2), (3, 1, 3), (6, 1, 2), (8, 1, 3), (1, 5, 2), (2, 5, 2).$$

Поскольку нормальная форма полностью описывает любую строку, для решения задачи нахождения всех кратных подстрок необходимо найти нормальные формы всех максимальных кратных подстрок, содержащихся в данной строке x . Однако напомним (см. упражнение 1.2.19), что в нормальной форме u^{r^*} любой строки образующая строка u не может быть кратной строкой.

Можно сказать, что тройка (i, p^*, r^*) — это способ **кодирования** кратных подстрок. В терминах такого кодирования в разделе 3.4 мы покажем, что строки Фибоначчи f_n содержат $\Theta(f_n \log f_n)$ максимальных кратных подстрок (здесь $f_n \equiv \equiv |f_n|$), каждая из которых кодируется одной тройкой (i, p^*, r^*) . Отсюда следует, что $\Omega(n \log n)$ является нижней границей сложности любого алгоритма поиска всех кратных подстрок данной строки x длиной n . В главе 12 мы покажем, что эта нижняя граница достигается на алгоритмах, решающих следующую задачу.

Задача 2.15 (Вычисление кратных подстрок). Вычислить все кратные подстроки данной строки x длиной n за время $O(n \log n)$ (раздел 12.1).

Обсуждение 2.3.1. Кодирование кратных подстрок, описанное выше, неявно предполагает, что имеется массив граней (см. обсуждение 1.3.2). Это позволяет нам (по соглашению, сделанному в обсуждении 1.3.2) хранить больше информации в памяти меньшего объема — с помощью массива граней в памяти объемом $\Theta(n)$ можно хранить информацию объемом $\Theta(n \log n)$. Поэтому информацию о кратных подстроках, общее количество которых имеет порядок $\Theta(n^2)$, можно с использованием нормальной формы представить (другими словами, сжать) в виде описанных выше троек, которых всего $\Theta(n \log n)$.

В случае массива граней мы не можем надеяться, что удастся уменьшить требуемый объем памяти до величины, меньшей $\Theta(n)$. Однако в случае кратных подстрок встает законный вопрос: можно ли найти такой способ кодирования всех кратных подстрок во “что-то такое”, что потребует объема памяти меньшего $\Theta(n \log n)$, например объема $\Theta(n \log \log n)$ или даже $\Theta(n)$? Выше было сказано, что нижняя граница времени вычисления всех кратных подстрок составляет $\Omega(n \log n)$, но эта величина не связана с требуемым объемом памяти. С другой стороны, очевидна тривиальная оценка $\Omega(n)$ необходимого объема памяти для записи всех кратных подстрок. В разделе 3.4 мы покажем, что в случае строк Фибоначчи f_n для описания всех кратных подстрок достаточно памяти порядка $\Theta(f_n)$. Отсюда вытекают следующие вопросы: можно ли добиться линейного порядка необходимого объема памяти для других типов строк, и если это возможно, то почему невозможен линейный порядок времени вычислений?

Мы рассмотрим эти интересные вопросы в главе 12, а данное обсуждение поможет найти правильные ответы. ■

Задача 2.15 — это “классическая” задача, связанная с повторяемостью подстрок; она рассматривалась во многих работах вплоть до 90-х годов прошлого столетия. И по мере исследования этой задачи идея кратных подстрок изменялась и получала развитие в самых разных направлениях. Например, большой интерес представляют несмежные повторяемые подстроки. Они нашли применение в анализе последовательностей ДНК. Другой пример: хранение, извлечение и анализ музыкальных текстов [58], где необходимо распознавать повторяемые мотивы или мелодии в отдельных или последовательных музыкальных фрагментах.

Обсуждение “расширения” понятия кратных подстрок начнем со следующего определения.

Определение 2.3.1. Раппортом (*repeat*)⁸ строки x называется кортеж

$$M_{x,u} = (p; i_1, i_2, \dots, i_r), r \geq 2,$$

⁸Слово *repeat* имеет множество значений, и более близким в данном случае был бы его перевод как “повторение”. Однако слово “повторение” в разных вариациях будет многократно использоваться в различных ситуациях. Поэтому считаем рациональным использовать термин *раппорт*, который также является переводом слова *repeat*. — *Примеч. ред.*

где $i_1 < i_2 < \dots < i_r$ и

$$\mathbf{u} = \mathbf{x}[i_1..i_1 + p - 1] = \mathbf{x}[i_2..i_2 + p - 1] = \dots = \mathbf{x}[i_r..i_r + p - 1].$$

Подстрока \mathbf{u} называется **повторяемой подстрокой** строки \mathbf{x} и **производящей строкой** для $M_{\mathbf{x},\mathbf{u}}$. Как и в случае кратных подстрок, назовем $p = |\mathbf{u}|$ **периодом** повторения, а r — **показателем**. Если кортеж охватывает все вхождения \mathbf{u} в строку \mathbf{x} , тогда раппорт $M_{\mathbf{x},\mathbf{u}}$ называется **полным** и обозначается как $M_{\mathbf{x},\mathbf{u}}^*$.⁹ ■

С учетом этого определения нашу задачу можно сформулировать как вычисление раппортов $M_{\mathbf{x},\mathbf{u}}^*$ для каждой повторяемой подстроки \mathbf{u} . Как мы увидим далее, задачу можно упростить путем введения понятия “продолжаемости”. Но сначала рассмотрим связь нового понятия раппорта с уже знакомым понятием кратных строк.

Очевидно, что раппорт является обобщением кратных строк с заимствованием соответствующих сопровождающих понятий производящей строки, периода и показателя. Чтобы показать связь между ними более четко, введем понятие **лакуны** как разности

$$g_j = i_{j+1} - i_j, \quad 1 \leq j \leq r - 1,$$

между последовательными элементами раппорта $M_{\mathbf{x},\mathbf{u}}$. Теперь можно классифицировать раппорты в соответствии со значениями лакун.

- Раппорт $M_{\mathbf{x},\mathbf{u}}$ назовем кратной строкой (или **последовательным раппортом** (tandem repeat)), если для всех лакун $g_j = p$.¹⁰ В этом случае $M_{\mathbf{x},\mathbf{u}}$ можно записать в сокращенной форме

$$(i, p, r) \equiv (p; i, i + p, \dots, i + (r - 1)p).$$

- Раппорт $M_{\mathbf{x},\mathbf{u}}$ назовем **расщеплением** (split), если для всех лакун $g_j > p$.¹¹
- Раппорт $M_{\mathbf{x},\mathbf{u}}$ назовем **покрытием** (overlap), если для всех лакун $g_j < p$.
- Раппорт $M_{\mathbf{x},\mathbf{u}}$ назовем **оболочкой** (cover), если для всех лакун $g_j \leq p$. В этом случае будем говорить, что $M_{\mathbf{x},\mathbf{u}}$ является оболочкой строки $\mathbf{x}[i_1..i_r + p - 1]$, а эта строка **имеет оболочку** или **\mathbf{u} -оболочку**.
- Если раппорт $M_{\mathbf{x},\mathbf{u}}$ не попадает ни в одну из перечисленных категорий, назовем его **смешанным**.

⁹Если отвлекаться от технических деталей и формы записи, то раппорт можно определить как кратную строку \mathbf{u}^r , где \mathbf{u} — подстрока строки \mathbf{x} . Сама строка-раппорт может быть подпоследовательностью или подстрокой строки \mathbf{x} либо не быть ни тем, ни другим (см. классификацию раппортов ниже). — *Примеч. ред.*

¹⁰В этом случае раппорт является подстрокой строки \mathbf{x} . — *Примеч. ред.*

¹¹В этом случае раппорт является подпоследовательностью строки \mathbf{x} . — *Примеч. ред.*

Чтобы пояснить приведенную классификацию раппортов, рассмотрим пример строки

$$x = \overset{1}{a} \overset{2}{b} \overset{3}{a} \overset{4}{a} \overset{5}{b} \overset{6}{a} \overset{7}{b} \overset{8}{a} \overset{9}{a} \overset{10}{b} \overset{11}{a} \overset{12}{a} \overset{13}{b} \overset{14}{a} \overset{15}{b} \overset{16}{a} \overset{17}{a} \overset{18}{b} \overset{19}{a} \overset{20}{b} \overset{21}{a}$$

Пусть $u = aba$. Можно определить следующие раппорты.

- Раппорты (3; 1, 4) и (3; 6, 9, 12) являются кратными строками.
- Раппорт (3; 1, 6, 12, 17) является расщеплением.
- Раппорт (3; 4, 6) является покрытием.
- Раппорт (3; 1, 4, 6, 12) является смешанным.
- Полный раппорт $M_{x,u}^* = (3; 1, 4, 6, 9, 12, 14, 17, 19)$ является оболочкой строки x .

Заметим, что из полного раппорта $M_{x,u}^*$ можно эффективно вычислить все повторения подстроки u , все покрытия и все подстроки максимальной длины, имеющие u -оболочки, — вся эта информация содержится в полном раппорте. Так в нашем примере просматривая полный раппорт, определяем, что

- раппорты (3; 1, 4), (3; 6, 9, 12) и (3; 14, 17) являются повторениями подстроки $u = aba$ в строке x ;
- раппорты (3; 4, 6), (3; 12, 14) и (3; 17, 19) являются покрытиями;
- существует только одна подстрока максимальной длины, имеющая оболочку, и эта подстрока — сама строка x .

В нашем примере рассмотрим также полный раппорт для $u = b$:

$$M_{x,b}^* = (1; 2, 5, 7, 10, 13, 15, 18, 20).$$

Просмотр этого раппорта позволяет обнаружить интересный факт: за каждой буквой b следует буква a .¹² Из этого наблюдения сразу получаем раппорт

$$M_{x,ba}^* = (2; 2, 5, 7, 10, 13, 15, 18, 20),$$

который в точности совпадает с раппортом $M_{x,b}^*$! Из последнего раппорта получаем

$$x[2..3] = x[5..6] = \dots = x[20..21] \quad \text{и} \quad x[2] = x[5] = \dots = x[20].$$

Аналогично просмотр раппорта $M_{x,ba}^*$ позволяет сделать заключение, что каждой подстроке ba предшествует буква a . Отсюда получаем раппорт

$$M_{x,aba}^* = (3; 1, 4, 6, 9, 12, 14, 17, 19).$$

¹²Поскольку все лакуны имеют значения, не меньшие 2. — Примеч. ред.

Описанный способ последовательного вычисления раппортов подводит нас к способу кодирования раппортов (подобно кодированию граней и кратных подстрок, см. обсуждения 1.3.2 и 2.3.1), который может уменьшить необходимый объем памяти и время вычислений. Чтобы точно описать этот способ кодирования, введем еще несколько определений.

Определение 2.3.2. Раппорт $M_{x,u} = (p; i_1, i_2, \dots, i_r)$ называется **продолжаемым влево** (left-extendible, LE) (соответственно, **продолжаемым вправо** (right-extendible, RE)), если

$$(p; i_1 - 1, i_2 - 1, \dots, i_r - 1) \quad (\text{соответственно, } (p; i_1 + 1, i_2 + 1, \dots, i_r + 1))$$

также является раппортом. Если раппорт $M_{x,u}$ не продолжаем ни влево, ни вправо, то он называется **непродолжаемым** (NE). Для краткости такие раппорты будем обозначать **LE**, **RE** и **NE** соответственно. ■

Интересно рассмотреть особый случай определения 2.3.2, когда раппорт $M_{x,u}$ является кратной строкой $(u[1..p])^r$, $r \geq 2$. Например, пусть $x = \dots a(aba)^r a \dots$ и порождающая строка $u = aba$ с периодом $p = 3$. Очевидно, что кратная строка $(aba)^r$ является раппортом типа LE (продолжаемым влево), поскольку $(aab)^r$ является раппортом с порождающей строкой aab с периодом $p = 3$. Из подобных соображений следует, что строка $(aba)^r$ является также раппортом типа RE (продолжаемым вправо). Этот пример показывает, что в случае кратной строки тот факт, что раппорт $M_{x,u}$ продолжаемый влево, говорит о том, что подстроке u^r в строке x предшествует буква $\lambda = u[p]$ и поэтому раппорт $M_{x,u}$ можно преобразовать в раппорт $M_{x,u'}$ с тем же периодом p , где $u' = \lambda u[1..p-1]$. Аналогично, если раппорт $M_{x,u}$ продолжаемый вправо, тогда он преобразуется в раппорт $M_{x,u''}$ с тем же периодом p , где $u'' = u[2..p]u[1]$. В разделе 1.4 подстроки u' и u'' определялись нами как *циклический сдвиг* подстроки u . В терминах троек Крочемора сказанное можно записать так: если кратная строка (i, p^*, r^*) является раппортом типа LE, тогда строка $(i-1, p^*, r^*)$ также будет кратной; если же кратная строка (i, p^*, r^*) является раппортом типа RE, тогда кратной будет строка $(i+1, p^*, r^*)$. Эти рассуждения подводят нас еще к одному важному понятию.

Определение 2.3.3. Для данной строковой последовательности x 4-местный кортеж (i, p^*, r^*, t) назовем **серией** (run) строки x , если выполняются следующие условия:

- а) $t \in 0..p^* - 1$;
- б) раппорт (i, p^*, r^*) является кратной строкой, не продолжаемой вправо;
- в) раппорт $(i+t, p^*, r^*)$ является кратной строкой, не продолжаемой вправо.

Константа t называется **хвостом** (tail) серии. ■

Понятие серии впервые, по-видимому, было введено Мейном (Main) в работе [168], где для этого понятия использовался менее краткий, но более выразительный термин *максимальная периодичность* (maximal periodicity). Для этого же понятия в работе [139] используется термин *максимальная кратная строка* (maximal repetition).

В обсуждении 2.3.1 упоминалось, что кратные подстроки в строках Фибоначчи можно вычислить за линейное время; в разделе 3.4 мы также увидим, что за такое же время для этих строк можно вычислить серии. Случай произвольных строк исследован в разделе 12.2: там показано, что количество серий в произвольной строке фактически является линейной функцией от длины строки. Этот результат открывает возможность вычисления всех серий по крайней мере для некоторых строк за время, не превышающее $O(n \log n)$. Таким образом, имеем еще одну задачу.

Задача 2.16 (Вычисление серий). Эффективно вычислить все серии для данной строки x (раздел 12.2). ■

Введенное выше понятие серии подводит к формулировкам задач, связанных с вычислением оболочек и раппортов. Обобщение кратных строк путем введения раппортов позволяет рассмотреть не только расщепленные кратные строки, но и подстроки максимальной длины, имеющие оболочки. Таким способом можно обобщить задачу 2.15 путем замены кратных строк (когда лакуны равны p) на подстроки, имеющие оболочки (когда лакуны не превышают p).

Задача 2.17 (Вычисление подстрок, имеющих оболочки). Вычислить все подстроки данной строки x , которые имеют непродолжаемые оболочки. ■

Известны три алгоритма [17, 42, 124], которые решают задачу, очень похожую на задачу 2.17, — задачу вычисления всех подстрок заданной строки, имеющих непродолжаемые вправо оболочки. К сожалению, эти алгоритмы очень сложные как в теоретическом плане, так и в практическом. Поэтому мы не будем исследовать их во всех деталях, но в разделе 13.2, где показано решение близкой задачи вычисления раппортов (задача 2.19), кратко их рассмотрим и обсудим возможное их усовершенствование.

Идея “оболочки” порождает интересную задачу, которая находит применение в сжатии данных: вычислить все оболочки данной строки. Оболочку строки можно рассматривать как обобщение производящей подстроки для кратных строк. С этой точки зрения оболочку строки иногда называют *квазипериодом*, а строку, имеющую оболочку, — *квазипериодической*.

Несколько алгоритмов для вычисления квазипериодов предложено в 90-х годах. В главе 13 мы рассмотрим один из них [160], который вычисляет внутренний паттерн, называемый “массивом оболочек”, — аналог массива граней, описанно-

го в разделе 1.3. Подобно массиву граней, массив оболочек можно вычислить за линейное время. И так же, как массив граней определяет все периоды всех префиксов строки x , так и массив оболочек определяет все квазипериоды всех префиксов строки x .

Задача 2.18 (Вычисление оболочек). Вычислить все оболочки заданной строки x (раздел 13.1). ■

В связи с этой задачей очевидна необходимость вычисления всех раппортов типа NE (т.е. непродолжаемых). Это приводит к следующему обобщению задачи 2.15.

Задача 2.19 (Вычисление раппортов). Вычислить все непродолжаемые раппорты в заданной строке x (раздел 13.2). ■

Как будет показано в разделе 13.2, эту общую задачу можно решить за время порядка $\Theta(n \log n)$, и ее решение в свою очередь может послужить основой для решения задач 2.15 и 2.17. Алгоритмы, решающие данную задачу, находят применение в вычислительной биологии.

На практике также находят применение алгоритмы, вычисляющие аппроксимирующие раппорты подстрок, отличающиеся не более, чем на расстояние k , где “расстояние” может принимать одну из форм, описанных в разделе 2.2.

Задача 2.20 (Вычисление аппроксимирующих раппортов). Для данной функции расстояния d и целого числа $k \geq 0$ в заданной строке x вычислить все аппроксимирующие раппорты, отличающиеся не более, чем на расстояние k (раздел 13.3). ■

Несмотря на свою значимость, задача вычисления аппроксимирующих раппортов изучена относительно слабо, особенно в сравнении с задачей вычисления аппроксимирующих паттернов (раздел 2.2). Причиной является исключительная сложность данной задачи — даже лучшие алгоритмы имеют время выполнения порядка $O(n^2)$. Трудность задачи, как показано ниже, проявляется уже на этапе точной ее формулировки.

Ранее мы упоминали о нарушении транзитивности при сравнении строк, содержащих символы замещения. Такая же проблема возникает при попытке определить “приближенное равенство” двух подстрок u_1 и u_2 одной строки x . Используя, например, расстояние Хемминга, имеем

$$d_H(ab, aa) = d_H(ab, bb) = 1.$$

Отсюда следует, что строки aa и bb можно считать “аппроксимирующими раппортами” для строки ab при условии, что “приближенное равенство” определяется как расстояние между строками, не превышающее $k = 1$. Но поскольку

$d_H(aa, bb) = 2$, при таком условии нельзя считать “равными” любые множества строк, содержащие подстроки или aa или bb . Таким образом, перед нами стоит нелегкий выбор — использовать аппроксимирующие раппорты только в паре (что увеличивает количество необходимых вычислений и все равно имеет на выходе неполный результат) либо использовать вместе только те аппроксимирующие раппорты, расстояние между которыми “мало” (что опять ведет к усложнению вычислений и порождает только частичный ответ).

Свойство продолжаемости, которое оказалось полезным для уменьшения объема вычислений при работе с точными раппортами, также находит применение при работе с аппроксимирующими раппортами. Как и выше, рассмотрим подстроки ab и aa некой строки x , расстояние между которыми равно 1, т.е. они “приближенно равны”. Предположим, что подстрока aa в строке x может появляться только как суффикс подстроки baa , а подстрока ab может появиться в двух вариантах: как подстрока строки bab или как подстрока строки aab (например, как в строках Фибоначчи). Поскольку $d_H(bab, baa) = 1$, то для подстрок bab аппроксимирующий раппорт $\{ab, aa\}$ в некотором смысле продолжаемый влево. Но для подстрок aab , где a предшествует подстроке ab , нельзя говорить о продолжаемости влево, так как $d_H(aab, baa) = 2 > k$. Таким образом, мы находимся в положении, когда с помощью аппроксимирующего раппорта $\{ab, aa\}$ можно определить только некоторые (но не все) вхождения подстроки ab в строку x .

Возможно, описанные трудности уменьшатся (или даже исчезнут), если вместо расстояния Хемминга использовать другую функцию расстояния. Например, $d_L(ab, abc) = 1$, и поэтому при использовании расстояния Левенштейна строки ab и abc могут быть аппроксимирующими раппортами для любого $k \geq 1$. Если аппроксимирующие раппорты определять исходя из функции расстояния, отличной от функции Хемминга, то исчезает (или по крайней мере уменьшается) проблема нахождения точных периодов. В общем случае использование различных функций расстояния порождает различные последствия для алгоритмов вычисления аппроксимирующих раппортов, поскольку эти алгоритмы учитывают специфику используемых функций расстояния. Например, если “приближенное равенство” двух строк u_1 и u_2 устанавливается при выполнении неравенства $d(u_1, u_2) \leq k$, то при использовании расстояния Хемминга все строки, длина которых не превышает k , будут приближенно равны между собой. Если же использовать расстояние Левенштейна, то строки u_1 и u_2 будут приближенно равны, если для них выполняется неравенство $|u_1| + |u_2| \leq k$.

Мы вернемся к обсуждению этих проблем в разделе 13.4. Сейчас сформулируем последнюю задачу, в сжатой форме подытоживающую это обсуждение.

Задача 2.21 (Вычисление аппроксимирующих кратных строк). Для данной функции расстояния d и целого числа $k \geq 0$ в заданной строке x вычислить

все аппроксимирующие кратные строки, отличающиеся не более, чем на расстояние k (раздел 13.4). ■

Упражнения 2.3

1. Докажите, что для любого целого $n \geq 0$ строка a^n содержит $\lfloor n^2/4 \rfloor$ вхождений различных квадратов.
2. Найдите максимальные кратные подстроки в строке Фибоначчи

$$f_6 = abaababaabaab$$

и запишите их в виде троек (i, p^*, r^*) . Подсчитайте количество таких троек. Аппроксимируйте этот результат на количество максимальных кратных подстрок в строке $f_7 = f_6(abaababa)$.

3. Строка x называется **слабо кратной**, если она представима в виде $x = u_1 u_2 \dots u_k$ для некоторого $k \geq 2$, где для любого $i \in 2..k$ подстрока u_i является перестановкой букв подстроки u_1 .
 - а) Покажите, что любая кратная строка также является слабо кратной.
 - б) Найдите все слабо кратные подстроки в строке f_6 .
4. Найдите все оболочки, если они существуют, для строк f_6 и f_7 .
5. Предложите определение оболочек u для петель $C(x)$. Совпадает ли множество оболочек петли $C(x)$ с множеством оболочек всех циклических сдвигов строки x ?
6. Охарактеризуйте множество оболочек петли $C((abc)^n)$.
7. На основе определения 2.3.2 докажите, что раппорт

$$M_{x,u} = (p; i_1, i_2, \dots, i_r)$$

будет продолжаемым влево последовательным раппортом только тогда, когда раппорт

$$M_{-1x,u'} = (p; i_1 - 1, i_2 - 1, \dots, i_r - 1)$$

будет продолжаемым вправо последовательным раппортом.

8. На основе определения 2.3.3 покажите, что для любого целого $t' \in 1..t - 1$ раппорт $(i + t', p^*, r^*)$ будет кратной строкой, продолжаемой как влево, так и вправо.

ГЛАВА 3

Такие разные строки

Там, где идеи терпят крах, становятся чрезвычайно необходимыми умелые слова.

— Иоганн Вольфганг Гете (1749–1837). Фауст

3.1 Проблема исключений и морфизмы

В этой главе мы рассмотрим три специальных типа строковых последовательностей, которые обладают интересными свойствами. Этим типам строк присущи следующие характеристические признаки.

- Строки можно классифицировать по свойствам их кратных подстрок.
- Строки можно сгенерировать путем применения специального вида подстановок, которые называются “морфизмами”.

Такие строки представляют интерес как экстремальные входные данные для некоторых строковых алгоритмов, особенно для тех, которые вычисляют характеристические паттерны. Но прежде чем вплотную заняться изучением таких строк, сначала рассмотрим более общую задачу.

Строковая последовательность x называется *исключаемой* паттерном p (частным или характеристическим), если отсутствуют вхождения паттерна p в строку x . В этом случае также будем говорить, что строка x *p -свободна*. Например, строка $f_5 = abaababa$ исключается паттерном $p = bb$, а строка $x = abacabca$ свободна от квадратов (т.е. в данной строке x нет квадратов — подстрок вида

u^2). Классическая проблема исключений в обобщенном виде формулируется как задача построения бесконечных строк, свободных от кратных подстрок.

Задача 3.1 (Построение (α, r) -исключений). Построить бесконечную строковую последовательность (либо бесконечную петлю) на алфавите размером α , свободную от кратных подстрок порядка r , но содержащую кратные подстроки порядков $2, 3, \dots, r - 1$. ■

Отметим, что в некоторых случаях эта задача тривиальна. Например, для $\alpha = 1$ и для любого $r > 1$ не существует строк (состоящих из повторений одной буквы) длины r или большей длины, свободных от кратных подстрок порядка r . Таким образом, невозможно построить $(1, r)$ -исключения. Упражнение 3.1.1 показывает, что также невозможно построить $(2, 2)$ -исключения, т.е. не существует бесконечных строковых последовательностей на бинарном алфавите, свободных от квадратов. Однако, как покажем далее, можно построить на бинарном алфавите бесконечные строки, свободные от кубов и кратных подстрок порядка 4; другими словами, можно построить $(2, 3)$ - и $(2, 4)$ -исключения.

Аксель Туе (Axel Thue) [220] первым исследовал проблему исключений, в частности, задачу построения $(2, 3)$ - и $(3, 2)$ -исключений. Мы рассмотрим эти построения: в разделе 3.2 строки на бинарном алфавите, свободные от кубов, и в разделе 3.3 строки на тернарном алфавите, свободные от квадратов. Наконец, в разделе 3.4 исследуем строки Фибоначчи (порождаемые числами Фибоначчи) и примеры построения $(2, 4)$ -исключений. Все эти примеры строковых последовательностей обладают замечательными свойствами, но они могут быть и представителями “плохих”, в определенном смысле, строк. Так, $(2, 3)$ -исключения Туе и строки Фибоначчи содержат “слишком много” кратных подстрок — фактически строки Фибоначчи содержат максимально (асимптотически) возможное количество кратных подстрок на множестве всех строковых последовательностей. Поэтому такие строки представляют наихудший случай входных данных для алгоритмов вычисления кратных строк. С другой стороны, $(3, 2)$ -исключения Туе вообще не содержат кратных подстрок, и, следовательно, такие алгоритмы будут выполняться на них впустую без получения конечного результата (в виде найденных кратных подстрок). Мы также покажем, что строки Фибоначчи представляют наихудший случай для алгоритмов, вычисляющих характеристические паттерны.

Интересное обобщение проблемы исключений (назовем это обобщение *задачей построения слабых (α, r) -исключений*) предложил Эрдёс (Erdős) [83] — построить, если возможно, бесконечные строковые последовательности на алфавите размера α , не содержащие слабо кратные подстроки порядка r (определение слабо кратных строк дано в упражнении 2.3.3). Поскольку любая кратная строка является слабо кратной, поэтому если задача построения “обычных” исключений

не имеет решения, то в этом случае не будет иметь решения и задача построения слабых исключений. Обратное утверждение, конечно, не верно.

В частности, Эрдёс поставил вопрос о минимальном значении α , при котором существуют бесконечные строки, не содержащие слабые (Абелевы) квадраты. В 1970 году Плесентс (Pleasant) [194] опубликовал решение задачи построения слабых (5, 2)-исключений, а недавно Керёнен (Ker nen) [132] нашел (путем машинного моделирования) решение задачи построения слабых (4, 2)-исключений. В упражнении 3.1.3 предложено доказать, что результат Керёнена наилучший из возможных.

Вернемся ко второму свойству рассматриваемых в этой главе строк, о котором упоминалось в начале раздела. Морфизмом называется отображение h , которое каждой букве λ из алфавита A ставит в соответствие строку $h(\lambda)$ из множества A^+ . Отображение h распространяется на любую строку $x = x[1..n]$ из множества A^+ путем использования следующего тождества

$$h(x) = h(x[1])h(x[2]) \dots h(x[n]). \quad (3.1)$$

Для полноты распространим отображение на множество A^* , положив, что для любого морфизма $h(\varepsilon) = \varepsilon$.

На основании равенства (3.1) любой морфизм h можно применять к исходной строке x_0 любое число раз, тем самым генерируя *последовательность итераций* $h^*(x_0)$ по следующему правилу

$$h^*(x_0) = \{h^0(x_0), h^1(x_0), h^2(x_0), \dots\},$$

где $h^0(x_0) \equiv x_0$ и для любого целого $k \geq 1$ $h^k(x_0) = h(h^{k-1}(x_0))$. Например, для *тождественного морфизма*, который задается как

$$h(\lambda) = \lambda, \forall \lambda \in A,$$

для любой строки $x \in A^*$ получим тождественную последовательность $h^*(x) = \{x, x, \dots\}$. Для алфавита $A = \{a, b\}$ другой простой морфизм $h(a) = a, h(b) = ab$ порождает следующие последовательности для различных начальных строк:

$$\begin{aligned} h^*(a) &= \{a, a, \dots\}, \\ h^*(b) &= \{b, ab, a^2b, \dots, a^k b, \dots\}, \\ h^*(ab) &= \{ab, a^2b, \dots, a^k b, \dots\}, \\ h^*(ba) &= \{ba, aba, a^2ba, \dots, a^k ba, \dots\} \end{aligned}$$

и т.д.

В следующих трех разделах данной главы изучаются бесконечные последовательности, порожденные следующими морфизмами.

Строки Туе, являющиеся (2, 3)-исключениями: $h(a) = ab, h(b) = ba.$ (3.2)

Строки Туе, являющиеся (3, 2)-исключениями: $h(a) = abcab,$
 $h(b) = acabcb,$
 $h(c) = acbcacb.$ (3.3)

Строки Фибоначчи: $h(a) = ab, h(b) = a.$ (3.4)

Каждый из этих морфизмов обладает дополнительным важным свойством, которое мы сейчас рассмотрим. Обозначим через $h^*(A)$ множество строк, содержащее все последовательности $h^*(\lambda), \lambda \in A$, т.е.

$$h^*(A) = \bigcup_{\lambda \in A} h^*(\lambda) \quad (3.5)$$

Будем говорить, что морфизм h является **кодом**, если для любой строки $x \in h^*(A)$ существует не более одной строки $y \in A^*$, такой, что $h(y) = x$. Очевидно, что для любой строки $x \in h^*(A) - A$ существует *не менее одной* такой строки y (что следует из равенства (3.5)), но если $x \in A$, то такой строки y может не существовать. Когда h является кодом и $h(y) = x$, назовем строку y **предшествующей** строке x и будем писать $y = h^{-1}(x)$, где h^{-1} — обратное отображение; если читатель позволит мне некоторую вольность языка, то в этом случае будем говорить об **обратном морфизме**.

Таким образом, код является морфизмом, имеющим обратное отображение. Однако следует помнить, что обратное отображение ограничено только множеством $h^*(A) - A$, для строк из множества $A^* - h^*(A)$ может не существовать предшествующих строк либо их может быть несколько. Например, морфизм

$$h : a \rightarrow a, \quad b \rightarrow ab, \quad c \rightarrow abc$$

на алфавите $\{a, b, c\}$ является кодом, хотя такие простые строки, как ba и c^2 , не имеют предшествующих строк. Аналогично морфизм

$$h : a \rightarrow a, \quad b \rightarrow ab, \quad c \rightarrow bab$$

также является кодом, но строка $abab$, которая не принадлежит множеству $h^*(A)$, имеет две различные предшествующие строки ac и b^2 .

Тот факт, что морфизмы являются кодами, играет ключевую роль при исследовании основных свойств строк Туе, являющихся (3, 2)-исключениями, и строк Фибоначчи.

Упражнения 3.1

1. Докажите, что не существует решения задачи построения $(2, 2)$ -исключений.
2. Для любого неотрицательного целого k на алфавите из $k + 1$ букв постройте свободные от квадратов строки длиной $n = 2^k$.
3. Докажите, что не существует решения задачи построения слабых $(3, 2)$ -исключений.
 - а) Пусть задан морфизм на алфавите $A = \{a, b, c, d\} : h(a) = b, h(b) = c, h(c) = d, h(d) = \varepsilon$. Для каждой строки $x \in A$ постройте $h^4(x)$ и $h^5(x)$.
 - б) Измените морфизм так, чтобы $h(d) = a$. Что в этом случае будет представлять собой строка $h^4(x)$?
4. Используя определение (3.1), покажите, что если $x = uv$, тогда выполняется равенство $h(x) = h(u)h(v)$.
5. Пусть заданы морфизм h и начальная строка x_0 . Будем говорить, что $x \equiv y$ только в том случае, когда и x и y являются элементами последовательности $h^*(x_0)$. Является ли это отношение “ \equiv ” отношением эквивалентности? Поясните свой ответ.
6. Пусть заданы алфавит A и морфизм h на этом алфавите. Возможно ли, чтобы последовательность $h^*(x_0)$ имела конечное число членов?
7. Является ли обратный морфизм морфизмом?
8. Определите, являются ли следующие морфизмы кодами. Поясните свой ответ.
 - а) $h : a \rightarrow aba, b \rightarrow ab, c \rightarrow aab$;
 - б) $h : a \rightarrow a, b \rightarrow ab, c \rightarrow ab$;
 - в) $h : a \rightarrow aba, b \rightarrow ab$.

3.2 Строки Туе, являющиеся $(2, 3)$ -исключениями

В этом разделе мы изучим строки Туе t_n , определенные с помощью морфизма (3.2): $h(a) = ab, h(b) = ba$. Положим

$$h^*(t_0) = \{t_0, t_1, t_2, \dots\},$$

где $t_n = h(t_{n-1})$ для любого $n \geq 1$. Пусть $t_0 = a$, тогда

$$\begin{aligned} t_1 &= ab, \\ t_2 &= abba, \\ t_3 &= abbabaab, \\ t_4 &= abbabaabbaabba, \\ &\vdots \end{aligned}$$

Очевидно, что $|t_n| = 2^n$. Заметим, что, по крайней мере, в первых нескольких строках t_n первые 2^{n-1} букв “сопряжены” со следующими 2^{n-1} буквами, т.е.

$$t_n[i + 2^{n-1}] = a \Leftrightarrow t_n[i] = b, \quad \forall i \in 1..2^{n-1}. \quad (3.6)$$

Если ввести запись $\bar{a} = b$ и $\bar{b} = a$ для обозначения сопряженности, тогда морфизм (3.2) для любой буквы $\lambda \in \{a, b\}$ можно переписать в виде

$$h : \lambda \rightarrow \lambda\bar{\lambda}. \quad (3.7)$$

Эту нотацию можно распространить и на произвольные строки из множества $\{a, b\}^*$: для данной строки x длиной n будем говорить, что строка y сопряжена со строкой x (запишем это как $y = \bar{x}$), если $|y| = n$ и $y[i] = \bar{x[i]}$ для всех $i \in 1..n$. Отметим, что согласно этому определению любая строка x имеет единственную сопряженную строку. Исключение составляет пустая строка ε , которая сопряжена сама с собой.

Пусть $h(x)$ обозначает строку, сформированную путем применения к строке x правила подстановки (3.7). Тогда

$$h(x) = x[1]\bar{x[1]}x[2]\bar{x[2]} \cdots x[n]\bar{x[n]},$$

и используя тот факт, что $\lambda = \bar{\bar{\lambda}}$, можем легко вычислить строку, сопряженную со строкой $h(x)$:

$$h(x) = x[1]\bar{x[1]}x[2]\bar{x[2]} \cdots x[n]\bar{x[n]} = h(\bar{x[1]\bar{x[2]} \cdots \bar{x[n]}) = h(\bar{x}). \quad (3.8)$$

Теперь можно формально доказать наличие в строках t_n отношения сопряженности (3.6).

Лемма 3.2.1. Для любого целого $n \geq 1$ $t_n = t_{n-1}\bar{t_{n-1}}$.

Доказательство. Мы знаем, что при $n = 1$ выполняется $t_1 = ab = a\bar{a} = t_0\bar{t_0}$. Предположим, что утверждение леммы выполняется для некоторого целого $n' = n - 1 \geq 1$. Далее методом математической индукции с использованием соотношений (3.1) и (3.8) получаем

$$t_n = h(t_{n-1}) = h(t_{n-2}\bar{t_{n-2}}) = h(t_{n-2})h(\bar{t_{n-2}}) = t_{n-1}\bar{t_{n-1}}.$$

Лемма доказана. ■

Из этой леммы вытекает, что правило подстановки (т.е. морфизм) для строк t_n *сохраняет префикс* — строка t_n имеет t_{n-1} своим префиксом. Как мы увидим в разделе 3.4, этим свойством обладают как строки Туе t_n , так и строки Фибоначчи f_n . В случае строк Туе сохранение префикса означает, что строка t_n имеет префиксы $t_i \bar{t}_i$, $i = 0, 1, \dots, n-1$, тогда как \bar{t}_n имеет префиксы $\bar{t}_i t_i$. Отсюда сразу следуют два утверждения.

Лемма 3.2.2.

- а) Строки t_n и \bar{t}_n не имеют общего непустого префикса;
- б) ни строка t_n , ни строка \bar{t}_n не имеют префиксом квадрат. ■

Мы оставляем полное доказательство этой леммы в качестве упражнения, а также доказательство аналогичной леммы, в утверждении которой “префикс” заменен на “суффикс”. Эти результаты будут использованы далее в доказательстве других лемм и теорем.

В следующей важной лемме встретим наших “старых друзей” — грани. Оказывается, что в строках Туе все грани являются строками Туе. Аналогичный результат также справедлив для строк Фибоначчи (лемма 3.4.4).

Лемма 3.2.3. Для любого целого $n \geq 2$ гранями строки t_n будут только строки t_i , где $i = n-2, n-4, \dots, n \bmod 2$.

Доказательство. Для доказательства леммы достаточно показать, что наибольшей гранью t_n является строка t_{n-2} . Поскольку

$$t_n = t_{n-1} \bar{t}_{n-1} = t_{n-2} \overline{t_{n-2} t_{n-2}} t_{n-2}, \quad (3.9)$$

то очевидно, что строка t_{n-2} является гранью строки t_n . Теперь надо показать, что не существует большей грани.

Предположим, что в строке t_n есть грань b , такая, что $|b| \geq 2^{n-1}$. В этом случае, как следует из равенства (1.1), строка t_n должна быть сильно периодической и поэтому должна иметь префиксом квадрат, что противоречит утверждению б) леммы 3.2.2. Поэтому грань b может быть только такой, что $|b| < 2^{n-1}$.

Поскольку обе строки b и t_{n-2} являются гранями строки t_n и b — наибольшая грань, то отсюда следует, что t_{n-2} должна быть гранью строки b . Из этого замечания и условия $|b| < 2^{n-1}$ вытекает, что строка b должна быть сильно периодической и поэтому должна иметь префиксом квадрат. Но b является префиксом строки t_n , которая не может иметь префиксом квадрат. Это противоречие доказывает, что грани, большей t_{n-2} , не существует. ■

Из последней леммы вытекает также “сопряженный” результат, что для любого $n \geq 2$ строка t_n может иметь гранями только строки \bar{t}_i , где $i = n-2, n-4, \dots$,

$n \bmod 2$. Отметим на будущее, что поскольку $t_{n+1} = t_n \overline{t_n}$ (соответственно, $\overline{t_{n+1}} = \overline{t_n t_n}$), лемма 3.2.3 эквивалентна утверждению о том, что префиксами строки t_n (соответственно, $\overline{t_n}$) могут быть только суффиксы строки $\overline{t_n}$ (соответственно, t_n), т.е. строки t_i (соответственно, $\overline{t_i}$), $i = n - 2, n - 4, \dots, n \bmod 2$.

Теперь продемонстрируем, как строки t_n можно использовать для решения задачи построения (2, 3)-исключений. Мы покажем, что хотя при любом $n \geq 2$ строка t_n имеет не менее одного квадрата, она не содержит кубов. Фактически, мы докажем следующий результат: любой квадрат в строке t_n является *непродолжаемым* в смысле определения, данного в разделе 2.3, т.е. докажем, что если для позиции i и некоторого целого $p \geq 1$ выполняется равенство

$$t_n[i..i + 2p - 1] = (t_n[i..i + p - 1])^2,$$

то ни подстрока $t_n[i - 1..i + 2p - 2]$, ни подстрока $t_n[i + 1..i + 2p]$ не будут квадратом.

Лемма 3.2.4. Пусть u^2 является квадратом в строке t_n для некоторого $n \geq 2$. Если $|u|$ — нечетное число, тогда строка u будет одним из элементов множества $\{a, b, aba, bab\}$.

Доказательство. Пусть $k = |u| \geq 1$. Вследствие морфизма (3.7) строка u^2 может иметь только одну из двух следующих форм:

а) $u^2 = \lambda_1 \overline{\lambda_1} \lambda_2 \overline{\lambda_2} \cdots \lambda_k \overline{\lambda_k}$, где $v = \lambda_1 \lambda_2 \cdots \lambda_k$ — подстрока строки t_{n-1} ;

б) $u^2 = \overline{\lambda_0} \lambda_1 \overline{\lambda_1} \cdots \lambda_{k-1} \overline{\lambda_{k-1}} \lambda_k$, где $v = \lambda_0 \lambda_1 \cdots \lambda_k$ — подстрока строки t_{n-1} .

В случае а), поскольку k нечетно, имеем

$$u = \lambda_1 \overline{\lambda_1} \cdots \overline{\lambda_{\lfloor k/2 \rfloor}} \lambda_{\lfloor k/2 \rfloor + 1} = \overline{\lambda_{\lfloor k/2 \rfloor + 1}} \lambda_{\lfloor k/2 \rfloor + 2} \cdots \lambda_k \overline{\lambda_k}.$$

Без потери общности можем положить, что $\lambda_1 = a$, тогда $\overline{\lambda_1} = b$. В этом случае

$$\begin{aligned} \overline{\lambda_{\lfloor k/2 \rfloor + 1}} &= \lambda_1 = a, & \lambda_{\lfloor k/2 \rfloor + 2} &= \overline{\lambda_1} = b, \\ \lambda_2 &= \overline{\lambda_{\lfloor k/2 \rfloor + 2}} = b, & \lambda_{\lfloor k/2 \rfloor + 3} &= \overline{\lambda_2} = a. \end{aligned}$$

Поскольку k нечетно, строка u должна начинаться и заканчиваться буквой a . С другой стороны, так как $\overline{\lambda_{\lfloor k/2 \rfloor + 1}} = a$, то, должно быть, $u[k] = \lambda_{\lfloor k/2 + 1 \rfloor} = b$. Получили противоречие с тем фактом, что строка u должна заканчиваться буквой a . Приходим к выводу, что форма а) представления строки u^2 невозможна.

Теперь предположим, что строка u^2 представима в форме б). Тогда

$$u = \overline{\lambda_0} \lambda_1 \overline{\lambda_1} \cdots \lambda_{\lfloor k/2 \rfloor} \overline{\lambda_{\lfloor k/2 \rfloor}} = \lambda_{\lfloor k/2 \rfloor + 1} \overline{\lambda_{\lfloor k/2 \rfloor + 1}} \cdots \lambda_{k-1} \overline{\lambda_{k-1}} \lambda_k.$$

Опять без потери общности положим $\lambda_k = a$. Тогда

$$\begin{aligned} \overline{\lambda_{\lfloor k/2 \rfloor}} &= \lambda_k = a, & \overline{\lambda_{k-1}} &= \lambda_{\lfloor k/2 \rfloor} = b, \\ \overline{\lambda_{\lfloor k/2 \rfloor - 1}} &= \lambda_{k-1} = a, & \overline{\lambda_{k-2}} &= \lambda_{\lfloor k/2 \rfloor - 1} = b. \end{aligned}$$

Для $k = 1$ строка u равна a или b , а для $k = 3$ — aba или bab . Однако для $k \geq 5$ строка u имеет префикс $ababa$ или $babab$. Более того, поскольку строка u начинается с буквы $\bar{\lambda}_0$, в строке t_n ей должна предшествовать буква λ_0 . Таким образом, в строке t_n должна присутствовать одна из подстрок $ababab$ или $bababa$. Но тогда в строке t_{n-1} обязательно будет присутствовать одна из “запретных” подстрок $\lambda_0\lambda_1\lambda_2$, равная aaa или bbb . Получили противоречие. Следовательно, $k < 5$. Лемма доказана. ■

Очевидно, что квадраты aa и bb действительно входят в строки t_n и \bar{t}_n при любом $n \geq 2$. Чтобы найти первые вхождения квадратов $(aba)^2$ и $(bab)^2$, рассмотрим строку

$$t_5 = abbabaabbaababbabaababbaabbabaab.$$

Мы найдем квадрат $(bab)^2$ в строке t_5 стоящим симметрично относительно позиции $2^4 - 1 = 15$ (вследствие симметрии сопряженный квадрат $(aba)^2$ центрирован относительно позиции 19). Тогда, вследствие леммы 3.2.1, можно утверждать, что все четыре различных квадрата, указанные в лемме 3.2.4, будут входить во все строки t_n при $n \geq 5$.

Также ясно, что если в строку t_n входит “нечетный” квадрат u^2 , тогда “четный” квадрат $(h^r(u))^2$ будет входить в строку t_{n+r} при любом положительном целом r . Теперь покажем, что в строку t_∞ входят *только* четные квадраты, являющиеся r -ми итерациями нечетных квадратов.

Рассмотрим квадрат u^2 , присутствующий в строке t_n , для которого $k = |u| \geq 2$ чётно. Строка u^2 (как и в случае нечетных квадратов) может иметь только одну из двух следующих форм:

- а) $u^2 = \lambda_1\bar{\lambda}_1\lambda_2\bar{\lambda}_2 \cdots \lambda_k\bar{\lambda}_k$, где $v = \lambda_1\lambda_2 \cdots \lambda_k$ — подстрока строки t_{n-1} ;
- б) $u^2 = \bar{\lambda}_0\lambda_1\bar{\lambda}_1 \cdots \lambda_{k-1}\bar{\lambda}_{k-1}\lambda_k$, где $v = \lambda_0\lambda_1 \cdots \lambda_k$ — подстрока строки t_{n-1} .

В случае а) имеем

$$u = \lambda_1\bar{\lambda}_1 \cdots \lambda_{k/2}\bar{\lambda}_{k/2} = \lambda_{k/2+1}\bar{\lambda}_{k/2+1} \cdots \lambda_k\bar{\lambda}_k.$$

Поэтому строка $(\lambda_1\lambda_2 \cdots \lambda_{k/2})(\lambda_{k/2+1}\lambda_{k/2+2} \cdots \lambda_k)$ является квадратом, также принадлежащим строке t_{n-1} . Аналогично исследуется случай б):

$$u = \bar{\lambda}_0\lambda_1\bar{\lambda}_1 \cdots \lambda_{k/2-1}\bar{\lambda}_{k/2-1}\lambda_{k/2} = \bar{\lambda}_{k/2}\lambda_{k/2+1}\bar{\lambda}_{k/2+1} \cdots \lambda_{k-1}\bar{\lambda}_{k-1}\lambda_k,$$

тогда строка

$$(\lambda_1\lambda_2 \cdots \lambda_{k/2})(\lambda_{k/2+1}\lambda_{k/2+2} \cdots \lambda_k) \tag{3.10}$$

является квадратом, также принадлежащим строке t_{n-1} .

Отсюда следует, что любой четный квадрат длиной $2k$, присутствующий в строке t_n , можно получить из квадрата длиной k , принадлежащего строке t_{n-1} . В свою очередь квадрат длиной k , принадлежащий строке t_{n-1} , если он четный,

можно получить из соответствующего квадрата строки t_{n-2} либо, если он нечетный, он принимает одну из четырех возможных форм, определенных леммой 3.2.4. Таким образом, мы доказали следующий результат.

Лемма 3.2.5. Каждый квадрат, присутствующий в строке t_∞ , представим в форме $(h^r(\mathbf{u}))^2$, где r — неотрицательное целое число и $\mathbf{u} \in \{a, b, aba, bab\}$. ■

Теперь рассмотрим нечетный квадрат \mathbf{u}^2 из строки t_n длиной $2k = 2$. В соответствии с утверждением б) леммы 3.2.2 квадрат \mathbf{u}^2 не может быть ни префиксом ни суффиксом строки t_n . Поэтому в t_n должна быть подстрока $v = buub$ длиной $2k + 2$ со свойствами

$$v[1] \neq \mathbf{u}[k], \quad v[2k + 2] \neq \mathbf{u}[1]. \quad (3.11)$$

Эти свойства показывают, что квадрат \mathbf{u}^2 непродолжаемый. Фактически, нет необходимости доказывать, что такими свойствами обладают все квадраты длиной $2k$, поскольку они получены путем применения преобразования h^r к строке $baab$, а все квадраты $h^r(aa)$ непродолжаемы. Этот факт предложено доказать в упражнении 3.2.7, а в упражнении 3.2.8 показано, что все квадраты $h^r(abaaba)$ также непродолжаемы. Все это доказывает основной результат данного раздела.

Теорема 3.2.6. Все квадраты в строках Туе t_n непродолжаемые. ■

Бесконечное множество строк $h^*(t_0)$, введенное в начале этого раздела, содержит все конечные строки t_n , но не содержит бесконечной строки Туе t_∞ . Однако заметим, что если строка t_∞ содержит куб, тогда этот куб должен быть конечным и должен быть подстрокой некоторой конечной строки Туе, но последнее запрещено теоремой 3.2.6. Таким образом, получаем следующий результат.

Теорема 3.2.7. Бесконечная строка Туе t_∞ является решением задачи построения (2, 3)-исключений. ■

Краткое и, возможно, более элегантное доказательство этого результата приведено в работе [39]. Однако наше более длинное доказательство позволило получить много дополнительной информации о структуре строки t_∞ .

В заключение этого раздела рассмотрим один важный вычислительный вопрос: сколько всего квадратов присутствует в строке t_n или, по крайней мере, сколько квадратов имеют форму $h^r(a)$ или $h^r(b)$. Условимся называть квадраты, которые можно записать в такой форме, **регулярными**, поскольку любой такой квадрат строится путем сочленения суффикса строки t_n (соответственно, \overline{t}_n) с префиксом строки \overline{t}_n (соответственно, t_n). Как следует из леммы 3.2.3, регулярные квадраты \overline{t}_i ($i = n - 2, n - 4, \dots, n \bmod 2$) формируются в центре строки $t_n = t_{n-1}\overline{t}_{n-1}$. Квадратов в строке t_n , которые сформированы из квадратов строк

t_{n-1} и $\overline{t_{n-1}}$, всего $\lfloor n/2 \rfloor$. Отсюда получаем рекуррентное соотношение

$$q(n) - 2q(n-1) = \lfloor n/2 \rfloor, \quad n \geq 2, \quad (3.12)$$

где $q(n)$ — общее количество регулярных квадратов в строке t_n . Предполагая, что n нечетное число, решаем уравнение (3.12) стандартным способом.

$$\begin{aligned} q(n) - 2q(n-1) &= (n-1)/2, \\ 2q(n-1) - 4q(n-2) &= 2(n-1)/2, \\ &\vdots \\ 2^{n-3}q(3) - 2^{n-2}q(2) &= 2^{n-3}(2/2), \\ 2^{n-2}q(2) - 2^{n-1}q(1) &= 2^{n-2}(2/2). \end{aligned}$$

С учетом того, что $q(1) = 0$, суммируем эти уравнения, в результате получаем

$$\begin{aligned} q(n) &= \frac{3}{2} [(n-1) + (n-3)2^3 + (n-5)2^4 + \dots + (2)2^{n-3}] = \\ &= 3 \sum_{i=0}^{(n-3)/2} (n-2i-1)2^{2i-1}. \end{aligned}$$

Отсюда нетрудно получить границы для $q(n)$:

$$3 \cdot 2^{n-3} \leq q(n) \leq 3 \cdot 2^{n-2}.$$

Следовательно, $q(n)$ имеет порядок $\Theta(2^n)$. Такой же порядок получим в случае четного n . Это доказывает следующую теорему.

Теорема 3.2.8. Количество $q(n)$ регулярных квадратов в строке t_n имеет порядок $\Theta(|t_n|)$. ■

Более сложным способом Джемси Симпсон (Jamie Simpson) нашел точное значение величины $q(n)$, его формулу предложено доказать в упражнении 3.2.10.

В разделе 3.4 мы увидим, что строки Фибоначчи f_n длиной f_n содержат кратные подстроки, количество которых имеет порядок $\Theta(f_n \log f_n)$, т.е. строки Фибоначчи асимптотически содержат “больше” кратных подстрок, чем строки Туе t_n .

Упражнения 3.2

1. Пусть морфизм (3.7) применяется к начальной строке $t'_0 = b$, порождая *обратные строки Туе* t'_n , $n = 0, 1, \dots$. Покажите, что $t'_n = \overline{t_n}$.

2. Напомним, что строка называется *палиндромом*, если она одинаково читается слева направо и справа налево. Покажите, что для любого целого n строка t_n является палиндромом.
3. Формула морфизма (3.7) и лемма 3.2.1 дают основу для ряда итерационных алгоритмов вычисления строк t_n . Разработайте несколько таких алгоритмов и сравните их эффективность.
4. Докажите лемму, “сопряженную” с леммой 3.2.2.
5. Почему при доказательстве леммы 3.2.3 мы рассматривали только грани длиной больше $|t_{n-2}|$? Завершите доказательство этой леммы.
6. В доказательстве леммы 3.2.5 есть один “неясный” момент: в строке t_{n-1} , кроме квадрата (3.10), присутствует также квадрат

$$(\lambda_0 \lambda_1 \cdots \lambda_{k/2-1})(\lambda_{k/2} \lambda_{k/2+1} \cdots \lambda_{k-1}).$$

Поэтому квадрат (3.10) будет продолжаемым влево, что противоречит теореме 3.2.6. Как вы проясните этот “неясный” момент?

7. Покажите, что свойства (3.11) выполняются также для всех строк вида

$$v = h^r(baab) = h^r(b)h^r(a^2)h^r(b),$$

где $u = h^r(a)$, $|u| = k$, $r \geq 0$. Отсюда будет следовать, что все квадраты u^2 непродолжаемые.

8. Покажите, что в строке t_n все квадраты вида $h^r((aba)^2)$ непродолжаемые.
9. Докажите, что если строка u^2 является квадратом, содержащимся в строке t_n , тогда $|u| = 2^r$ или $3 \cdot 2^r$ для некоторого целого $r \geq 0$.
10. Методом математической индукции докажите формулу Джемси Симпсона для $q(n)$, $n \geq 0$,

$$q(n) = 2^{n+1}/3 - n/2 - \varepsilon_n,$$

где $\varepsilon_n = 2/3$, если n четно, и $\varepsilon_n = 5/6$, если n нечетно.

11. Можно ли вычислить количество *всех* квадратов, включая регулярные, находящихся в строке t_n ?
12. Покажите, что морфизм (3, 2) является кодом.
13. В работе [39] предложено генерировать строку t_∞ следующим образом: сначала записываются все неотрицательные целые числа в двоичной форме

$$0, 1, 10, 11, 100, 101, 110, \dots,$$

затем в этих числах цифры складываются по модулю 2:

$$0, 1, 1, 0, 1, 0, 0, \dots,$$

после этого удаляются запятые и 0 заменяется на букву a , а 1 — на букву b . Все, строка t_∞ готова! Докажите или опровергните, что таким способом действительно можно построить строку t_∞ .

3.3 Строки Туе, являющиеся (3, 2)-исключениями

В этом разделе мы изучим строки Туе $\tau\tau_n$, которые в определенном смысле противоположны строкам Туе t_n , изученным в предыдущем разделе: строки τ_n вообще не содержат кратных подстрок. Эти строки формируются с помощью морфизма (3.3):

$$h(a) = abcab, \quad h(b) = acabcb, \quad h(c) = acbcacb.$$

Положим $h^*(\tau_0) = \{\tau_0, \tau_1, \tau_2, \dots\}$, где $\tau_n = h(\tau_{n-1})$ для любого целого $n \geq 1$. Если $\tau_0 = a$, тогда получим

$$\begin{aligned} \tau_1 &= abcab, \\ \tau_2 &= (abcab)(acabcb)(acbcacb)(abcab)(acabcb), \\ &\vdots \end{aligned}$$

Каждая из трех компонент $abcab$, $acabcb$, $acbcacb$ морфизма h называется **блоком**. Отметим, что любой блок B можно записать в форме IT , где строки I и T имеют префиксом букву a и суффиксом строку, построенную на алфавите $\{b, c\}$. Назовем строки I и T соответственно **инициатором** и **терминатором**. Отметим на будущее, что все инициаторы и терминаторы различны:

$$\{abc, ac, acbc; ab, abcb, acb\}.$$

Фактически, множество инициаторов и терминаторов в точности совпадает с множеством всех строк длиной 2 или 3, имеющих префикс a и свободный от квадратов суффикс, содержащий только буквы $\{b, c\}$.

Полезность инициаторов и терминаторов проявляется в том, что если в какой-либо строке из множества $h^*(\tau_0)$ найдется инициатор, тогда за ним следует терминатор. И, аналогично, если в такой строке найден терминатор, то ему предшествует инициатор. Кроме того, каждая строка τ_n ($n > 0$) из множества $h^*(\tau_0)$ является сочленением блоков, т.е. тех же инициаторов и терминаторов. Поэтому для любого блока B в строке τ_n с помощью обратного морфизма h^{-1} можно определить предшествующую подстроку $h^{-1}(B)$.

Лемма 3.3.1. Строка Туе τ_n будет свободной от квадратов только тогда, когда будет свободной от квадратов строка τ_{n-1} .

Доказательство. Очевидно, что строка τ_{n-1} будет свободной от квадратов, если строка τ_n свободна от квадратов. Теперь предположим, что при некотором $n \geq 1$ в строке τ_n существует квадрат $u = u_L u_R$.

Если \mathbf{u}_R имеет инициатор \mathbf{I}_i в качестве префикса для некоторого $i \in 1..3$, тогда для некоторого $k \geq 0$ подстрока \mathbf{u}_L представима в виде

$$\mathbf{u}_L = (\mathbf{I}_i \mathbf{T}_i)(\mathbf{I}_{j_1} \mathbf{T}_{j_1}) \cdots (\mathbf{I}_{j_k} \mathbf{T}_{j_k}).$$

Подстроки \mathbf{u}_L и \mathbf{u}_R состоят из целого числа блоков, и поэтому определено обратное отображение $h^{-1}(\mathbf{u}_L)$, а также $h^{-1}(\mathbf{u}_L \mathbf{u}_R)$, которое будет квадратом, принадлежащим строке τ_{n-1} .

Подобным образом, если \mathbf{u}_R содержит терминатор \mathbf{T}_i в качестве префикса, для некоторого $k \geq 0$ имеем

$$\mathbf{I}_i \mathbf{u}_L = (\mathbf{I}_i \mathbf{T}_i)(\mathbf{I}_{j_1} \mathbf{T}_{j_1}) \cdots (\mathbf{I}_{j_k} \mathbf{T}_{j_k}) \mathbf{I}_i = \mathbf{v}_L \mathbf{I}_i.$$

Тогда $\mathbf{v}_R = \mathbf{T}_i(\mathbf{I}_{j_1} \mathbf{T}_{j_1}) \cdots (\mathbf{I}_{j_k} \mathbf{T}_{j_k})$ будет префиксом подстроки \mathbf{u}_R . Поэтому $\mathbf{v}^2 = \mathbf{v}_L \mathbf{I}_i \mathbf{v}_R$ будет квадратом, где $\mathbf{v} = \mathbf{v}_L$ состоит из целого числа блоков. Следовательно, существует обратное отображение $h^{-1}(\mathbf{v})$, которое порождает квадрат, принадлежащий строке τ_{n-1} .

Теперь предположим, что \mathbf{u}_R не имеет буквы a в качестве префикса. Поскольку не существует подстрок-квадратов, которые не содержали бы букву a , эта буква будет хотя бы по одному разу входить как в подстроку \mathbf{u}_L , так и в подстроку \mathbf{u}_R . Поэтому или в некотором инициаторе \mathbf{I}_i или в некотором терминаторе \mathbf{T}_i выполняется равенство $\mathbf{u}_R[1] = \mathbf{u}_L[1]$. Далее рассуждения, подобные приведенным выше, приводят к заключению, что в строке τ_n присутствует квадрат \mathbf{v}^2 , такой, что существует $h^{-1}(\mathbf{v}^2)$, порождающее квадрат в строке τ_{n-1} .

Мы рассмотрели все возможные случаи присутствия квадратов в строке τ_n , которые приводят к присутствию квадратов в строке τ_{n-1} . Лемма доказана. ■

Из леммы 3.3.1 немедленно следует теорема.

Теорема 3.3.2. Если начальная строка τ_0 свободна от квадратов, то свободными от квадратов будут все строки из множества $h^*(\tau_0)$. ■

Для окончания анализа этого типа строк Туге осталось заметить, что если бесконечная строка τ_∞ содержит квадрат, то в этом случае квадраты должны содержать некоторые конечные строки τ_n , что противоречит теореме 3.3.2. Следовательно, справедлива такая теорема.

Теорема 3.3.3. Бесконечная строка Туге τ_∞ является решением задачи построения (3, 2)-исключений. ■

Морфизм (3.3) был первым исследуемым преобразованием, который генерирует свободные от квадратов строки на алфавите из трех букв [29], но существуют и другие подобные морфизмы. Один из них задается следующими формулами [39, 75].

$$h(a) = acb, \quad h(b) = c, \quad h(c) = ab. \quad (3.13)$$

Существует очень “симпатичный” способ показать, что этот морфизм действительно порождает бесконечные строки, свободные от квадратов. Имея строку T_{∞}

$$t_{\infty} = abbabaabbaababba \dots,$$

сформируем новую строку u_{∞} на алфавите $\{0, 1, 2\}$ путем подсчета в строке t_{∞} букв b между последовательными буквами a . Получим

$$u_{\infty} = 2102012 \dots$$

Нетрудно заметить, что если в строке u_{∞} существует квадрат, тогда в строке t_{∞} существует продолжаемый квадрат, что противоречит теореме 3.2.6. Следовательно, строка u_{∞} свободна от квадратов. Наконец, отметим, что строка u_{∞} порождена морфизмом

$$h(2) = 210, \quad h(0) = 1, \quad h(1) = 20. \quad (3.14)$$

Подстановка $\{0, 1, 2\} \rightarrow \{b, c, a\}$ показывает, что данный морфизм совпадает с морфизмом (3.13). В упражнении 3.3.5 показаны детали этого соответствия.

В другом морфизме, который также решает задачу построения $(3, 2)$ -исключений [155], отображение $h(b)$ получается из отображения $h(a)$, а $h(c)$ — из $h(b)$ путем циклической перестановки букв алфавита:

$$\begin{aligned} h(a) &= abcbaabcba, \\ h(b) &= bcacbacacb, \\ h(c) &= cabacabac. \end{aligned} \quad (3.15)$$

Морфизмы на больших алфавитах, порождающие строки, свободные от квадратов, также представляют интерес. Рассмотрим, например, известную головоломку “Ханойские башни”: имеется три стержня, пронумерованные от 1 до 3, и n дисков разного диаметра. Первоначально все диски собраны на стержне 1 в убывающем порядке их диаметров (наименьший диск находится на верху пирамиды). Цель головоломки заключается в том, чтобы переместить диски (по одному за один раз) на другой стержень так, чтобы они опять располагались в первоначальном порядке. Разрешается помещать только меньший диск на больший либо на свободный стержень — другие перемещения дисков запрещены. Возможно всего шесть типов перемещения: со стержня 1 на стержень 2, со стержня 2 на стержень 3, со стержня 3 на стержень 1 и в обратном порядке. Закодируем эти перемещения буквами a, b, c и $\bar{a}, \bar{b}, \bar{c}$ соответственно. На алфавите $\{a, b, c, \bar{a}, \bar{b}, \bar{c}\}$ построим морфизм

$$\begin{aligned} h(a) &= a\bar{c}, & h(b) &= c\bar{b}, & h(c) &= b\bar{a}. \\ h(\bar{a}) &= ac, & h(\bar{b}) &= cb, & h(\bar{c}) &= ba. \end{aligned} \quad (3.16)$$

Этот морфизм решает задачу построения $(6, 2)$ -исключений и, следовательно, находит минимальную последовательность перемещений дисков в нашей головоломке, поскольку, очевидно, эта минимальная последовательность не должна содержать повторяющихся подпоследовательностей. Подробно этот морфизм рассмотрен в работе [8].

Упражнения 3.3

1. Покажите, что морфизм (3.3) является кодом.
2. Разработайте алгоритм, который на основании морфизма (3.3) по начальной строке τ_0 генерирует заданное количество строк τ_n .
3. Приведите полное и подробное доказательство леммы 3.3.1.
4. Покажите, что шесть блоков I_i и T_i , $i = 1, 2, 3$, можно попарно упорядочить 120 разными способами. Исходя из этого ответьте, сколько можно определить морфизмов, порождающих строки, свободные от квадратов?
5. Чтобы убедиться в том, что морфизм (3.13) действительно порождает строки, свободные от квадратов, докажите следующее.
 - а) Если строка u_∞ содержит квадрат, тогда строка t_∞ содержит подстроку $x^2\lambda$ с наибольшей гранью длиной 1.
 - б) Морфизм (3.14) порождает строку u_∞ как предел последовательности строк из множества $h^*(0)$.
6. Покажите, что морфизм (3.13) порождает строки длиной $3 \cdot 2^i$, $i \geq 0$, циклические сдвиги которых свободны от квадратов. Сравните эту задачу с “исследовательской задачей” из упражнения 1.4.4.
7. Херенди Тамаш (Herendi Tamás) заметил, что морфизм (3.13) не обладает одним хорошим свойством, которое присуще морфизму Туе, а именно: морфизм (3.13) не отображает *любую* строку, свободную от квадратов, в строку, также свободную от квадратов. Приведите пример, подтверждающий наблюдение Тамаша.
8. Докажите, что морфизм (3.15) решает задачу построения $(3, 2)$ -исключений.
9. Докажите, что морфизм (3.16) решает задачу построения $(6, 2)$ -исключений. (Решение приведено в работе [8].)

3.4 Строки Фибоначчи

Эту главу мы завершим изучением строк Фибоначчи, порождаемых морфизмом (3.4): $h(a) = ab$, $h(b) = a$. Как и в предыдущих разделах, введем множество $h^*(f_0) = \{f_0, f_1, f_2, \dots\}$, где $f_n = h(f_{n-1})$ для любого целого $n \geq 1$ и для

определенности положим $f_0 = b$. Далее мы увидим, что длины $f_n = |f_n|$ этих строк являются числами Фибоначчи. Не будет преувеличением сказать, что числа Фибоначчи “окупили” комбинаторику и не только — они привлекают такой научный интерес, что им посвящен отдельный журнал *The Fibonacci Quarterly*, где публикуются работы, изучающие их свойства и свойства других итерационных числовых последовательностей.

Приведем несколько первых строк Фибоначчи.

$$\begin{aligned} f_0 &= b, \\ f_1 &= a, \\ f_2 &= ab, \\ f_3 &= aba, \\ f_4 &= abaab, \\ f_5 &= abaababa, \\ f_6 &= abaababaabaab, \\ &\vdots \end{aligned} \tag{3.17}$$

Длины этих строк равны соответственно $f_0 = 1, f_1 = 1, f_2 = 2, f_3 = 3, f_4 = 5, f_5 = 8, f_6 = 13$. Строки Фибоначчи связаны соотношением, аналогичным соотношению между числами Фибоначчи.

Лемма 3.4.1. Строки Фибоначчи удовлетворяют рекуррентному соотношению

$$f_0 = b, \quad f_1 = a, \quad f_n = f_{n-1}f_{n-2}, \quad n \geq 2. \tag{3.18}$$

Доказательство. леммы легко получить методом математической индукции (см. упражнение 3.4.2). ■

Отсюда получаем, что длины строк f_n удовлетворяют соотношению

$$f_n = f_{n-1} + f_{n-2}, \quad \forall n \geq 2, \tag{3.19}$$

которое является рекуррентным соотношением для чисел Фибоначчи [135].

Из соотношения (3.18) следует, что, подобно строкам Туе, строки Фибоначчи f_n при $n > 0$ сохраняют префикс. Это свойство строк Фибоначчи позволяет ввести *бесконечную строку Фибоначчи* f_∞ как такую, префиксом которой будет любая строка $f_n, n > 0$.

В разделе 3.2 мы видели, что длина строк Туе t_n равна 2^n , т.е. длина удваивалась при каждом применении отображения h . Поскольку длина строк Фибоначчи “почти” удваивается при увеличении номера строки на 1, естественно предположить, что длина f_n будет экспонентой от n . Если это действительно так, тогда

существует действительное число $\phi > 1$, такое, что $f_n \in \Omega(\phi^n)$. Чтобы доказать это, заметим, что из соотношения (3.19) следуют неравенства $f_{n-1} \geq \phi^{n-2}$ и $f_{n-2} \geq \phi^{n-3}$. Отсюда с необходимостью получаем

$$f_n \geq \phi^{n-2} + \phi^{n-3} = \phi^{n-3}(\phi + 1). \quad (3.20)$$

Если число ϕ удовлетворяет уравнению

$$\phi^2 = \phi + 1, \quad (3.21)$$

тогда автоматически получим неравенство $f_n \geq \phi^{n-1}$ и, следовательно, аналогичные неравенства для f_{n-1} и f_{n-2} . Уравнение (3.21) — это простое квадратное уравнение относительно ϕ , положительный корень которого имеет значение

$$\phi = (1 + \sqrt{5}) / 2 \approx 1,618034. \quad (3.22)$$

Это число, интересное само по себе, часто называется “золотым значением”¹ [135]. Для проверки значения (3.22) заметим, что

$$f_0 = 1 \geq \phi^{-1} = (\sqrt{5} - 1) / 2 \approx 0,618034 \quad \text{и} \quad f_1 = 1 \geq \phi^0 = 1.$$

Из соотношений (3.20) и (3.21) следует, что $f_2 \geq \phi^1$, тогда $f_3 \geq \phi^2$ и т.д. Отсюда по индукции получаем основной результат

$$f_n \geq \phi^{n-1}, \quad (3.23)$$

где число ϕ задается формулой (3.22). В упражнении 3.4.4 предложено доказать, что $f_n \leq \phi^n$.

Завершив изложение предварительных сведений о строках Фибоначчи, вернемся к основным свойствам этих строк. Установим два факта, касающиеся строк Фибоначчи, которые на первый взгляд кажутся противоречащими друг другу.

- В терминах троек (i, p^*, r^*) , введенных в связи с задачей 2.15, покажем, что в строке Фибоначчи f_n число кратных подстрок имеет порядок $\Theta(f_n \log f_n)$. (Таким образом, строки Фибоначчи являются “наихудшими” входными данными для алгоритмов вычисления кратных подстрок в строковых последовательностях.)
- С другой стороны, мы также покажем, что *все* кратные подстроки в строках Фибоначчи можно найти за время порядка $\Theta(n)$, используя для этого кодировку (i, p^*, r^*, t) . Эта кодировка основана на “сериях” кратных подстрок (см. формальное определение серий в разделе 2.3). Если говорить не формально, то серия — это кратная строка $u[1..p^*]^{r^*}$, следующая за (возможно, пустым) собственным префиксом $u[1..t]$ подстроки $u[1..p^*]$.

¹В отечественной литературе это значение обычно называется “коэффициентом золотого сечения”. — *Примеч. ред.*

К этим фактам “из жизни” строк Фибоначчи мы придем кружной дорогой, “посетив по пути некоторые интересные места” свойств этих строк. Как упоминалось ранее, строки Фибоначчи должны быть решением задачи построения (2, 4)-исключений — мы покажем, что они содержат кубы, но не содержат кратных подстрок четвертого порядка.

Начнем обобщение идеи строк Фибоначчи следующим способом. Напомним, что морфизм $h : A \rightarrow A^*$ (см. формулы (3.4)) оперирует двумя буквами a и b : $h(a) = ab$, $h(b) = a$. Расширим это отображение на множество A^* , т.е. будем считать, что это отображение h действует $A^* \rightarrow A^*$, оперируя при этом двумя произвольными строками x и y из множества A^* :

$$h(x) = xy, h(y) = x. \quad (3.24)$$

Таким образом, морфизм (3.4) будет частным случаем отображения (3.24) при $x = a$ и $y = b$. Аналогично (3.17) можно вычислить

$$\begin{aligned} h^2(y) &= xy, \\ h^3(y) &= xyx, \\ h^4(y) &= xyxxy, \\ &\vdots \end{aligned}$$

Теперь определим n -ю *обобщенную строку Фибоначчи* как

$$f_n(x, y) = h^n(y), \quad (3.25)$$

где $h^n(y)$ — $(n + 1)$ -й член последовательности $h^*(y) = \{y, h(y), h^2(y), \dots\}$. Несколько первых обобщенных строк имеют вид

$$\begin{aligned} f_0(x, y) &= y, \\ f_1(x, y) &= x, \\ f_2(x, y) &= xy, \\ f_3(x, y) &= xyx, \\ f_4(x, y) &= xyxxy, \\ f_5(x, y) &= xyxxyxyx, \\ f_6(x, y) &= xyxxyxyxxyxy, \\ &\vdots \end{aligned}$$

В общем случае

$$f_n(x, y) = f_{n-1}(x, y)f_{n-2}(x, y), \quad n \geq 2. \quad (3.26)$$

Определим *бесконечную обобщенную строку Фибоначчи* $f_\infty(x, y)$ как строку, содержащую все строки $f_n(x, y)$, $n \geq 0$, в качестве префиксов.

Поскольку $h^n(y) = h^{n-k}(h^k(y))$, из определения (3.25) следует, что

$$f_n(x, y) = f_{n-k}(h^k(x), h^k(y)).$$

Тогда, положив $x = a$ и $y = b$ и заметив, что $h^k(x) = h^{k+1}(y)$, получим начальный результат, относящийся к обычным строкам Фибоначчи.

Лемма 3.4.2. Для любых целых k и n , таких, что $0 \leq k \leq n - 1$, справедливо соотношение $f_n = f_{n-k}(f_{k+1}, f_k)$. ■

Этот результат позволяет выразить любую строку Фибоначчи f_n через выбранные строки f_{k+1} и f_k . Например, строку f_7 можно выразить через строки f_3 и f_2 :

$$f_7 = f_5(f_3, f_2) = (aba)(ab)(aba)(aba)(ab)(aba)(ab)(aba).$$

Лемма 3.4.2 применима и в бесконечном случае:

$$f_\infty = f_\infty(f_{n+1}, f_n) = f_{n+1}f_n f_{n+1}f_n f_{n+1}f_n f_{n+1}f_n \cdots \quad (3.27)$$

Это равенство используется в упражнении 3.4.11 для доказательства двух следующих лемм.

Лемма 3.4.3. Для любого целого $n \geq 2$ выполняется равенство $f_n^2 = f_{n+1}f_{n-2}$. ■

Лемма 3.4.4. Для любого целого $n \geq 3$ строка f_n имеет грани f_i для $i = n - 2, n - 4, \dots, 2 - (n \bmod 2)$. ■

Из этих лемм и свойства сохранения префиксов в строках Фибоначчи следует, что при $n \geq 3$ любое вхождение строки f_{n+1} в строку f_∞ порождает в этой строке квадрат f_n^2 . Строка f_n может располагаться или после подстроки $f_{n+1}f_n$ или после подстроки $f_{n+1}^2 = f_{n+1}f_n f_{n-1}$. Повторное применение леммы 3.4.1 показывает, что

$$f_n f_{n+1} f_n = f_n^3 f_{n-3} f_{n-2}. \quad (3.28)$$

Таким образом, любое вхождение строки f_n в выражение (3.27) возможно только в виде куба f_n^3 . Кроме того, для любого $n \geq 3$ строки f_n и f_n^2 в действительности будут *оболочками* строки f_∞ в том смысле, как они (оболочки) определены в формулировке задачи 2.18. Поскольку f_n — оболочка f_∞ , тогда любое другое²

²Здесь имеются в виду вхождения строки f_n в строку f_∞ , отличные от квадратов и кубов. — *Примеч. ред.*

вхождение строки f_n в строку f_∞ , которое не содержится в кратных подстроках выражения (3.27), должно быть некоторым циклическим сдвигом $R_j(f_n)$, $j > 0$. Но по теореме 1.4.2 это возможно только в случае, если f_n является кратной строкой, что, конечно, не соответствует действительности. Таким образом, мы рассмотрели все возможные случаи вхождения строки f_n (и, следовательно, f_n^2 и f_n^3) в строку f_∞ .

Зная о вхождениях f_n^3 в строку f_∞ , можем сделать предположение о существовании квадратов любого циклического сдвига $R_j(f_n)$, $j > 0$. Это предположение доказано в упражнении 3.4.12.

Теперь мы можем показать, что в строке f_∞ не существует четвертых степеней от f_n . На основе леммы 3.4.1 можно записать

$$f_n^4 = f_n^3 f_{n-2} f_{n-3} f_{n-2}.$$

Сравнивая это выражение с выражением (3.28), делаем заключение, что f_n^4 может существовать только тогда, когда выполняется равенство $f_{n-3} f_{n-2} = f_{n-2} f_{n-3}$. Как показано в упражнении 3.4.13, такое равенство невозможно, поэтому f_n^4 не существует.

Далее встает естественный вопрос: можно ли утверждать, что *все* квадраты в строке f_∞ являются квадратами конечных строк Фибоначчи или их циклических сдвигов? Для ответа на этот вопрос расширим понятие обратного морфизма, введенного в разделе 3.1. Там обратное отображение было определено только на множестве строк $h^*(A) - A$. В случае строк Фибоначчи мы можем определить обратное отображение h^{-1} на более коротких подстроках, если учтем, что подстрока ab может появиться в строке f_∞ только как результат отображения $h : a \rightarrow ab$, тогда как подстрока aa может появиться только в том случае, если первая буква a является самым левым элементом строковой последовательности и получена в результате отображения $h : b \rightarrow a$. С учетом этих замечаний определим обратный морфизм для строк Фибоначчи следующим образом:

$$h^{-1} : ab \rightarrow a, \quad a \not\rightarrow b, \tag{3.29}$$

где запись $a \not\rightarrow b$ означает, что отображение $a \rightarrow b$ возможно только в том случае, если за буквой a не следует (справа) буква b (другими словами, отображение $a \rightarrow b$ применяется только тогда, когда за буквой a следует другая буква a либо эта буква a является самым правым элементом строки). Поскольку морфизм Фибоначчи (3.4) не может породить подстроки bb , поэтому отображение (3.29) действительно определяет алгоритм, который должен выполняться над строкой f_n слева направо, проверяя каждое вхождение буквы a и выполняя обратное отображение h^{-1} в зависимости от того, какая буква следует за данной буквой a . Поскольку этот алгоритм порождает строку, “предшествующую” строке f_n , то отсюда заключаем, что морфизм Фибоначчи (3.4) является кодом.

Отметив, что подстрока a^3 не может появиться в строках Фибоначчи (см. упражнение 3.4.9), мы можем сформулировать следующее утверждение.

Теорема 3.4.5. Квадрат u^2 может входить в строку f_∞ только в том случае, если строка u является циклическим сдвигом строки f_n при некотором целом $n \geq 0$.

Доказательство. Мы уже видели, что для любого циклического сдвига $R_j(f_n)$ строки f_n , где $j \in 0..f_n - 1$, строка u^2 обязательно входит в строку f_∞ . Это доказывает достаточность утверждения теоремы.

Для доказательства необходимости утверждения теоремы предположим существование в строке f_∞ квадрата u^2 , где строка u не будет сдвигом ни какой строки f_n . Рассмотрим минимальную из таких строк u . Пусть n — наименьшее целое, такое, что u является собственной подстрокой строки f_n . Далее покажем, что при этих предположениях строка $h^{-1}(f_n) = f_{n-1}$ содержит квадрат v^2 , где $|v| < |u|$ и строка v не является сдвигом ни какой строки Фибоначчи. Это противоречие докажет утверждение теоремы.

Возможны четыре варианта строки u .

- $u = b \cdots b$. Этот вариант не реализуем, поскольку в этом случае строка u^2 будет содержать подстроку bb .
- $u = a \cdots b$. Этот вариант также не реализуем, поскольку в данном случае строка u обязана иметь вид $u = a \cdots ab$, и тогда для строки $v = h^{-1}(u)$ будет выполняться неравенство $|v| < |u|$.
- $u = a \cdots a$. В этом случае $u^2 = a \cdots aa \cdots a$, и поскольку подстрока a^3 невозможна, имеем $u^2 = a \cdots baab \cdots ba$. Так как и подстрока bb невозможна, то с неизбежностью $u = aba \cdots aba$. Тогда для строки $v = h^{-1}(u)$ будет выполняться неравенство $|v| < |u|$.
- $u = b \cdots a$. Поскольку подстрока bb невозможна, то строка u должна располагаться после буквы a . В этом случае существует подстрока

$$ab \cdots \lambda ab \cdots \lambda a = (ab \cdots \lambda)^2 a$$

некоторой буквы $\lambda \in \{a, b\}$. Таким образом, этот вариант сводится к варианту $u = a \cdots a$ или $u = a \cdots b$.

Теорема доказана. ■

Из этой теоремы сразу получаем, что ни один циклический сдвиг любой строки Фибоначчи не может входить в f_∞ в четвертой степени, поскольку

$$(R_j(f_n))^4 = ((R_j(f_n))^2)^2,$$

но, как мы видели выше, квадрат $(R_j(f_n))^2$ не может быть строкой Фибоначчи.

Отметим, что строка f_7 содержит куб $(aba)^3$. Поэтому имеет место следующая теорема.

Теорема 3.4.6. Все строки Фибоначчи f_n при $n \geq 7$ содержат кубы, но ни одна строка Фибоначчи не содержит кратных подстрок четвертого порядка. ■

Мы показали в разделах 3.2 и 3.3, что если f_∞ содержит кратные подстроки четвертого порядка, тогда такие же подстроки содержат и строки f_n . Но по теореме 3.4.6 строки f_n не могут содержать кратные подстроки четвертого порядка, поэтому справедлива следующая теорема.

Теорема 3.4.7. Бесконечная строка Фибоначчи f_∞ является решением задачи построения $(2, 4)$ -исключения. ■

Теперь мы располагаем достаточными данными для получения оценки числа кратных подстрок, содержащихся в строке f_n . Чтобы получить приближенную нижнюю границу для этого числа, заметим, что для любого целого $k \in 3..n - 4$ количество кратных подстрок определяется вхождением строк f_k в строку f_n . Однако при этом в конечную сумму не включаются квадраты a^2 и $(ab)^2$, что, как покажем ниже, не влияет на асимптотику оценок.

В упражнении 3.4.15 показано, что буква bf_{n-2} раз входит в строку f_n . Тогда из леммы 3.4.2 следует, что строка $f_k f_{n-k-2}$ раз входит в f_n . Поскольку только одно из этих вхождений является суффиксом строки f_n , значит, не менее $f_{n-k-2} - 1$ вхождений строки f_k порождают кубы f_k^3 . Далее отметим, что каждый куб порождает по крайней мере f_k различных квадратов $(R_j(f_k))^2$ ($j = 0, 1, \dots, f_k - 1$), формируемых из циклических сдвигов строки f_k . Следовательно, общее число кратных подстрок, определяемых $f_{n-k-2} - 1$ вхождением кубов f_k^3 в строку f_n , не меньше величины $f_k(f_{n-k-2} - 1)$.

Ранее мы уже показали, что $f_n \in \Theta(\phi^n)$, поэтому $f_k f_{n-k-2} \in \Theta(\phi^k \phi^{n-k-2})$ и, независимо от значения k ,

$$f_k(f_{n-k-2} - 1) \in \Theta(f_{n-2}). \quad (3.30)$$

Поскольку оценка (3.30) справедлива для любого $k \in 3..n - 4$, то отсюда мы делаем заключение, что в строке f_n содержится не менее $\Theta((n - 2)f_{n-2})$ кратных подстрок. Так как $n \approx \log_\phi f_n$ и $f_{n-2} \approx f_n/\phi^2$, приходим к следующей теореме.

Теорема 3.4.8. Количество $q(n)$ кратных подстрок в строке f_n имеет порядок $\Omega(f_n \log f_n)$. ■

Поскольку, как упоминалось ранее, существуют алгоритмы (см. раздел 12.1), вычисляющие все кратные подстроки в любой строке длиной n за время $\Theta(n \log n)$, тогда нижняя граница из теоремы 3.4.8 будет также верхней границей и, кроме того, даст асимптотически наибольшее число кратных подстрок, которые могут содержаться в любой строке длиной n .

В начале этого раздела уже говорилось о том, что несмотря на доказанную оценку из теоремы 3.4.8, “действительное” число кратных подстрок в строках

f_n имеет линейный порядок $\Theta(f_n)$. Это не противоречие, если вспомнить обсуждения 1.3.2 и 2.3.1: кодирование информации — это ключ к эффективности строковых алгоритмов.

Вернемся к рассуждениям, приведенным выше. Очевидно, что описать все кратные подстроки способна кодировка в виде тройки (i, p^*, r^*) , введенной в разделе 2.3: для любой позиции i , с которой начинается квадрат или куб подстроки $R_j(f_k)$, присутствует одна кратная строка — $(i, f_k, 2)$ или $(i, f_k, 3)$. Но, как показано в наших рассуждениях выше, каждый куб f_k^3 порождает дополнительно $f_k - 1$ кратных строк, которые, впрочем, нет необходимости описывать точно. Рассмотрим, например, кратную подстроку $f_4^3 = (abaab)^3$, которая начинается с позиции 9 в строке

$$f_8 = abaababaabaababaababaababaababaababaab.$$

Этот куб и все кратные строки, входящие в его циклические сдвиги, можно описать с помощью *серий*, кодируемых “четверками”, введенными в разделе 2.3:

$$(i, p^*, r^*, t) = (9, 5, 3, 1),$$

где конечный элемент $t = 1$ показывает общее количество циклических сдвигов строки f_5 , исключая саму строку f_5 , которая является кубом. В данном случае четверка $(9, 5, 3, 1)$ определяет кубы $(abaab)^3$ и $(baaba)^3$ и квадраты $(abaab)^2$, $(ababa)^2$ и $(baba)^2$. Предлагаю доказать читателю, что с помощью четверок (i, p^*, r^*, t) можно описать серии кубов и квадратов в любой строке f_n .

Итак, для определения всех вхождений кратных подстрок строки f_k в строку f_n ($k = 1, 2, \dots, n - 2$) необходимо получить значения только одной четверки (i, p^*, r^*, t) . Как мы уже видели, строка f_k входит в строку f_n f_{n-k-2} раз. Поэтому для строки f_n необходимо определить следующее количество четверок:

$$q'(n) = \sum_{k=1}^{n-2} f_{n-k-2} = \sum_{k=0}^{n-3} f_k.$$

По индукции легко показать, что $\sum_{k=0}^{n-3} f_k < f_{n-1}$. Это доказывает такую теорему.

Теорема 3.4.9. Количество $q'(n)$ серий кратных подстрок в строке f_n не превышает f_{n-1} . ■

Мы снова вернулись к нашей интересной ситуации. Теорема 3.4.8 показывает, что в строках Фибоначчи очень много кратных строк. Несмотря на этот факт, теорема 3.4.9 говорит, что если эти кратные строки “запаковать” в серии, то количество серий будет иметь линейный порядок относительно длины строки. Здесь встает закономерный вопрос для любопытных: для любой ли строковой последовательности x можно представить все ее кратные подстроки в виде линейного (по объему) множества серий? Мы вернемся к этому интересному вопросу в главе 12.

Отметим, что *обе* кодировки (i, p^*, r^*) и (i, p^*, r^*, t) , определяющие кратные подстроки путем указания их позиции i в строке, требуют представления строки в виде массива, как сказано в обсуждении 1.2.1.

В заключение этого раздела рассмотрим применение отображения (3.24) для точной локализации квадратов f_n^2 в строке f_∞ для любого целого $n \geq 3$. Покажем, что для этого можно построить алгоритм, который вычисляет все кратные подстроки в строках Фибоначчи за линейное время. Положим в формулах (3.24) $x = f_{n+1}$, $y = f_n$ для любого неотрицательного целого n . Теперь x и y являются строками на алфавите, состоящем из чисел Фибоначчи. Рассмотрим строку $f_\infty(f_{n+1}, f_n)$, например, для $n = 3$:

$$f_\infty(f_4, f_3) = 5355353553553553553 \cdot \dots \quad (3.31)$$

Обозначим через $\Sigma_{i,n}$ сумму первых i значений в строке $f_\infty(f_{n+1}, f_n)$. Очевидно, что $\Sigma_{0,n} = 0$ для любого n , а для примера (3.31) имеем

$$\Sigma_{1,3} = 5, \quad \Sigma_{2,3} = 8, \quad \Sigma_{5,3} = 21, \quad \Sigma_{9,3} = 46.$$

Используя эти обозначения, можно указать позиции в строке f_∞ , где находятся строки f_{n+1} и f_n^2 . Сначала докажем следующую теорему.

Теорема 3.4.10. Для любого целого $i \geq 0$

- а) $f_{n+1} = f_\infty[\Sigma_{i,n} + 1.. \Sigma_{i,n} + f_{n+1}]$, $n \geq 2$,
- б) $f_n^2 = f_\infty[\Sigma_{i,n} + 1.. \Sigma_{i,n} + 2f_n]$, $n \geq 3$.

Доказательство. Поскольку строка f_∞ свободна от квадратов b^2 , тогда любое вхождение строки f_n в выражение (3.27) следует за вхождением строки $f_{n+1} = f_{n-1}f_{n-2}f_{n-1}$, т.е. строка f_{n+1} присутствует при каждом вхождении строки f_n в строку f_∞ . Поэтому утверждение а) теоремы следует непосредственно из леммы 3.4.2 и равенства (3.27). Вследствие леммы 3.4.3 и того факта, что при $n \geq 3$ строка f_{n-2} является префиксом как строки f_{n+1} , так и строки f_n , делаем заключение, что квадрат f_n^2 встретится в тех же позициях, что и строки f_n и f_{n+1} . Это доказывает утверждение б) теоремы. ■

Некоторым “недостатком” последней теоремы можно считать условие $n \geq 3$ в утверждении б), что не позволяет локализовать квадраты $f_1^2 = a^2$ и $f_2^2 = (ab)^2$. Однако следующий результат показывает, что эти квадраты также учитываются.

Лемма 3.4.11. Для $n = 1, 2$ квадраты f_n^2 входят в строку f_∞ только как подстроки строк f_{n+2}^2 .

Доказательство этой леммы предложено дать в упражнении 3.4.16. ■

С учетом последней леммы теорема 3.4.10 позволяет эффективно локализовать все вхождения квадратов f_n^2 в строку f_∞ для любого $n \geq 1$. Как мы видели, вхождения строк f_n в выражение (3.27) порождает кубы f_n^3 и, следовательно, квадраты и кубы сдвигов $R_j(f_k)$ ($j = 0..n-1, n \geq 3$). Использование для серий кодировки (i, p^*, r^*, t) позволяет вычислить все кратные подстроки в строке f_n , причем за линейное время. Алгоритм, реализующий этот подход, описан в [123]. В работе [89] приведены точные формулы для количества квадратов и количества *различных* квадратов в строках Фибоначчи. Отметим работу [139], в которой доказан несколько неожиданный результат, что количество серий (см. обсуждение 2.3.3) в строках Фибоначчи меньше количества различных квадратов в этих строках.

Обсуждение 3.4.1. Строки Фибоначчи обладают многими замечательными свойствами, но, наверное, еще более замечательным фактом является то, что большинство важных свойств этих строк можно обобщить с помощью строк Штурма.

Строкой Штурма (Sturmian string) называется бесконечная строковая последовательность, в которой для любого целого $k \geq 0$ существует $k+1$ различных подстрок длиной k . Строки Штурма строятся на бинарном алфавите. Любой конечный префикс бесконечной строки Штурма также называется строкой Штурма. Значение $k+1$ часто называют *сложностью* строки. Бесконечная строка x имеет сложность k только тогда, когда она представима в виде $x[1..k]^\infty$. Поэтому строки Штурма являются строками наименьшей возможной сложности.

Строки Штурма представляют большой интерес. Поскольку они имеют малую сложность, то естественно предположить, что они содержат много кратных подстрок. И это действительно так [92]: все конечные строки Штурма длиной n содержат $\Theta(n \log n)$ кратных подстрок, которые можно представить в виде $\Theta(n)$ серий, вычисляемых за время порядка $\Theta(n)$ (как и в случае строк Фибоначчи). Более того, для любого целого $r \geq 4$ существует строка Штурма, которая является решением задачи построения $(2, r)$ -исключения, точно так же как строки Фибоначчи являются решением задачи построения $(2, 4)$ -исключений.

Подобно строкам Фибоначчи, строки Штурма можно определить с помощью морфизма, при использовании которого можно установить различные свойства этих строк [92]. Этот морфизм также позволяет распознать за линейное время, является ли данная строка строкой Штурма [37]. В недавней работе [91] показано, что эти свойства распространяются на обобщение строк Штурма, называемых *двух-паттерными*.

Строки Штурма можно определить разными способами, которые показывают их с разных сторон; им посвящено много работ [9, 34, 33, 162, 166, 180, 178, 215]. Существуют другие классы строк с малой сложностью, изучаемые в работах [179, 10, 197]. ■

Упражнения 3.4

1. Мы упоминали строки Фибоначчи в двух предыдущих главах. Какие свойства данных строк использовались в примерах этих глав?
2. Дайте другое определение строк f_n для $n \geq 2$ на основе понятий “суффикс” и “префикс”. Затем докажите лемму 3.4.2.
3. Предложите рекурсивный алгоритм вычисления первых 20 строк Фибоначчи, построенный на основе соотношений (3.18). Затем исключите рекурсию и напишите “обычный” итерационный алгоритм. Какой из этих алгоритмов проще для понимания?
4. Методом математической индукции докажите, что $f_n \leq \phi^n$. Отсюда из неравенства (3.23) будет следовать, что f_n имеет порядок $\Theta(\phi^n)$.
5. Докажите точную формулу $f_n = (\phi^{n+1} - (1 - \phi)^{n+1}) / \sqrt{5}$. Эту формулу также можно переписать в виде $f_n = [\phi^{n+1} / \sqrt{5}]$, где $[x]$ обозначает целую часть числа x .
6. **Обратные строки Фибоначчи** определяются с помощью соотношений

$$f'_0 = b, \quad f'_1 = a, \quad f'_n = f'_{n-2} f'_{n-1}, \quad \forall n \geq 2.$$

Покажите, что морфизм $h' : a \rightarrow ba, b \rightarrow a$ порождает обратные строки Фибоначчи при начальной строке $f'_0 = b$.

7. Докажите, что $h(f_\infty) = f_\infty$.
8. Покажите, что строка f_n не содержит собственных префиксов b , префиксов aa и суффиксов aa .
9. Покажите, что строка f_n не содержит подстрок $b^2, a^3, (ab)^3$ и $(ba)^3$.
10. Предложите алгоритм, реализующий “обратный морфизм” (3.29) для строк f_n .
11. а) Докажите лемму 3.4.3.
 б) Докажите лемму 3.4.4. Можно ли доказать более сильный результат (подобный лемме 3.2.3), что других граней, кроме указанных в условии леммы, в строках f_n нет?
 в) Используя леммы 3.4.3 и 3.4.4, докажите утверждение, сделанное в тексте, что любое вхождение строки f_{n+1} (соответственно, f_n) в строку f_∞ порождает квадрат f_n^2 (соответственно, f_n^3).
12. Покажите, что любое вхождение куба f_n^3 в строку f_∞ порождает квадраты $(R_j(f_k))^2$ для любого целого $j = 1, \dots, f_k - 1$.
13. Докажите, что $f_{n-1} f_n \neq f_n f_{n-1}$.

14. В тексте мы “бездоказательно” утверждали, что нет строк Фибоначчи, которые были бы кратными строками. Докажите или опровергните это утверждение.
15. а) Докажите, что в строке f_n содержится f_{n-1} букв a и f_{n-2} букв b .
 б) Докажите, что f_n будет четным числом только тогда, когда $n = 3k + 2$ для некоторого неотрицательного целого k .
 в) Используя утверждения а) и б), предложите другое доказательство предыдущего упражнения.
16. Докажите лемму 3.4.11.
17. Рассуждения, ведущие к теореме 3.4.8, можно значительно сократить, если установить более точную нижнюю границу для числа кратных подстрок, присутствующих в строке f_n . Можете ли вы получить *точную* формулу для числа кратных подстрок в строке f_n ? (Попробуйте использовать для этого лемму 3.4.11!)
18. Аналогично предыдущему упражнению, попробуйте усовершенствовать доказательство теоремы 3.4.9 путем вывода точной формулы для числа серий кратных подстрок, присутствующих в строке f_n .
19. Чтобы локализовать все кратные подстроки в строке f_∞ , предложите алгоритм, который бы за фиксированное время находил значение $f_\infty[i]$ для любого заданного целого $i > 0$. Можете ли вы предложить формулу для вычисления этого значения?
20. Покажите, как можно построить свободную от кратных подстрок строку длиной 2^k на алфавите из $k + 1$ букв.
21. Покажите, что “морфизм Штурма” $h' : a \rightarrow aab, b \rightarrow ab$ порождает строки $f'_n = R_{n-1}(f_n)$, n нечетно, для начальной строки $f'_1 = a$.
22. Покажите, что если для некоторого целого $k > 0$ строка x имеет сложность k , тогда $x = x[1..k]^r$ для некоторого целого $r \geq 2$.

ГЛАВА 4

Строковые алгоритмы и тестовые данные

Используйте мягкие слова и твердые аргументы.

— Аноним

В данной книге основное внимание уделяется алгоритмам над строковыми последовательностями. Поэтому естественно перед обсуждением специальных строковых алгоритмов рассмотреть некоторые наиболее общие свойства таких алгоритмов. Этой теме посвящен раздел 4.1. Одним из обсуждаемых свойств алгоритмов является их корректность, что обычно проверяется на определенных тестовых данных. Специальные строковые последовательности, введенные в главе 3, будут прекрасными тестовыми данными для некоторых строковых алгоритмов, однако также необходимы тестовые данные и более общего вида. В разделах 4.2 и 4.3 показаны практичные и эффективные методы генерирования множеств тестовых данных, эти методы разработаны специально для доказательного тестирования определенных строковых алгоритмов. Некоторые из этих методов генерирования тестовых данных реализованы в экспериментальных пакетах программного обеспечения [161].

4.1 Хорошие строковые алгоритмы

В этом разделе мы рассмотрим следующие шесть свойств строковых алгоритмов, которые определяют их практическую полезность.

- Корректность.
- Скорость вычислений и требуемая память (временная и пространственная сложность алгоритма).
- Отсутствие возврата к ранее просмотренным данным.
- Независимость от алфавита.
- Возможность предварительной подготовки данных.
- Кодирование выходного результата.

Корректность

Нет сомнения, что корректность является необходимым свойством любого алгоритма, — некорректный алгоритм бесполезен. В худшем случае некорректный алгоритм, реализованный в компьютерной программе, может быть просто опасным, если, например, он управляет тормозной системой автомобиля или полетом самолета. Доказательство корректности программных продуктов — даже если есть четкое определение понятия корректности — большая нерешенная проблема “инженерии программного обеспечения”. Для математических алгоритмов, представленных в книге, корректность доказывается относительно легко по сравнению с доказательствами корректности для многих других типов алгоритмов, но это не означает, что эти доказательства совсем простые. Несколько лет назад один научный работник реализовал на языке C четыре ранее опубликованных алгоритма для вычисления канонических форм петель (см. определение 1.4.3); среди этих алгоритмов был один мой. Обнаружилось, что три из этих алгоритмов, включая и мой, оказались некорректными!

Доказательство корректности алгоритма проходит два основных этапа. На первом этапе предлагается математическое “доказательство” корректности; на втором этапе выполняется проверка алгоритма на тестовых входных данных, для которых известен правильный выходной результат. Доказательства корректности, выполненные на каждом из этих этапов, ненадежны. Математические доказательства ненадежны из-за “человеческого фактора”, а выполнение алгоритма на тестовых данных не может доказать “полную” корректность, поскольку тестовые данные никогда не бывают исчерпывающими. Однако лучше иметь ненадежные средства, чем не иметь средств вообще. В разделах 4.2 и 4.3 я покажу, как можно генерировать полные тестовые данные, подходящие для многих строковых алгоритмов, в виде строковых последовательностей с широким диапазоном их длин и при этом с невысокой стоимостью самого процесса генерирования данных.

Временная и пространственная сложность

Временная и пространственная сложность — основные характеристики любого алгоритма, и, конечно, алгоритм, время выполнения которого в наихудшем случае

имеет порядок $\Theta(n)$ на строках длиной n , предпочтительнее алгоритма с временем выполнения в наихудшем случае порядка $O(n^2)$. Но не все так очевидно и ясно с этими характеристиками. Например, простой алгоритм 2.2.1, вычисляющий вхождения частного паттерна длиной m в строку длиной n и описанный в разделе 2.2, имеет время выполнения, не лучшее, чем у какого-нибудь плохого алгоритма. Однако для этого алгоритма временная граница $O(mn)$ достигается только на одном патологическом варианте строковых последовательностей, который не встречается на практике, тогда как на “практических” строках время выполнения имеет порядок $O(m + n)$. Таким образом, этот алгоритм фактически является “хорошим” алгоритмом, поскольку он прост и в среднем выполняется быстрее, чем алгоритм КМП (см. раздел 7.1). Кроме того, он элегантно обходит возможные возвраты к просмотренным ранее позициям в строке для их повторной проверки.

Пример вычисления деревьев суффиксов (см. раздел 2.1), требующий времени порядка $\Theta(n \ln \alpha)$ на алфавите объемом α и памяти порядка $\Theta(n\alpha)$, показывает, что “хороший” порядок относительно n не гарантирует быстрого выполнения или малого требуемого объема памяти, поскольку константы пропорциональности в асимптотических оценках могут быть чрезмерно большими, особенно для больших алфавитов. Аналогично тщательно разработанные алгоритмы сравнения с паттернами, где сведено к теоретическому минимуму количество парных буквенных сравнений [97, 52], почти не используются на практике, поскольку они чрезвычайно сложны, вследствие чего усложняется процесс вычисления, что, в свою очередь, приводит к увеличению значений констант пропорциональности в их асимптотических оценках.

Если сравнивать временную и пространственную сложности алгоритмов, то временная сложность считается более важным критерием при выборе алгоритма. Это отображает тот факт, что стоимость хранения единицы данных обычно ниже стоимости единицы времени процесса вычислений. Таким образом, временная сложность алгоритма часто “перевешивает” все другие характеристики алгоритма.

Отсутствие возврата к ранее просмотренным данным

Просмотр строковых последовательностей естественно выполнять слева направо. Поэтому чрезвычайно привлекательным свойством строковых алгоритмов считается их *онлайновость*. Это означает, что если в процессе вычислений уже просмотрена позиция i в строковой последовательности, то алгоритм уже не вернется к этой позиции i снова. Другими словами, онлайн-алгоритм никогда *не выполняет возврата* к уже просмотренным позициям.

Существует также более слабое определение онлайн-алгоритмов. Например, алгоритм называется онлайн-алгоритмом с *окном* размера k , если существует такая положительная константа k , что для выполнения вычислений в текущей позиции i требуется не более k строковых позиций из интервала $i - k + 1..i$. Еще одно по-

добное слабое определение позволяет считать алгоритм онлайнным, если значение, вычисляемое алгоритмом для любой позиции i , никогда не пересчитывается, даже если потребуется заново обратиться к позициям $i' < i$.

Из любого сделанного выше определения онлайнности следует, что строковый онлайнный алгоритм при решении определенной задачи для строки $x = x[1..n]$ эффективно решает ту же задачу для каждого префикса $x[1..i]$ этой строки. Таким образом, можно сказать, что онлайнный алгоритм фактически решает n задач по цене одной задачи.

Онлайнным алгоритмом является алгоритм КМП, упоминавшийся выше. Отметим, что алгоритм 1.3.1 вычисления массива граней также является онлайнным (в слабом определении этого понятия): массив граней $\beta[1..i]$ для подстроки $x[1..i]$ вычисляется после просмотра значения $x[i]$, и далее этот массив не пересчитывается.

Независимость от алфавита

В разделе 1.2 показано, что строковые последовательности в общем случае могут быть определены на произвольных множествах элементов. Поэтому желательно, чтобы строковые алгоритмы эффективно выполнялись на строках, определенных на любых алфавитах, т.е. чтобы они были независимы от алфавита. Таким алгоритмом является, например, алгоритм 1.3.1 вычисления массива граней — он имеет асимптотически оптимальное линейное время выполнения при условии, что сравнение (равны или не равны) двух любых букв алфавита выполняется за константное время, и при этом алгоритм не использует никакие другие свойства алфавита.

С другой стороны, часто можно уменьшить временную сложность алгоритма, если использовать некоторые особые свойства алфавита. Поэтому если данное свойство алфавита часто встречается на практике, то разрабатываются специальные версии алгоритма, которые оказываются более эффективными (по сравнению с исходным алгоритмом) на строковых последовательностях, определенных на таких алфавитах. Чтобы более точно описать эту проблему, рассмотрим следующую простую задачу.

Задача 4.1 (Подсчет вхождений). Подсчитать количество вхождений каждой буквы в данную строку x длиной n . ■

Рассмотрим эту задачу для трех различных алфавитов A .

- **Произвольный (неупорядоченный) алфавит:** для произвольных букв $\lambda, \mu \in A$ равенство $\lambda = \mu$ проверяется за фиксированное время (примером такого алфавита может служить множество китайских иероглифов).
- **Упорядоченный алфавит:** для произвольных букв $\lambda, \mu \in A$ некоторое определенное отношение упорядочения $\lambda < \mu$ проверяется за фиксированное

время (примером такого алфавита может служить множество действительных чисел или множество слов английского языка).

- **Индексированный алфавит:** определен массив T , такой, что для произвольной буквы $\lambda \in A$ значение $T[\lambda]$ можно получить за константное время (примером такого алфавита может служить подмножество символов ASCII).

Отметим, что индексированный алфавит A обязательно должен быть конечной мощности $\alpha = |A|$. Тогда на индексированном алфавите задачу 4.1 можно решить непосредственным подсчетом: сначала целочисленный массив $T[1..n]$ полагается равным нулю, затем для каждой позиции $i \in 1..n$ входной строки x значение $T[x[i]]$ должно увеличиваться на единицу. Имея дополнительный список кодировки букв, ненулевые элементы массива T можно вывести за время, пропорциональное количеству различных букв в строке x . Таким образом, для решения задачи требуется время порядка $\Theta(\alpha + n)$.

В случае упорядоченного алфавита можно построить сбалансированное дерево поиска (такое, как AVL-дерево¹ или 2-3-дерево²) путем просмотра входной строки x слева направо: каждый узел дерева идентифицируется парой (λ, c) , где λ — буква из строки x , а c — счетчик, который принимает значение 1 при создании узла и затем увеличивается на единицу при нахождении каждого нового вхождения буквы λ в строку x . Время, необходимое для решения задачи 4.1, в данном случае имеет порядок $O(n \log n)$. Отметим, что если объем алфавита α значительно больше величины n , то такой алгоритм может выполняться быстрее, чем соответствующий алгоритм для строк, определенных на индексированном алфавите.

В случае неупорядоченного алфавита необходимо выполнить $\Theta(n)$ шагов, на каждом из которых следует выполнить порядка $\Theta(n)$ попарных сравнений букв. Поэтому решение задачи 4.1 в данном случае требует $O(n^2)$ времени. Однако нетрудно заметить, что решение может быть найдено проще и быстрее, если алфавит будет, например, двоичным. В этой связи отметим, что в цифровых моделях компьютерных вычислений буквы неупорядоченного алфавита можно закодировать в виде бинарных строк, для которых применимы отношения порядка. Таким образом, на практике всегда можно считать, что алфавит упорядочен.

На этой простой (и типичной) задаче показано, что выбор подходящего алгоритма очень сильно зависит от алфавита, которым определяется исследуемая строковая последовательность. Обсуждение интересного вопроса зависимости ал-

¹AVL-деревом (названным так в честь его авторов Г. М. Адельсона-Вельского и Е. М. Ландиса) называется сбалансированное по высоте дерево двоичного поиска, у которого для каждого узла высоты его правого и левого поддеревьев отличаются не более чем на единицу. — *Примеч. ред.*

²2-3-деревом называется дерево поиска специального вида, когда из каждого внутреннего узла выходит от двух до трех ребер и все пути от корня дерева до любого листа имеют одинаковую длину. Кроме того, во внутренние узлы обычно помещаются определенные записи с информацией об этих узлах или путях от корня до узла. — *Примеч. ред.*

горитма от алфавита на этом не заканчивается; оно будет продолжено в других разделах книги.

Необходимость предварительной подготовки данных

Как мы видели в разделе 2.1, вычисление дерева суффиксов можно рассматривать как процесс предварительной подготовки данной строки x : это дерево затем используется для многократных поисков в строке x , при этом для каждого такого поиска его время пропорционально длине искомого паттерна. Аналогично в алгоритме КМП (см. раздел 7.1) на этапе предварительной подготовки вычисляется массив граней паттерна, который затем используется в процессе поиска, позволяя исключить обратный просмотр. Таким образом, возможность алгоритма использовать предварительную подготовку данных можно считать дополнительным положительным качеством этого алгоритма.

Продемонстрируем преимущества предварительной подготовки на примере задачи 4.1, сформулированной в предыдущем подразделе. Предположим, что алфавит A фиксирован (что известно заранее) и конечен. Также предположим, что необходимо подсчитать вхождения букв не для одной строки, а для $m \geq 1$ строк, определенных на этом алфавите A . Для простоты и без потери общности можно считать, что все m строк имеют одинаковую длину n . Тогда на этапе предварительной подготовки за время $\Theta(\alpha)$ алфавит можно закодировать в виде хеш-таблицы T так, чтобы отдельные буквы были доступны за фиксированное или почти константное время. Значения хеш-таблицы $T[\lambda]$ для каждой буквы λ должны содержать счетчик строк $j \in 0..m$, показывающий номер строки, в которой последний раз встречалась буква λ . Первоначально счетчик j устанавливается равным нулю для всех λ . Значения $T[\lambda]$ также должны содержать счетчик вхождений c . После завершения этапа предварительной подготовки начинается процесс просмотра m строк. Для каждой позиции i в j -й строке x_j вычисляется $T[x_j[i]]$; если значение строкового счетчика равно j , тогда значение счетчика вхождений c увеличивается на единицу, иначе строковый счетчик устанавливается равным j , а счетчик c — равным 1. Среднее время этапа предварительной подготовки должно “раскладываться” на все m строк. Поэтому время обработки одной строки составляет $\Theta(\alpha/m + n)$ независимо от используемого алфавита. Очевидно, такой алгоритм, использующий предварительную подготовку, во многих случаях предпочтительнее других рассмотренных выше алгоритмов для решения задачи 4.1, поскольку предполагает, что алфавит обработан заранее.

Отметим, что все рассмотренные алгоритмы для подсчета количества вхождений букв в заданные строки являются онлайн-овыми.

Кодирование выходного результата

Как мы видели в разделе 2.3, при вычислении кратных строк (задача 2.15) ключевую роль играет кодирование (способ представления) выходного результа-

та, которое позволяет уменьшить объем выходных данных, а также сложность алгоритма от $O(n^2)$ до $O(n \log n)$. Например, для отображения в выходных результатах кратной строки a^6 наиболее эффективно в нормальной форме $(i, 1, 6)$, при этом предполагается, что вхождения кратных подстрок a^2 , a^3 и a^4 также определены.

В общем случае, как показано в обсуждениях 1.3.2 и 2.3, неясно, на основании какого критерия принимать решение об использовании кодирования выходных данных. Рассмотрим, например, массив граней β для некоторой строки x . Обычно считается, что массив граней является подходящим способом кодирования (представления) граней. Однако, как следует из леммы 1.3.1, для того чтобы действительно иметь список всех граней, необходимо вычислить $\beta[n], \beta^2[n], \dots$. Напомним: из леммы 3.2.3 следует, что строки Туе, являющиеся $(2, 3)$ -исключениями, а также строки Фибоначчи имеют количество граней, пропорциональное логарифму от длин строк. Поэтому в общем случае требуется $\Theta(\log n)$ времени (или даже больше) для формирования списка граней строки x на основании вычисленного массива граней β . Скептики правы, когда утверждают, что, вычислив массив граней, мы в действительности не имеем решения задачи, поскольку требуются еще дополнительные значительные вычислительные усилия, чтобы получить окончательное решение задачи. В пользу массивов граней приводятся такие практические соображения: как правило, основной интерес представляют грани наибольшей длины (которые легко получить из массивов граней), а грани небольшой длины рассматриваются только после того, как получены грани наибольшей длины.

Упражнения 4.1

1. Разработайте в деталях четыре алгоритма вычисления количества вхождений букв в заданные строки, кратко описанные в этом разделе.
2. Оцените временную сложность всех четырех алгоритмов, описанных в данном разделе. Определите классы строковых последовательностей, для которых тот или иной алгоритм предпочтительнее других.
3. Покажите, что буквы индексированного алфавита A объемом α можно однозначно отобразить за линейное время в множество натуральных чисел $N_\alpha = \{1, 2, \dots, \alpha\}$. Отсюда будет следовать утверждение, что любой индексированный алфавит без потери общности можно считать множеством натуральных чисел N_α .
4. В тексте раздела утверждалось (но не доказывалось), что строки Туе, являющиеся $(2, 3)$ -исключениями, имеют $\Theta(\log n)$ граней, а строки Фибоначчи — $\Omega(\log n)$ граней, где n — длина строки. Докажите эти утверждения, используя леммы 3.2.3 и 3.4.4.

4.2 Разные паттерны

Как сделать одну строку отличной от другой? Это фундаментальный вопрос методов генерирования данных, предназначенных для тестирования строковых алгоритмов: как сгенерировать как можно больше “различных” строк, избегая при этом дублирования не “различных” строк.

Традиционный подход (согласно определениям раздела 1.2) — считать две строки различными только тогда, когда они различны хотя бы в одной позиции. Согласно этому определению, на конечном алфавите A объемом α можно построить в точности α^n различных строк длиной n , тем самым предлагая для тестирования исчерпывающие данные в виде α^n различных строк длиной n на алфавите A .

Однако такой подход не является единственным или наилучшим единственным для многих строковых алгоритмов. Рассмотрим простую задачу подсчета вхождений букв в некоторую заданную строку x , описанную в разделе 4.1. На неупорядоченном алфавите алгоритмы, решающие эту задачу, просматривают каждую букву строки x по отдельности, выполняя операцию сравнения каждой буквы с буквами алфавита. Для таких алгоритмов результат выполнения будет одинаковым независимо от того, будет ли входная строка иметь вид $x = aab$ или $y = ddc$. Нужно, чтобы тестовые данные содержали одну из таких строк, но совсем не обязательно тестировать обе эти строки.

Этот пример порождает идею *эквивалентности относительно паттерна*, где две строки x и y будут считаться эквивалентными, если они соответствуют некоторому “паттерну”. Чтобы сформулировать эту идею более точно, введем такое определение.

Определение 4.2.1. Две строки x и y называются *p -эквивалентными*, если выполняются следующие условия:

- а) для всех целых $n \geq 0$ $|x| = |y| = n$;
- б) для всех целых i и j , $1 \leq i \leq j \leq n$, $x[i] = x[j] \Leftrightarrow y[i] = y[j]$.

Если строки x и y не являются p -эквивалентными, то они называются *p -различными*. ■

В этом разделе, следуя работе [183], мы исследуем свойства p -различных строк. В частности, мы покажем, как подсчитать количество и как вычислить все p -различные строки длиной n . Покажем, что таких строк относительно немного. Далее будет описано, как вычислить такие строки за константное время в расчете на одну строку.

В следующем разделе мы введем понятие эквивалентности относительно граней (или b -эквивалентность), которое еще ближе обычному понятию эквивалентности, чем p -эквивалентность. Это приведет к определению b -различных строк,

которых, как будет показано, еще меньше, чем p -различных строк. Такие строки также можно вычислить за константное время в расчете на одну строку.

Таким образом, p -различные и b -различные строки могут быть тестовыми данными для алгоритмов, которым они подходят, при этом они эффективны (дают исчерпывающий набор тестовых данных) и результативны (их можно быстро сгенерировать).

Первым важным фактом, относящимся к p -эквивалентности, является то, что это отношение действительно является отношением эквивалентности между строковыми последовательностями, и поэтому оно разбивает все множество строк длиной n на классы эквивалентности p -различных строк (см. упражнение 4.2.2). Отсюда следует, что для каждого класса эквивалентности можно определить единственный представитель этого класса. Чтобы увидеть это, введем счетный бесконечный *стандартный алфавит*

$$\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_\alpha, \dots\} \quad (4.1)$$

и подалфавиты $\Lambda_\alpha = \{\lambda_1, \lambda_2, \dots, \lambda_\alpha\}$ для каждого целого $\alpha \geq 1$. Предполагаем, что буквы алфавита Λ упорядочены естественным образом

$$\lambda_1 < \lambda_2 < \dots < \lambda_\alpha < \dots$$

Вследствие такого упорядочения без потери общности можно считать, что стандартный алфавит состоит из натуральных чисел. Тогда для любой заданной строки x на алфавите Λ , определим *p -каноническую строку x^** , соответствующую строке x , как лексикографически наименьшую строку на алфавите Λ , p -эквивалентную строке x . Очевидно, что строка x^* является уникальным представителем своего класса p -эквивалентности. Нетрудно показать, что любая строка x^* удовлетворяет следующему условию.

(P) Для любого целого положительного j , если буква λ_j входит в строку x^* , тогда следующая буква λ_{j+1} не может иметь другого вхождения в эту строку ранее самого первого (самого левого) вхождения буквы λ_j .

Обозначим через $p'[\alpha, n]$ количество p -канонических строк длиной n , сформированных с использованием всех букв подалфавита Λ_α (каждая буква подалфавита должна хотя бы один раз входить в каждую строку). Мы подсчитаем это количество с помощью бесконечного двумерного массива, называемого p' -массивом.

Теорема 4.2.2. Для любых положительных целых α и n

- а) $p'[1, n] = 1$;
- б) если $\alpha > n$, тогда $p'[\alpha, n] = 0$;
- в) $p'[\alpha, \alpha] = 1$;
- г) если $\alpha \geq 2$ и $n \geq 2$, тогда $p'[\alpha, n] = p'[\alpha - 1, n - 1] + \alpha p'[\alpha, n - 1]$.

Доказательство.

- а) Если $\alpha = 1$, тогда имеется только одна p -каноническая строка $\mathbf{x}^* = \lambda_1^n$.
- б) Если $\alpha > n$, тогда на основании свойства (P) можно утверждать, что не существует p -канонической строки, которая содержала бы букву λ_α .
- в) Из свойства (P) следует, что существует только одна p -каноническая строка, которая содержит α различных букв: $\mathbf{x}^* = \lambda_1 \lambda_2 \dots \lambda_\alpha$.
- г) Обозначим через $\pi_1 = p'[\alpha - 1, n - 1]$ количество различных p -канонических строк длиной $n - 1$, которые содержат в точности $\alpha - 1$ букв подалфавита $\Lambda_{\alpha-1}$. Обозначим множество этих строк как $S_1 = \{\mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_{\pi_1}^*\}$. Тогда все строки вида

$$\mathbf{x}_i^* \lambda_\alpha, \tag{4.2}$$

где целое i удовлетворяет неравенствам $1 \leq i \leq \pi_1$, будут различными и p -каноническими.

Подобным образом обозначим через $\pi_2 = p'[\alpha, n - 1]$ количество различных p -канонических строк длиной $n - 1$, которые содержат в точности α букв подалфавита Λ_α . Обозначим множество этих строк как $S_2 = \{\mathbf{y}_1^*, \mathbf{y}_2^*, \dots, \mathbf{y}_{\pi_2}^*\}$. Тогда все α строк вида

$$\{\mathbf{y}_i^* \lambda_1, \mathbf{y}_i^* \lambda_2, \dots, \mathbf{y}_i^* \lambda_\alpha\}, \tag{4.3}$$

где целое i удовлетворяет неравенствам $1 \leq i \leq \pi_2$, будут различными и p -каноническими. Поскольку конечные буквы в этих строках входят в них не менее двух раз, эти строки не совпадают ни с одной строкой вида (4.2). Поэтому справедливо неравенство

$$p'[\alpha, n] \geq p'[\alpha - 1, n - 1] + \alpha p'[\alpha, n - 1].$$

Пусть строка \mathbf{x}^* будет p -канонической длины n , содержащей все буквы подалфавита Λ_α . Представим строку \mathbf{x}^* как $\mathbf{x}^* = \mathbf{y}^* \lambda_i$. Если буква λ_i входит в строку \mathbf{y}^* , тогда $\mathbf{y}^* \in S_2$, и поэтому строка \mathbf{x}^* будет одной из строк вида (4.3). С другой стороны, из свойства (P) следует, что буква λ_α не может входить в строку \mathbf{y}^* . Поэтому при $i = \alpha$ имеем $\mathbf{y}^* \in S_1$, а строка \mathbf{x}^* является одной из строк вида (4.2). Эти рассуждения приводят к неравенству

$$p'[\alpha, n] \leq p'[\alpha - 1, n - 1] + \alpha p'[\alpha, n - 1].$$

Из последнего неравенства и противоположного ему, полученного выше, следует утверждение з) теоремы. ■

Рекуррентное соотношение утверждения *з*) теоремы 4.2.2 хорошо известно в математике: с учетом начальных значений, определяемых утверждениями *а*)–*в*) этой теоремы, это соотношение определяет так называемые числа Стирлинга второго рода $\sigma_n^{(\alpha)}$ [135]. Таким образом, для всех положительных целых чисел n и α имеем

$$p'[\alpha, n] = \sigma_n^{(\alpha)}. \tag{4.4}$$

В табл. 4.1 приведены значения верхнего левого угла p' -массива. Суммы $p[n]$ значений по столбцам равны общему количеству p -канонических строк длиной n . Отметим, что числа $p[n]$ также хорошо известны — они называются числами Белла [211].

Таблица 4.1. Значения p' -массива для строк длиной $n \leq 9$ и $\alpha \leq 9$

n	1	2	3	4	5	6	7	8	9
$\sigma_n^{(1)}$	1	1	1	1	1	1	1	1	1
$\sigma_n^{(2)}$	0	1	3	7	15	31	63	127	255
$\sigma_n^{(3)}$	0	0	1	6	25	90	301	966	2025
$\sigma_n^{(4)}$	0	0	0	1	10	65	350	1701	7770
$\sigma_n^{(5)}$	0	0	0	0	1	15	140	1050	6951
$\sigma_n^{(6)}$	0	0	0	0	0	1	21	266	2646
$\sigma_n^{(7)}$	0	0	0	0	0	0	1	28	462
$\sigma_n^{(8)}$	0	0	0	0	0	0	0	1	36
$\sigma_n^{(9)}$	0	0	0	0	0	0	0	0	1
$p[n]$	1	2	5	15	52	203	877	4140	20147

На примере значений табл. 4.1 можно проиллюстрировать другую зависимость между числами Стирлинга и нашими числами p' . Числа Стирлинга $\sigma_n^{(\alpha)}$ обычно определяют как число способов разбиения множества S из n элементов на α непустых непересекающихся подмножеств, объединение которых составляет множество S . Чтобы увидеть, как это определение соответствует числам $p'[\alpha, n]$, рассмотрим случай, когда $n = 4$ и $\alpha = 2$. Если выписать для этого случая семь строк ($p'[2, 4] = \sigma_4^{(2)} = 7$) и затем образовать $\alpha = 2$ подмножества *позиций* одинаковых букв в этих строках, то увидим, что декомпозиция множества $\{1, 2, 3, 4\}$ на подмножества обладает следующими свойствами.

- Такая декомпозиция единственная (поскольку все строки различные).
- Декомпозиция непустая (поскольку в строки входят все α букв).
- Образованные подмножества не пересекаются (поскольку каждой позиции соответствует только одна буква).

Приведем эту декомпозицию.

1234	
aaab	{1, 2, 3}{4}
aaba	{1, 2, 4}{3}
aabb	{1, 2}{3, 4}
abaa	{1, 3, 4}{2}
abab	{1, 3}{2, 4}
abba	{1, 4}{2, 3}
abbb	{1}{2, 3, 4}

Легко видеть, что объединение в приведенном примере по отдельности левых и правых множеств порождает два непустых и непересекающихся множества. Доказательство этого факта в общем случае оставим для упражнения 4.2.4.

Теорема 4.2.2 предлагает итерационный способ вычисления чисел $p'[\alpha, n]$. Также различные формулы для этих чисел можно получить, используя связь чисел Стирлинга второго рода с биномиальными коэффициентами и с числами Стирлинга первого рода [135]. Отметим, что для любого фиксированного алфавита объемом α сумма $\sum_{i=1}^{\alpha} p'[i, n]$ равняется количеству p -различных строк длиной n , содержащих не более α различных букв. Поскольку для значений n , больших α , почти все такие строки содержат в точности α букв, отсюда следует, что

$$\lim_{n \rightarrow \infty} \left(\sum_{i=1}^{\alpha} p'[i, n] / \frac{\alpha^n}{\alpha!} \right) = 1. \tag{4.5}$$

Если рассматривать различные, в обычном понимании, строки длиной n , сформированные из не более α букв, то таких строк всего α^n . Формула (4.5) показывает, что p -различных строк, также сформированных из не более α букв, асимптотически в $\alpha!$ раз меньше. Например, для $\alpha = 5$ и $n = 9$ имеем $5^9 = 1953125$ обычных различных строк, тогда как, суммируя первые пять значений в 9-м столбце табл. 4.1, получаем 17002 p -различных строк. Отношение 1953125 к 17002 дает значение, близкое к 115, что примерно равно $5!$.

Конечно, наибольший интерес представляет сумма $p[n] = \sum_{i=1}^n p'[i, n]$, которую можно вычислить различными способами, в частности с помощью рекуррентного соотношения

$$p[n] = \sum_{j=0}^{n-1} C_{n-1}^j p[j], \quad n \geq 2 \tag{4.6}$$

или как бесконечную сумму

$$p[n] = e^{-1} \sum_{j \geq 1} \frac{j^{n-1}}{(j-1)!}, \quad (4.7)$$

в которой фактически надо оценить только $\Theta(n)$ слагаемых, поскольку, как легко показать,

$$p[n] - 1 > e^{-1} \sum_{j=1}^{2n} \frac{j^{n-1}}{(j-1)!}. \quad (4.8)$$

Итак, можем подсчитать количество p -различных строк путем вычисления величин $p'[\alpha, n]$ и $p[n]$. Теперь покажем, как генерировать p -канонические строки, которые можно использовать в качестве тестовых данных для строковых алгоритмов. Возвращаясь к доказательству теоремы 4.2.2, видим, что

- добавление буквы λ_α к p -каноническим строкам увеличивает их количество на величину $p'[\alpha - 1, n - 1]$;
- добавление букв $\lambda_1, \lambda_2, \dots, \lambda_\alpha$ к p -каноническим строкам увеличивает их количество на величину $p'[\alpha, n - 1]$.

Это простое замечание дает основу для итерационного алгоритма вычисления всех $p[n]$ p -канонических строк длиной n . Как показано ниже, эти строки можно генерировать за фиксированное время в расчете на одну строку путем построения дерева T_n . Фактически, мы возвращаемся к обсуждению задачи 2.4 из раздела 2.1, поскольку дерево T_n является синтаксическим деревом, но у него помечаются не ребра, а узлы.

Узлы дерева T_n помечаются парой (λ, α) , где λ — буква алфавита Λ , а α — количество различных букв, которые можно найти на пути от корня дерева до текущего узла. Дерево T_1 состоит из одного корневого узла, помеченного как $(\lambda_1, 1)$, а для всех целых $n \geq 2$ деревья T_n формируются из деревьев T_{n-1} путем добавления к листу, конечному узлу с меткой (λ, α) , новых конечных узлов с метками

$$(\lambda_1, \alpha), (\lambda_2, \alpha), \dots, (\lambda_\alpha, \alpha), (\lambda_{\alpha+1}, \alpha + 1).$$

На рис. 4.1 показано полное дерево T_3 , где буквы λ_i заменены на индексы i .

Из теоремы 4.2.2 следует, что дерево T_n имеет ровно $p[n]$ конечных узлов (листьев) и буквы, стоящие на пути от корня до этих узлов, дают все $p[n]$ p -канонические строки x^* длиной n . Таким способом генерируются эти строки на основе построенного дерева T_n . Еще раз отметим, что дерево T_n формируется из дерева T_{n-1} путем добавления $p[n]$ концевых узлов — операция, требующая фиксированного времени для одного узла или всего $\Theta(p[n])$ времени. Из формулы (4.6) следует очевидное неравенство $p[n] \geq 2p[n-1]$, откуда заключаем, что все дерево T_n можно построить за время $\Theta(p[n])$. Это доказывает следующую теорему.

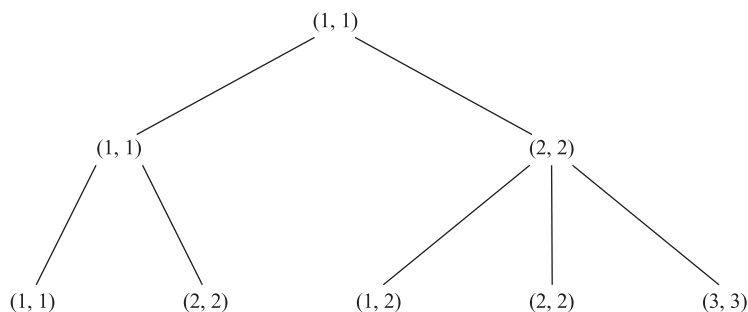


Рис. 4.1. Дерево T_3 , представляющее p -различные строки длиной 3

Теорема 4.2.3. Для любого положительного целого n все p -канонические строки длиной n можно вычислить за время порядка $\Theta(p[n])$, для чего необходима память объемом $\Theta(p[n])$. ■

Отметим, что данный алгоритм генерирования p -канонических строк является онлайнным в слабом значении этого понятия, как обсуждалось в разделе 4.1, — поскольку дерево T_n строится на основе дерева T_{n-1} , для решения задачи для данного n необходимо сначала решить задачу для $n - 1$. Такое онлайнное решение эквивалентно построению дерева T_n “в ширину”. Это означает, что решение легко расширяемо: если возникнет необходимость в дереве T_{n+1} , то оно может быть сформировано тем же алгоритмом за один шаг из дерева T_n . Как мы увидим в следующем разделе, эффективное решение “в ширину” иногда может привести к определенным трудностям.

Результат, подобный теореме 4.2.3, можно установить и для p -канонических строк длиной n , использующих не более α букв алфавита Λ . Напомним, что таких строк всего $p'[\alpha, n]$. В этом случае для генерирования таких строк используется поддереву дерева T_n , в котором пути длиной n заканчиваются в конечных узлах, помеченных метками вида $(*, \alpha)$.

В заключение этого раздела отметим, что дерево T_n можно проходить (просматривать) различными способами, получая при этом различное упорядочение p -канонических строк. Например, при прямом прохождении T_n получаем список строк, упорядоченных в лексикографическом порядке.³ И, соответственно, по-

³При прохождении дерева в прямом порядке сначала посещается корень дерева, затем в таком же порядке посещаются все узлы самого левого поддерева, далее посещаются узлы следующих слева направо поддеревьев. При прохождении дерева в обратном порядке обход начинается с самого правого листа самого правого поддерева данного дерева, затем в таком же порядке посещаются все узлы этого самого правого поддерева, далее узлы следующих справа налево поддеревьев. Последним посещается корень дерева. — *Примеч. ред.*

лучим обратный порядок строк, если проходить дерево T_n в обратном порядке. Поскольку каждую строку можно “исключать” из списка после ее генерирования, для генерирования строк на основе прохождения дерева T_n в любом порядке фактически требуется память объемом порядка $\Theta(n)$. Так как $p[n] \leq 2^n$ (что следует из формулы (4.6)), то этот объем памяти эквивалентен объему $\Theta(\log p[n])$.

Упражнения 4.2

1. Переформулируйте определение 4.2.1 таким образом, чтобы исключить из него требование равенства длин строк x и y .
2. Отношение эквивалентности обладает свойствами рефлексивности, симметричности и транзитивности. Докажите, что отношение p -эквивалентности действительно является отношением эквивалентности.
3. Докажите свойство (P).
4. Используя “обычное” определение чисел Стирлинга, покажите, что равно числу всех p -канонических строк длиной n , включающих ровно α букв.
5. Докажите формулы (4.6)–(4.8).
6. Постройте дерево T_4 и проверьте, что оно имеет $p[4] = 15$ листьев (конечных узлов).
7. Разработайте рекурсивный алгоритм построения дерева T_n для заданного значения n .
8. Покажите, что поддерево дерева T_n , соответствующего $p'[\alpha, n]$ p -каноническим строкам длиной n , содержащим не более α букв алфавита Λ , можно построить за время $O(\alpha p'[\alpha, n])$, используя память объемом $O(\alpha p'[\alpha, n])$.
9. В тексте без доказательства утверждалось, что при прямом прохождении дерева T_n получается список p -канонических строк, упорядоченных в лексикографическом порядке. Докажите или опровергните это утверждение.

4.3 Разные грани

Как мы видели в разделах 1.2 и 1.3, основное преимущество массивов граней заключается в возможности описать периодическую структуру строковых последовательностей. Поэтому не вызывает удивления тот факт, что массивы граней играют важную роль во многих алгоритмах, которые исследуют эту структуру или выходной результат которых зависит от этой структуры. Например, алгоритм КМП (см. раздел 7.1), который находит частные паттерны в строках (задача 2.8), зависит от вычисления массивов граней для этих паттернов. Алгоритмы, вычисляющие все оболочки для данной строки (задача 2.18), также сильно зависят от массивов граней: эффективные алгоритмы требуют первоначальной замены

строки ее массивом граней, тогда последующие вычисления выполняются с этим массивом, а не с самой строкой.

Подготовка тестовых данных для этих и других алгоритмов показывает важность генерирования строк, порождающих различные массивы граней. Введем определение.

Определение 4.3.1. Две строки x и y называются **эквивалентными относительно граней** (b -эквивалентными), если их массивы граней совпадают. В противном случае строки x и y называются b -различными. ■

Чтобы пояснить это определение, рассмотрим следующие строки:

$$x = ababb, \quad x' = ababc, \quad x'' = abaacb.$$

Строки x и x' p -различны, но b -эквивалентны, поскольку имеют одинаковый массив граней $\beta = 00120$. Таким образом, мы видим, что из b -эквивалентности не вытекает p -эквивалентность. С другой стороны, строки x и x'' b -различны, поскольку строка x'' имеет массив граней 00100 , но они также и p -различны. Очевидно, что из b -различия строк с необходимостью вытекает их p -различие. Это простое наблюдение сформулируем в виде леммы.

Лемма 4.3.2. Если две строки p -эквивалентны, то они также b -эквивалентны.

Доказательство этой леммы предложено дать в упражнении 4.3.1. ■

Как следует из этой леммы, каждый класс b -эквивалентных строк состоит из одного или нескольких классов p -эквивалентных строк. Следовательно, количество классов b -эквивалентных строк меньше количества классов p -эквивалентных строк. Поскольку p -эквивалентность является отношением эквивалентности, мы снова можем определить b -канонические представители x^* для каждого класса эквивалентности как лексикографически наименьшую строку из данного класса, определенную на стандартном алфавите Λ . Согласно этому определению, каждый класс b -эквивалентных строк на алфавите Λ имеет бесконечную мощность, но с учетом леммы 4.3.2 без потери общности мы можем упростить ситуацию, если ограничим этот класс только p -каноническими строками. Тогда, например, класс S p -канонических b -эквивалентных строк, соответствующих массиву граней $\beta = 00100$, вполне обозрим:

$$S = \{\lambda_1\lambda_2\lambda_1\lambda_3\lambda_2, \lambda_1\lambda_2\lambda_1\lambda_3\lambda_3, \lambda_1\lambda_2\lambda_1\lambda_3\lambda_4\},$$

Для этого класса b -каноническим элементом будет строка $x^* = \lambda_1\lambda_2\lambda_1\lambda_3\lambda_2$.

Для того чтобы сделать вычисление b -канонических строк эффективным, надо ввести структуры, подобные деревьям, используемым в разделе 4.2, тогда канонические строки длиной n можно будет вычислить на основе канонических строк

длиной $n - 1$. Для этого необходимо понять, как элемент $\beta[n]$ массива граней β можно вычислить, зная только предшествующие ему элементы $\beta[1..n - 1]$. Напомним, что из леммы 1.3.1 следует, что для любой строки $x = x[1..n]$ возможные значения $\beta[n]$ являются элементами убывающей последовательности

$$\langle 1\beta^1[n - 1] + 1, \beta^2[n - 1] + 1, \dots, \beta^k[n - 1] + 1, 0 \rangle, \quad (4.9)$$

где $\beta^1[n - 1] \equiv \beta[n - 1]$, $\beta^i[n - 1] = \beta[\beta^{i-1}[n - 1]]$ ($i \in 2..k$) и k — наименьшее целое, такое, что $\beta^k[n - 1] = 0$. Докажем более строгий результат о том, что множество значений, действительно принимаемых величиной $\beta[n]$, не зависит от префикса $x[1..n - 1]$.

Лемма 4.3.3. Для любого целого $n \geq 2$ значения, принимаемые $\beta[n]$, зависят только от $\beta[1..n - 1]$ и размера алфавита.

Доказательство. Предположим, что есть две строки x и y длиной $n - 1$, определенные на алфавите размером α и имеющие одинаковый массив граней $\beta[1..n - 1]$. Предположим далее, что для некоторой буквы λ и некоторого целого m строка $x' = x\lambda$ имеет массив граней $\beta = \beta[1..n] = \beta[1..n - 1]m$ и при этом не существует такой буквы μ , что строка $y' = y\mu$ имела бы тот же самый массив граней β . Тогда число m должно принимать одно из значений множества (4.9).

Сначала рассмотрим случай, когда $m = \beta^i[n - 1] + 1$ для некоторого $i \in 1..k$. Поскольку $\beta^i[n - 1] = m - 1$, то отсюда следует, что

$$y[1..m - 1] = y[n + 1 - m..n - 1].$$

Так как $\beta^i[n] \neq m$, то, положив $y[n] = y[m]$, для некоторого $m' > m$, получим

$$y[1..m'] = y[n + 1 - m'..n].$$

Но это означает, что

$$y[1..m' - 1] = y[n + 1 - m'..n - 1].$$

Поэтому $\beta^i[n - 1] = m' - 1 > m - 1$, что противоречит условию $\beta^i[n - 1] = m - 1$. Таким образом, лемма справедлива для любого целого $m = \beta^i[n - 1] + 1$.

Теперь пусть $m = 0$. Тогда для любой из α букв $y'[n] = \mu$ получаем значение $\beta[n] > 0$, тогда как в случае $x'[n] = \lambda$ имеем $\beta[n] = 0$. Следовательно, существует целое $m' > 0$, такое, что для $y'\beta[n] = m'$, но для $x'\beta[n] \neq m'$, что противоречит уже рассмотренному случаю.

Отсюда следует вывод, что массив β будет массивом граней для некоторой строки x' только тогда, когда он будет массивом граней и для строки y' . ■

Этот фундаментальный результат дает возможность (см. работы [90, 79]) вычислить массив $\beta[1..n]$ на основе массива $\beta[1..n-1]$ без непосредственного обращения к исходной строке. Как можно использовать этот результат, мы покажем далее, а сейчас докажем, что любой префикс b -канонической строки $x_n^* = x^*[1..n]$ также будет b -канонической строкой.

Лемма 4.3.4. Для любого $n \geq 1$ в b -канонической строке $x_n^* = x_{n-1}^* \lambda$, где λ — некоторая буква стандартного алфавита, префикс x_{n-1}^* также является b -канонической строкой.

Доказательство. Предположим, что строка $x_n^* = x_{n-1}^* \lambda$ имеет массив граней $\beta_n = \beta[1..n]$ и префикс x_{n-1}^* не является b -канонической строкой. Обозначим через $\beta_{n-1} = \beta[1..n-1]$ массив граней строки x_{n-1}^* . Поскольку x_{n-1}^* не является b -канонической строкой, то существует строка $y_{n-1} < x_{n-1}^*$ с тем же массивом граней β_{n-1} . Тогда вследствие леммы 4.3.3, существует строка $y_n = y_{n-1} \lambda'$ с массивом граней β_n , но при этом $y_n < x_n^*$. Последнее неравенство противоречит условию, что строка x_n^* является b -канонической. Лемма доказана. ■

Из этой леммы вытекает, что все b -канонические строки x_n^* можно сформировать из b -канонических строк x_{n-1}^* . Это дает нам надежду, что, как и в случае с p -каноническими строками, можно сгенерировать b -канонические строки x_n^* на основе древовидной структуры, построенной для генерирования строк x_{n-1}^* . Но прежде чем перейти к построению такой структуры, нам необходимо выяснить, как на основе массива граней β_{n-1} можно вычислить различные массивы β_n .

Лемма 4.3.5. Пусть b -канонической строке x_{n-1}^* на стандартном алфавите Λ соответствует массив граней β_{n-1} . Тогда массив β_{n-1} будет префиксом в точности k различных массивов граней β_n только в том случае, когда строка $x_{n-1}^* \lambda_k$ будет b -канонической с массивом граней $\beta_{n-1} 0$.

Доказательство. Сначала предположим, что строка $x_{n-1}^* \lambda_k$ является b -канонической и имеет только пустую грань. Поскольку для каждой b -канонической строки соответствующий ей массив граней лексикографически наименьший, поэтому не существует буквы λ_i , $i < k$, такой, что строка $x_{n-1}^* \lambda_i$ имеет только пустую грань. Следовательно, для всех $i \in 1..k-1$ строки $x_{n-1}^* \lambda_i$ имеют различные непустые грани.

Теперь предположим, что для некоторого $i > k$ b -каноническая строка $x_{n-1}^* \lambda_i$ имеет наибольшую грань длиной $m > 0$, тогда $\beta_n = \beta_{n-1} m$. (Фактически, $m \geq 3$, поскольку $m \geq i > k \geq 2$.) Из леммы 4.3.4 следует, что строка x_{n-1}^* имеет b -канонический префикс $x_m^* = x_{m-1}^* \lambda_i$, где x_{m-1}^* также b -каноническая строка. Далее, поскольку строка $x_{n-1}^* \lambda_k$ имеет только пустую грань, то и префикс $x_{m-1}^* \lambda_k$ также имеет только пустую грань. Тогда для любого положительного целого $k' \leq k$ строка $x_{m-1}^* \lambda_{k'}$ является b -канонической с непустой гранью. Таким образом, мы свели

первоначальную задачу с конечными положительными целыми числами $n - 1$ и k к такой же задаче с другими конечными числами $m - 1$ и k' . Такое сведение теоретически можно продолжать бесконечно, что убеждает нас в невозможности существования такого $i > k$, что строка $x_{n-1}^* \lambda_i$ имеет непустую грань. Следовательно, существует ровно k различных массивов граней β_n , что и доказывает достаточность утверждения леммы.

Для доказательства необходимости утверждения леммы предположим, что существует ровно k различных массивов граней β_n . Но тогда один из них обязательно должен иметь вид $\beta_{n-1} 0$, а этому массиву, как мы уже убедились, должна соответствовать строка $x_{n-1}^* \lambda_k$. ■

Отметим, что лемма 4.3.5 не обязана выполняться для конечного алфавита Λ_α — возможно, она будет выполняться, если алфавит достаточно большой. Например, на алфавите $\Lambda_3 = \{\lambda_1, \lambda_2, \lambda_3\}$ b -каноническая строка $x_7^* = \lambda_1 \lambda_2 \lambda_1 \lambda_3 \lambda_1 \lambda_2 \lambda_1$ имеет массив граней $\beta_7 = 0010123$, но не существует строки $x_8^* = x_7^* \lambda$ с массивом граней вида $\beta_8 = 00101230$.

Леммы 4.3.3–4.3.5 составляют основу алгоритма генерирования b -канонических строк длиной n : к каждой b -канонической строке x_j^* ($j = 1, 2, \dots, n - 1$) добавляется по одной букве λ стандартного алфавита $\{\lambda_1, \lambda_2, \dots\}$ до тех пор, пока для некоторого k получим строку $x_j^* \lambda_k$, имеющую только пустую грань. Строки $x_j^* \lambda_1, x_j^* \lambda_2, \dots, x_j^* \lambda_k$ составляют полное множество b -канонических строк, которые можно получить на основе строки x_j^* .

Для реализации этого алгоритма надо уметь генерировать синтаксические деревья T'_n , подобные деревьям T_n , используемым для аналогичных целей в разделе 4.2. В дереве T'_n узлы также будут помечены парой (λ, β) , где $\lambda \in \Lambda$, а β — элемент массива граней для буквы λ в строке, определяемой метками узлов на пути от корня дерева до текущего узла. Дерево T'_1 состоит только из корневого узла с меткой $(\lambda_1, 0)$, а для $n \geq 2$ дерево T'_n формируется из дерева T'_{n-1} путем добавления к каждому листу (конечному узлу) этого дерева новых конечных узлов с метками $(\lambda_1, \beta_1), (\lambda_2, \beta_2), \dots, (\lambda_k, \beta_k)$.⁴ Таким образом, каждый узел дерева T'_n определяет b -каноническую строку совместно с ее массивом граней. Если обозначить через $b[n]$ количество b -канонических строк длиной n , то легко видеть, что количество листьев в дереве T'_n в точности равно $b[n]$. Итак, все $b[n]$ b -канонических строк (и соответствующих им массивов) можно получить, добавляя $b[n]$ первые потомки (сыновья⁵) к листьям дерева T'_{n-1} .

⁴Здесь важно отметить, что для каждого листа дерева T'_{n-1} значение k будет свое, не обязательно совпадающее с аналогичными значениями для других листьев. — *Примеч. ред.*

⁵“Детско-родительская” терминология широко применяется в теории графов, особенно при описании деревьев, и не должна вызывать затруднений у читателя. Напомним, что в предисловии автор предполагает, что читатель знаком с теорией графов и алгоритмами, выполняемыми над деревьями (алгоритмами прохождения и построения деревьев и т.п.). — *Примеч. ред.*

Обсудим, стоит ли строить дерево T'_n методом построения “в ширину” (так же как строится дерево T_n из предыдущего раздела), когда для вычисления $b[n]$ конечных узлов дерева T'_n предварительно уже построено дерево T'_{n-1} . Однако при построении дерева T'_n таким способом возникает проблема: чтобы вычислить сына данного узла N_j надо выполнить вычисления над массивом граней, для чего необходимо иметь доступ к строке x_j^* и ее массиву граней β_j , что, в свою очередь, требует просмотра всех узлов на пути от корня дерева до узла N_j . На такой просмотр потребуется время порядка $\Theta(j)$. Эта дополнительная работа неприятна, если мы хотим вычислить сынов узла N_j за константное время в расчете на одного сына — вместо этого каждый из k сынов узла N_j вычисляется в среднем за время порядка $\Theta(j/k)$.

Чтобы избежать просмотра пути до каждого узла N_j , можно строить дерево T'_n методом построения “в глубину”. В этом случае правый брат⁶ текущего узла N_j добавляется к дереву T'_n только в том случае, когда в это дерево введены все сыновья узла N_j . В этом случае до завершения всех вычислений в узле N_j требуется хранить в памяти строки x^* и соответствующие им массивы β , необходимые для этих вычислений. Например, при построении дерева T'_3 , показанного на рис. 4.2, узлы вводятся в порядке $ABDECFG$, а строки x^* и массивы β принимают соответственно такие значения:

$$\begin{aligned} x^* &= \lambda_1, \lambda_1\lambda_1, \lambda_1\lambda_1\lambda_1, \lambda_1\lambda_1\lambda_2, \lambda_1\lambda_2, \lambda_1\lambda_2\lambda_1, \lambda_1\lambda_2\lambda_2; \\ \beta &= 0, 01, 012, 010, 00, 001, 000. \end{aligned}$$

В этом методе построения дерева “в глубину” реализуется симметричный способ прохождения дерева, когда каждый его узел посещается не более двух раз. Подробности этого метода построения дерева предложено показать в упражнении 4.3.5.

Напомним, что при доказательстве теоремы 1.3.3 мы показали, что для любой строки длиной n внутренний цикл **while** в алгоритме 1.3.1 (это алгоритм вычисления массива граней) выполняется не более $n - 2$ раза. В контексте построения дерева T'_n это означает, что любой путь длиной n в этом дереве обязательно заканчивается листом. Поэтому для вычисления массива граней для n строк (префиксов), соответствующих этому пути, необходимо время, не превышающее $\Theta(n)$. Тогда для вычисления каждого узла на этом пути в среднем затрачивается $\Theta(n)/n$ времени, т.е. константное (фиксированное) время. Поскольку в дереве T'_n содержится $b[n]$ конечных узлов, на вычисление каждого из которых требуется константное время, то отсюда получаем следующее утверждение, аналогичное утверждению теоремы 4.2.3.

⁶“Братями” называются узлы дерева, имеющие общего “родителя”. — *Примеч. ред.*

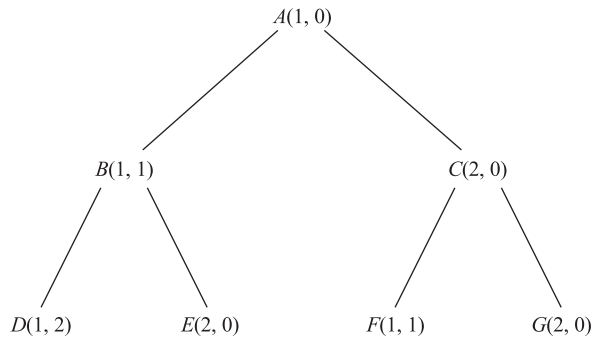


Рис. 4.2. Дерево T'_3 для вычисления всех b -различных строк, длина которых не превышает 3

Теорема 4.3.6. Для любого положительного целого n все b -канонические строки длиной n можно вычислить за время порядка $\Theta(b[n])$, для чего необходима память объемом $\Theta(b[n])$. ■

Простая модификация описанного алгоритма позволит вычислить все b -канонические строки длиной n , определенные на *конечном* алфавите Λ_α . Вместо условия введения новых сынов некоторого узла до тех пор, пока последний элемент массива граней будет равен нулю, применяется новое условие — новые сыны некоторого узла вводятся до тех пор, пока *или* последний элемент массива граней будет равен нулю *или* будут использованы все α букв. Как показано в работе [90], такой алгоритм можно легко модифицировать для того, чтобы он генерировал специальные деревья, определяемые заданным массивом граней. Однако такой алгоритм имеет существенный недостаток — он не является онлайн-овым, т.е. дерево T'_{n+1} нельзя легко сгенерировать на основании ранее построенного дерева T'_n . От этого недостатка данного алгоритма недавно удалось избавиться путем введения некоторых усовершенствований и упрощений в “механике” алгоритма [79].

Далее рассмотрим b' -массив, аналогичный p' -массиву из раздела 4.2. Для положительных целых чисел α и n $b'[\alpha, n]$ обозначает количество b -канонических строк длиной n , сформированных с использованием всех α букв алфавита Λ_α . Обозначим через $b[n]$ сумму элементов в n -м столбце b' -массива, т.е. $b[n] = \sum_{\alpha \geq 1} b'[\alpha, n]$. В общем случае найти замкнутую формулу для вычисления элементов b' -массива труднее, чем для p' -массива. Следующая теорема содержит информацию, которая поможет установить верхние границы для значений $b'[\alpha, n]$ и $b[n]$.

Теорема 4.3.7. Для любых положительных целых α и n

- а) если $\alpha > \lceil \log_2(n + 1) \rceil$, тогда $b'[\alpha, n] = 0$;
- б) $b'[1, n] = b'[\alpha, 2^{\alpha-1}] = 1$;

в) $b'[2, n] = p'[2, n] = 2^{\alpha-1} - 1;$

г) для любого положительного целого j и $\alpha > 2$

$$b'[\alpha, 2^{\alpha-1} + j] < p'[\alpha, \alpha + j].$$

Доказательство.

а) Доказательство этого утверждения проведем методом математической индукции. Сначала покажем, что утверждение выполняется для $n = 1$. Так как одна буква λ_1 является b -канонической строкой, то $p'[1, 1] = 1$. Тогда для любого $\alpha > 1$ $p'[\alpha, 1] = 0$, поскольку не существует строк длиной 1, сформированных из $\alpha > 1$ букв.

Теперь предположим, что утверждение справедливо для некоторого положительного целого α и для всех n , удовлетворяющих неравенствам $2^{\alpha-1} \leq n \leq 2^\alpha - 1$. Покажем, что в этом случае утверждение справедливо и для всех значений n' , удовлетворяющих неравенствам $2^\alpha \leq n' \leq 2^{\alpha+1} - 1$.

Из определения b' -массива следует, что предположение индукции эквивалентно утверждению о том, что для всех n , удовлетворяющих неравенствам $2^{\alpha-1} \leq n \leq 2^\alpha - 1$, для формирования b -канонических строк x_n , соответствующих любому массиву граней β_n , используется не более α букв $\lambda_1, \lambda_2, \dots, \lambda_\alpha$ (перечислены в возрастающем порядке). Поэтому в любой b -канонической строке $x_{n'}$, $n' \geq 2^\alpha$, буква $\lambda_{\alpha+1}$ не может встретиться в позиции, номер которой меньше 2^α .

Нам необходимо показать, что для любого n' , удовлетворяющего неравенствам $2^\alpha \leq n' \leq 2^{\alpha+1} - 1$, не существует b -канонической строки $x_{n'}$, содержащей букву $\lambda_{\alpha+2}$. Предположим противное — некоторая строка $x_{n'}$ содержит букву $\lambda_{\alpha+2}$ в качестве завершающей строку, т.е. $x_{n'}^* = x_{n'-1}^* \lambda_{\alpha+2}$. Такая строка может существовать только в том случае, если каждая из строк

$$\{x_{n'-1}^* \lambda_1, x_{n'-1}^* \lambda_2, \dots, x_{n'-1}^* \lambda_{\alpha+1}\}$$

будет b -канонической с непустой гранью. Рассмотрим строку $y = x_{n'-1}^* \lambda_{\alpha+1}$. Пусть буква $\lambda_{\alpha+1}$ в этой строке впервые встречается в позиции j . По предположению индукции $j \geq 2^\alpha$, следовательно, длина наибольшей грани строки y должна превышать $n'/2$. Но тогда $y[j - (n' - j)] = \lambda_{\alpha+1}$, что противоречит тому, что буква $\lambda_{\alpha+1}$ впервые встречается в позиции j . Отсюда следует, что строка $x_{n'-1}^* \lambda_{\alpha+1}$ не может иметь непустую грань. Поэтому на основании леммы 4.3.5 делаем заключение, что не существует b -канонической строки $x_{n'}$, содержащей букву $\lambda_{\alpha+2}$, что и требовалось доказать.

б) Значение $b'[1, n] = 1$ соответствует строке λ_1^n , значение $b'[\alpha, 2^{\alpha-1}] = 1 -$ строкам

$$X = \{\lambda_1, \lambda_1 \lambda_2, \lambda_1 \lambda_2 \lambda_1 \lambda_3, \lambda_1 \lambda_2 \lambda_1 \lambda_3 \lambda_1 \lambda_2 \lambda_1 \lambda_4, \dots\}.$$

Обозначим через x_α произвольную строку из множества X , $\alpha = 1, 2, \dots$, и представим ее как $x_\alpha = x'_\alpha \lambda_\alpha$. Нетрудно заметить, что любую строку из множества X можно генерировать последовательно с помощью соотношения $x_\alpha = x_{\alpha-1} x'_{\alpha-1} \lambda_\alpha$, где по определению $x_0 = x'_0 = \varepsilon$. Все строки x_α имеют только пустые грани. Кроме того, при $\alpha \geq 2$ они обладают тем свойством, что строки $x'_\alpha \lambda_i$, $1 \leq i \leq \alpha$, имеют непустые грани. Из леммы 4.3.5 следует, что b -канонические строки длиной $n = 2^{\alpha-1}$ должны заканчиваться буквой λ_α . Но существует только одна b -каноническая строка с таким свойством. Поэтому $b'[\alpha, 2^{\alpha-1}] = 1$.

- в) Заметим, что b -канонические строки x_n^* формируются из строк x_{n-1}^* путем присоединения к последним не менее двух различных букв. Присоединение к x_{n-1}^* буквы всегда приводит к появлению грани максимальной длины $\beta^* \geq 1$. Из леммы 4.3.5 следует, что только присоединение буквы λ_2 приведет к построению строки x_n^* , имеющей пустую грань. Но при $\alpha = 2$ к строкам x_{n-1}^* можно присоединить не более двух букв. Поэтому в данном случае каждая p -каноническая строка является также b -канонической. Поскольку лемма 4.3.2 утверждает, что каждая b -каноническая строка с необходимостью является p -канонической, то отсюда получаем требуемый результат.
- г) Для $\alpha \geq 3$ рассмотрим одну из b -канонических строк x_α длиной $n = 2^{\alpha-1}$, введенных при доказательстве утверждения б) данной теоремы. Такая строка является p -канонической и в дереве T_n порождает α p -канонических строк путем присоединения букв $\lambda_1, \lambda_2, \dots, \lambda_\alpha$. Однако только две b -канонические строки можно сформировать на основе x_α : $x_\alpha \lambda_1$ и $x_\alpha \lambda_2$. Но тогда поддерево с корнем в узле, соответствующем, например, строке $x_\alpha \lambda_3$, будет принадлежать дереву T_n , но не будет входить в дерево T'_n . Опять принимая во внимание тот факт, что каждая строка является p -канонической, приходим к выводу, что для любого $j > 1$ выполняется неравенство $b'[\alpha, n + j] < p'[\alpha, \alpha + j]$. Утверждение доказано. ■

В табл. 4.2 приведены значения верхнего левого угла b' -массива. Чтобы еще раз оценить преимущества от использования p - или b -канонических строк в качестве тестовых данных вместо различных в обычном понимании строк, отметим, например, что на 4-буквенном алфавите существует $4^9 = 262\,144$ различных строк, тогда как $p'[4, 9] = 7770$ (см. табл. 4.1) и из табл. 4.2 имеем $b'[4, 9] = 4$.

В общем случае мы уже видели (равенство (4.5)), что сумма $\sum_{i=1}^{\alpha} p'[i, n]$ асимптотически меньше величины α^n в $\alpha!$ (или больше) раз. Утверждение з) теоремы 4.3.7 говорит о том, что значения $b'[\alpha, n]$ в свою очередь значительно меньше значений $p'[\alpha, n]$. Фактически, это утверждение дает верхнюю границу для $b[n]$, выраженную посредством элементов p' -массива: для любого положительного це-

лого n имеем

$$b[n] \leq \sum_{k=1}^{k^*} \sigma_{n-2^{k-1}+k}^{(k)}, \quad (4.10)$$

где $k^* = \lceil \log_2(n+1) \rceil$. Это неравенство будет строгим уже при $n \geq 3$.

До настоящего времени не найдена замкнутая формула для $b[n]$, но в работе [G01] доказан впечатляющий результат, что $b[n] \leq f_{2n+1} < \phi^{2n}$, где f_{2n+1} — длина строк Фибоначчи f_{2n+1} , а $\phi = (1 + \sqrt{5})/2 \approx 1,62$ — отношение золотого сечения [135].

Таблица 4.2. Значения b' -массива для строк длиной $n \leq 10$ и $\alpha \leq 4$

n	1	2	3	4	5	6	7	8	9	10
$b'[1, n]$	1	1	1	1	1	1	1	1	1	1
$b'[2, n]$	0	1	3	7	15	31	63	127	255	511
$b'[3, n]$	0	0	0	1	4	15	46	134	370	997
$b'[4, n]$	0	0	0	0	0	0	0	1	4	16
$b[n]$	1	2	4	9	20	47	110	263	630	1525

Упражнения 4.3

1. Докажите лемму 4.3.2.
2. Докажите, что отношение b -эквивалентности действительно является отношением эквивалентности.
3. Постройте дерево T'_5 и убедитесь, что оно содержит 20 узлов.
4. Докажите по индукции, что $b[n] \geq \sum_{1 \leq i \leq n} b[i]$.
5. Напишите рекурсивный алгоритм построения дерева T'_n с константным временем для каждого узла.
6. Докажите неравенство (4.10).

ЧАСТЬ II

Вычисление внутренних паттернов

И Слово стало плотью и обитало с нами, полное благодати
и истины; и мы видели славу Его. . .

– Иоанн 1.14

ГЛАВА 5

Деревья для строковых последовательностей

... все слова,
но нет свершений!

– Филип Мессинджер (1583–1640). Парламент любви

Мы начинаем изучение внутренних паттернов с рассмотрения деревьев двух типов — деревьев граней и деревьев суффиксов, которые естественны для представления строковых последовательностей и возникают во многих приложениях. Мы отложим до раздела 13.1 рассмотрение другого внутреннего паттерна, дерева оболочек, который, подобно дереву граней, содержит информацию о периодических свойствах строк. В этой главе мы также рассмотрим две древовидные структуры суффиксов: ориентированные ациклические графы слов и массивы суффиксов.

5.1 Деревья граней

Как мы видели в разделе 1.2, наибольшие грани, присутствующие в строках, позволяют представить строки в нормальной форме. В разделе 1.3 показано, что все грани префиксов, т.е. массив граней, можно вычислить за линейное время, а сам массив используется для определения периодов и нормальных форм этих префиксов. В главе 1 понятие граней использовалось в доказательстве леммы периодичности (теорема 1.3.5) и для определения циклических сдвигов некратных

строк, что в главе 3 пригодилось при исследовании различных типов строк. В главе 4 эти понятия помогли в процессе разработки данных для тестирования строковых алгоритмов. Мы увидим, что грани естественны в алгоритмах вычисления частных паттернов (главы 7 и 8), кратных подстрок (подраздел 12.1.2) и оболочек (раздел 13.1).

Таким образом, массивы граней для внутренних паттернов представляют несомненный интерес, а как указывалось в разделе 2.1, эти массивы естественно представлять в виде древовидных структур. Здесь мы ограничимся только примером такой структуры для строк Фибоначчи, которая уже использовалась при пояснении определения 2.3.1.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
$f_7 =$	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>
$\beta =$	0	0	1	1	2	3	2	3	4	5	6	4	5	6	7	8	9	10	11	7	8

На рис. 5.1 показан массив граней β для данного примера, представленный в виде дерева. Отметим, что грани строк $f_7[1..i]$ определяются предками узлов i на этом дереве. Например, гранями строки $f_7[1..18]$ являются подстроки $f_7[1..10]$, $f_7[1..5]$, $f_7[1..2]$ и $f_7[1..0] = \varepsilon$.

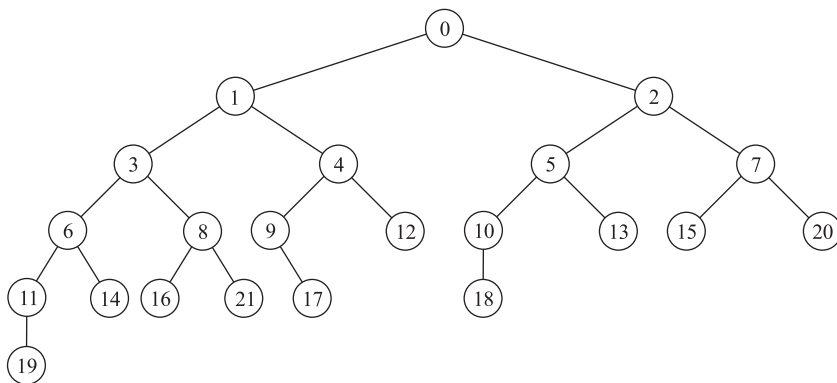


Рис. 5.1. Дерево граней для f_7

Упражнения 5.1

1. Эффективное прохождение дерева граней требует линейного времени, если проходить дерево снизу вверх, т.е. от сынов к родителям. Разработайте экономичную (по критерию объема памяти) структуру данных для представления этого дерева, которая дополнительно позволяла бы для любого родителя получать доступ к его сынам за константное время в расчете на

одного сына. (“Экономичность по критерию объема памяти” здесь обозначает, что общая требуемая память не должна превышать величины $3(n + 1)$.) Затем измените алгоритм 1.3.1 с учетом новой структуры данных таким образом, чтобы время выполнения алгоритма было линейным относительно длины строк.

2. Покажите, что дерево граней для строк Фибоначчи всегда будет бинарным, т.е. каждый узел будет иметь не более двух сынов.
3. Покажите, что в дереве граней поддерево с корнем в узле $i \geq 1$ определяет все раппорты подстроки $x[1..i]$ в строке x .

5.2 Деревья суффиксов

Как указывалось в разделе 2.1, деревья суффиксов имеют многочисленные применения в задачах со строковыми последовательностями [12]. Наиболее часто они используются на “подготовительном этапе” для того, чтобы заменить исходную входную строку x на дерево суффиксов T_x . Дерево T_x можно эффективно применить для поиска вхождений заданной подстроки u , определенной на алфавите A , в строку x с временными затратами порядка $O(|u| \log \alpha)$, что особенно эффективно, когда $|u|$ значительно меньше $|x|$, $\alpha = |A|$ мало, а строка x фиксированна. Такие условия выполняются, например, в последовательностях ДНК.

Как показано в разделе 2.1 (и упражнении 2.1.7), для хранения дерева T_x необходима память объемом $\Theta(nf(\alpha))$, где $f(\alpha)$ — некоторая функция от α или малая по величине константа, зависящая от природы алфавита. В этом разделе мы покажем, что дерево T_x можно построить за время порядка $O(n \log n)$, когда алфавит упорядочен, и за время $\Theta(n)$, если алфавит индексирован (вернитесь к разделу 4.1 для объяснения такого различия).

В этом разделе представлены три алгоритма построения дерева суффиксов, которые будут описаны после предварительного подраздела, где поговорим о деревьях суффиксов “в общем”. Два из этих алгоритмов, разработанные Мак-Крейтом (McCreight, [174]) и Укконеном (Ukkonen, [225]), выполняются на упорядоченных алфавитах, а алгоритм Фараха (Farach, [86]) работает на индексированных алфавитах. Существует еще один алгоритм построения дерева суффиксов — в разделе 12.1.1 мы обнаружим, что такое дерево появляется как побочный продукт алгоритма вычисления непродолжаемых вправо раппортов. Алгоритм Мак-Крейта построен на основе более раннего алгоритма Винера [230]. Эти алгоритмы имеют одинаковое время выполнения, порядка $O(n \log \alpha)$, однако алгоритм Винера более сложный и требует большего объема памяти, поэтому здесь он рассматриваться не будет. Однако отметим, что модифицированная версия алгоритма Винера [49] позволяет производить поиск обратных подстрок в заданной строке x . В недав-

ней работе [11] предложен алгоритм построения особого типа дерева суффиксов, который находит применение в анализе текста и последовательностей ДНК.

В этом и следующем разделах предполагается, что объем α алфавита конечен, но это не ограничивает возможностей применения описанных алгоритмов: если алфавит бесконечен, то в качестве значения α можно принять число различных букв, фактически принимающих участие в формировании данной исследуемой строки x . Однако подсчет количества различных букв, входящих в строку x , также требует определенных затрат, что показано в разделе 4.1 на примере задачи 4.1.

5.2.1 Предварительные сведения о деревьях суффиксов

Алгоритмы, описанные далее в этом подразделе, используют специальные способы эффективного прохождения дерева суффиксов после его построения. В этом подразделе мы покажем способы обхода дерева суффиксов и представим предварительные соображения об их роли в этих алгоритмах.

Начнем с определения *наибольшего общего префикса* u для двух данных строк x_1 и x_2 — это наибольшая строка u , которая является префиксом обеих строк x_1 и x_2 . Для наибольшего общего префикса введем обозначение $u = \text{LCP}(x_1, x_2)$.¹ Пусть строки x_1 и x_2 являются суффиксами некой строки x , т.е. $x_1 = x[i..n]$ и $x_2 = x[j..n]$. В этом случае наибольший общий префикс для таких строк будем обозначать как $\text{LCP}(i, j)$. Например, для строки Фибоначчи

$$\begin{array}{ccccccccccccccc}
 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\
 f_6 = & a & b & a & a & b & a & b & a & a & b & a & a & b
 \end{array}$$

имеем $\text{LCP}(baababaabaab, baabaab) = \text{LCP}(2, 7) = baaba$, тогда как $\text{LCP}(3, 7) = \varepsilon$ и $\text{LCP}(5, 7) = ba$.

Этот простой пример позволит сделать нам элементарное, но очень важное наблюдение. Пусть строка u определяет узел дерева T_x и пусть этот узел является корнем минимального поддеревя, содержащего суффиксы $x[i..n]$ и $x[j..n]$ строки x . В этом случае u будем называть *наименьшим общим предком* этих суффиксов и будем обозначать как $u = \text{LCA}(i, j)$.² Поскольку строка u определяет наибольший путь от корня дерева T_x до соответствующего узла, отсюда немедленно получаем определяющее свойство дерева суффиксов: для любой пары суффиксов $x[i..n]$ и $x[j..n]$ выполняется равенство

$$\text{LCA}(i, j) = \text{LCP}(i, j). \tag{5.1}$$

Мы покажем практическое применение этого соотношения при описании алгоритма Фараха (подраздел 5.2.4).

¹LCP обозначает Longest Common Prefix, т.е. наибольший общий префикс. — *Примеч. ред.*

²LCA обозначает Lowest Common Ancestor, т.е. наименьший общий предок. — *Примеч. ред.*

После построения дерева суффиксов для строки x необходимо для каждого суффикса $x[i..n]\$$ определить узлы $LCP(i, j)$, в которых путь по дереву T_x , определяющий префикс u , ветвится на пути к суффиксу $x[i..n]\$$ и на пути к другим суффиксам $x[j..n]\$$. В алгоритме Мак-Крейта, где суффиксы вычисляются в естественном порядке $i = 1, 2, \dots, n$, наибольший общий префикс вычисляется на множестве *наибольших* суффиксов $x[j..n]\$, j < i$, как строка, имеющая *максимальную* длину. В этом случае $LCP(i, j)$ определяет самый низкий узел в дереве T_x , в котором путь к суффиксу $x[i..n]\$$ имеет ответвление к любому ранее вычисленному суффиксу. Поэтому разобьем суффикс, определяющий конечный узел i ($1 \leq i \leq n + 1$) дерева T_x , на две части:³

$$x[i..n]\$ = \text{head}(i) \text{tail}(i),$$

где $\text{head}(1) = \varepsilon$ и для $i \geq 2$ $\text{head}(i) = LCP(i, j)$ максимальной длины, вычисленный по всем $j : 1 \leq j \leq i - 1$.⁴ Как упоминалось в разделе 2.1 и как видно на рис. 5.2, маркером окончания строки является символ $\$,$ показывающий, что подстрока $\text{tail}(i)$ не пустая. Это соглашение выполняется и в случае $i = n + 1$ — здесь строка $x[n + 1..n]\$$ дает $\text{head}(n + 1) = \varepsilon$ и $\text{tail}(n + 1) = \$$.

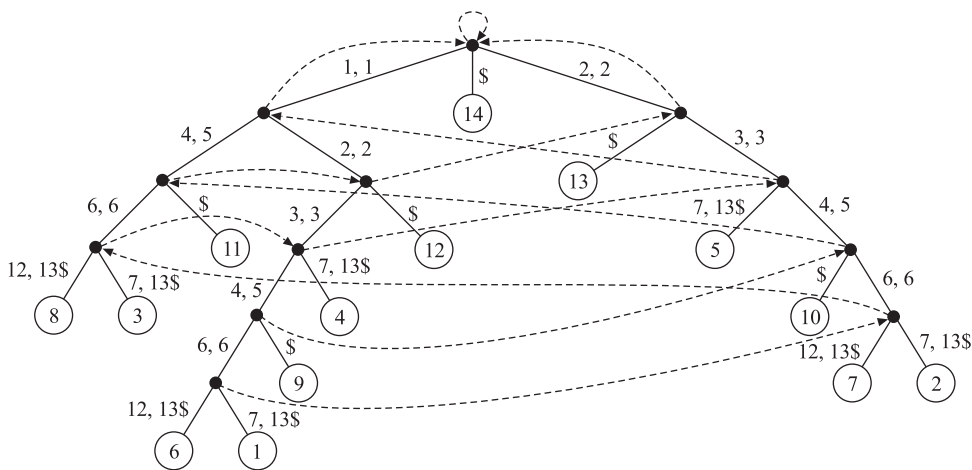


Рис. 5.2. Связное дерево суффиксов для строки Фибоначчи f_6

³Здесь и далее автор часто отождествляет узел и подстроку, соответствующую этому узлу. Поэтому можно встретить выражения, подобные тому, что один узел является префиксом другого узла или что одна подстрока является предком другой подстроки и т.п. Надеемся, что такая “взаимоподмена” понятий не вызовет у читателя затруднений (хотя вполне обоснованное раздражение вызвать может). — Примеч. ред.

⁴Название “head” означает “голова” или “передняя часть”, а “tail” — “хвост” или “задняя часть”. — Примеч. пер.

Заметим, что подстрока $\text{head}(i)$ никогда не будет соответствовать конечному узлу. Интуитивно ясно, что $\text{head}(i)$ является отправной точкой (префиксом), от которой начинается поиск в T_x суффикса $x[i..n]$. В табл. 5.1 приведены значения подстрок $\text{head}(i)$ для суффиксов строки Фибоначчи f_6 (см. также рис. 5.2).

Таблица 5.1. Значения подстрок $\text{head}(i)$ для суффиксов строки Фибоначчи f_6

i	$\text{head}(i)$	i	$\text{head}(i)$
1	ε	8	$aaba$
2	ε	9	$abaab$
3	a	10	$baab$
4	aba	11	aab
5	ba	12	ab
6	$abaaba$	13	b
7	$baaba$	14	ε

Очевидно, что подстроке $\text{tail}(i)$ в дереве T_x соответствует путь от внутреннего узла, определяемого $\text{head}(i)$, до конечного узла. В упражнении 5.2.1 предложено показать, что $\text{tail}(i)$ *всегда* соответствует конечному узлу. Например, на рис. 5.2 видно, что подстроке $\text{tail}(5) = baabaab$ соответствует конечный узел 7.

Следующая лемма, используемая в алгоритме Мак-Крейта, показывает, как на основе известного значения $\text{head}(i)$ можно локализовать в дереве T_x строку $\text{head}(i + 1)$. Если говорить неформально, то эта лемма утверждает, что наибольший собственный суффикс строки $\text{head}(i)$ является префиксом строки $\text{head}(i + 1)$ и поэтому будет предком строки $\text{head}(i + 1)$ в дереве T_x .

Лемма 5.2.1. Пусть $\text{head}(i) = x[i..i + h]$ для некоторых целых i , $1 \leq i \leq n$, и h , $-1 \leq h \leq n - i$. Тогда подстрока $x[i + 1..i + h]$ является префиксом строки $\text{head}(i + 1)$.

Доказательство. Утверждение леммы тривиально, когда $i + h \leq i$, поскольку в этом случае строка $x[i + 1..i + h] = \varepsilon$ является префиксом любой строки, в том числе и строки $\text{head}(i + 1)$. Поэтому рассмотрим случай, когда $i + h > i$. Обозначим $\lambda = x[i]$ и $u = x[i + 1..i + h]$, тогда $\text{head}(i) = \lambda u$. По определению подстрок head существует номер $j \in 1..i - 1$ такой, что строка λu является префиксом как суффикса $x[i..n]$, так и суффикса $x[j..n]$. Но тогда строка u также является префиксом обоих суффиксов $x[i + 1..n]$ и $x[j + 1..n]$ и поэтому должна быть префиксом строки $\text{head}(i + 1)$. ■

Отметим, что все подстроки $\text{head}(i)$ для строки f_6 , приведенные в табл. 5.1, удовлетворяют лемме 5.2.1.

На основе доказанной леммы можно ввести другой способ обхода дерева T_x , который будет востребован в алгоритмах Мак-Крейта и Укконена. Определим **суффиксную функцию** s для любого узла u дерева T_x следующим образом:

- для $u = \varepsilon$, $s(u) = \varepsilon$;
- для $u \neq \varepsilon$ (в этом случае $u = \lambda v$ с непустой буквой $\lambda \in A$) $s(u) = v$.

Суффиксная функция является указателем от любого непустого (конечного или внутреннего) узла u дерева T_x к узлу $v = s(u)$, который является наибольшим собственным суффиксом строки u . В упражнении 5.2.3 предложено доказать, что если u является узлом дерева T_x , то узлом этого дерева будет и $s(u)$. Отметим, что если u является конечным узлом $i \leq n$, тогда $s(u)$ также является конечным узлом $i + 1$. Если же u — конечный узел $n + 1$, тогда $s(u) = \varepsilon$. Когда $u = \text{head}(i)$, то, вследствие леммы 5.2.1, $s(u)$ будет префиксом $\text{head}(i + 1)$.

Дерево суффиксов, включающее связи между узлами, определяемые суффиксной функцией s , называется **связным деревом суффиксов**. На рис. 5.2 представлено связанное дерево суффиксов для строки Фибоначчи f_6 (показаны связи только для внутренних узлов дерева).

Теперь рассмотрим механизм поиска непустой строки v , при этом поиск начинается из заданного узла u дерева T_x (узел u назовем **префиксным узлом**). Другими словами, для данного префиксного узла u , который уже локализован в дереве T_x , необходимо найти префикс uv строки x . Такая ситуация представлена на рис. 5.3, где узел u имеет k различных сынов uu_i , $i = 1, 2, \dots, k$, и $u_i = \lambda_i u'_i$ с попарно различными буквами λ_i . Мы ищем совпадение u_i со строкой $v = \lambda v'$.

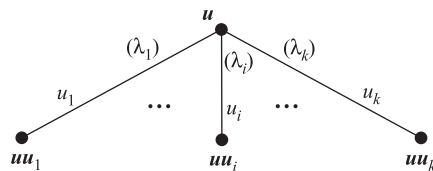


Рис. 5.3. Сканирование дерева суффиксов

Если мы не знаем, является ли в действительности uv префиксом строки x , то в этом случае поочередно просматриваются метки ребер λ_i , и при обнаружении совпадения $\lambda = \lambda_i$ продолжается побуквенное сравнение v' с метками ребер, соответствующих строке u'_i . Такой процесс назовем **медленным сканированием**.

Если известно, что uv действительно является префиксом строки x , то в таком случае надо только определить точную позицию uv в строке x (или ее положение в дереве T_x). Для этого не обязательно производить полное побуквенное сравнение v' и u'_i , поскольку после нахождения совпадения $\lambda = \lambda_i$ и определения u_i имеются две возможности.

1. $|v| > |u_i|$. Тогда $v = u_i v''$ для некоторой непустой строки v'' , и поиск совпадений продолжается рекурсивно для v'' с новым префиксным узлом $u u_i$.
2. $|v| \leq |u_i|$. Тогда $v = u_i[1..|v|]$, и поиск заканчивается.

Поскольку значение $|u_i|$ можно получить из метки ребра, поэтому определить, какая из двух описанных возможностей в действительности имеет место, можно за константное время. Таким образом, время поиска строки v не пропорционально $|v|$, а значительно меньше и зависит от количества просмотренных узлов. Поэтому такой процесс назовем *быстрым сканированием*. Как увидим далее, использование быстрого сканирования является ключевым фактором уменьшения времени выполнения алгоритмов Мак-Крейта и Укконена.

5.2.2 Алгоритм Мак-Крейта

Алгоритм Мак-Крейта (будем для краткости называть его алгоритмом МК) начинается с построения дерева T'_1 , содержащего только два узла: корень ε и конечный узел 1, соответствующий суффиксу $x\$ = x[1..n]\$$. Эти узлы связаны ребром, помеченным $x\$$. Дерево T'_1 также содержит суффиксную связь для корня $s(\varepsilon) = \varepsilon$.

Далее для любого $i = 1, 2, \dots, n$ алгоритм МК вычисляет новое дерево T'_{i+1} на основе уже построенного дерева T'_i путем добавления суффикса $x[i + 1..n]\$$. Эти суффиксы добавляются в убывающем порядке их длин, и поскольку конечных узлов $i + 2, i + 3, \dots, n + 1$ нет в дереве T'_{i+1} , полное дерево суффиксов T_x будет построено только на шаге $n + 1$, т.е. $T_x = T'_{n+1}$.

На i -й итерации в алгоритме МК на основе дерева T'_i выполняются следующие операции.

1. Добавляется конечный узел $i + 1$.
2. Добавляется внутренний узел $\text{head}(i + 1)$.
3. Вычисляется метка $\text{tail}(i + 1)$ для ребра, идущего от узла $\text{head}(i + 1)$ до конечного узла $i + 1$.
4. Устанавливается суффиксная связь $s(\text{head}(i))$.

Из последнего пункта видно, что каждое дерево T'_{i+1} является связным, но суффиксные связи вычисляются только для внутренних узлов. После завершения дерева T'_{i+1} суффиксная связь может быть не определена только для внутреннего узла $\text{head}(i + 1)$.

Представленные выше операции 1 и 3 тривиальны, если локализован узел $\text{head}(i + 1)$, операция 4, как увидим далее, может быть реализована в процессе вычисления $\text{head}(i + 1)$. Поэтому основное внимание уделим выполнению операции 2. Но сначала напомним следующие факты.

■ Строка $s(\text{head}(i))$ является префиксом (предком в дереве T_x) строки $\text{head}(i + 1)$ (лемма 5.2.1).

■ Обозначим через $\text{parent}(\mathbf{u})$ родителя в дереве T_x непустого узла \mathbf{u} , тогда или

$$\text{parent}(\mathbf{u}) = s(\mathbf{u}) = \varepsilon, \quad (5.2)$$

или $s(\text{parent}(\mathbf{u}))$ является собственным префиксом строки $s(\mathbf{u})$ (доказано в упражнении 5.2.3).

■ В дереве T'_{i+1} для непустой строки $\text{head}(i)$ вычислена суффиксная связь $s(\text{parent}(\text{head}(i)))$.

На основе этих фактов заключаем, что можем использовать узел $s(\text{parent}(\text{head}(i)))$ в качестве отправной точки для поиска $\text{head}(i + 1)$. Из этого могут быть два исключения.

■ Если $x[i]$ является самым левым (первым) вхождением некоторой буквы λ в строку x , тогда $\text{head}(i) = \varepsilon$. В этом случае отправной точкой для поиска $\text{head}(i + 1)$ будет корень дерева, при этом, очевидно, $\text{parent}(\text{head}(i)) = \varepsilon$.

■ Если $\text{head}(i) \neq \varepsilon$, но тем не менее выполняется равенство (5.2), тогда $\text{head}(i) = \lambda^k$ для некоторого положительного целого k . В данном случае отправной точкой для поиска $\text{head}(i + 1)$ будет узел λ^k . (Заметим, что при $k = 1$ отправной точкой опять будет корень дерева.)

Из дальнейшего рассмотрения эти два случая исключаем. Обозначим $\mathbf{u} = \text{parent}(\text{head}(i))$ и $\mathbf{v} = \text{label}(\mathbf{u}, \text{head}(i))$ — метка (label) ребра, проходящего от узла \mathbf{u} до узла $\text{head}(i)$. В этих обозначениях можно записать $\text{head}(i) = \mathbf{u}\mathbf{v}$, так что строка $\mathbf{w} = s(\mathbf{u})\mathbf{v}$ должна быть префиксом строки $\text{head}(i + 1)$. Тогда, даже если строка $\text{head}(i + 1)$ не определена как узел дерева T'_i , подстрока \mathbf{v} должна быть просмотрена на некотором нисходящем пути из $s(\mathbf{u})$. Заметим, что здесь выполняются условия, необходимые для быстрого сканирования, приведенные в разделе 5.2.1. Поэтому можем формально выразить поиск строки \mathbf{v} , который начинается от префиксного узла $s(\mathbf{u})$, как

$$\mathbf{w} \leftarrow \text{fastscan}(s(\mathbf{u}), \mathbf{v}),$$

где быстрое сканирование рассматривается как функция fastscan , которая определяет положение в дереве T'_{i+1} строки $\mathbf{w} = s(\mathbf{u})\mathbf{v}$.

Возможны две ситуации, в зависимости от того, является или нет строка \mathbf{w} узлом дерева T'_i . Если строка \mathbf{w} уже представлена узлом дерева T'_i , тогда на нисходящем пути от узла \mathbf{w} существует непустой префикс подстроки $x[|\mathbf{w}| + 1..n]$, при этом наибольшим таким префиксом будет строка $\text{head}(i + 1)$, которую можно вычислить с помощью медленного сканирования. Эту операцию представим в виде функции slowscan (медленное сканирование):

$$\text{head}(i + 1) \leftarrow \text{slowscan}(\mathbf{w}, \text{tail}(i)).$$

Если после вычислений окажется, что строка $\text{head}(i + 1)$ еще не представлена узлом дерева, то такой узел надо создать. Отметим, что возможна ситуация, когда $\text{head}(i + 1) = w$.

Теперь рассмотрим случай, когда в дереве T'_i нет узла, представляющего строку w . Сначала предположим, что $\text{head}(i + 1) \neq w$, тогда ребро, выходящее из узла w , должно вести к непустому префиксу строки $\text{tail}(i)$. Но, по определению $\text{head}(i)$, из узла, представляющего $\text{head}(i)$, должно выходить не менее двух ребер, ведущих к различным префиксам, одним из которых будет $\text{tail}(i)$. Поскольку $w = s(\text{head}(i))$, тогда из узла w также должно выходить не менее двух ребер, ведущих к различным префиксам. Получили противоречие. Отсюда заключаем, что возможно только равенство $\text{head}(i + 1) = w$, которое не требует дальнейшего поиска.

Основываясь на представленных рассуждениях, запишем в явном виде алгоритм Мак-Крейта.

Алгоритм 5.2.1 (Алгоритм Мак-Крейта)

▷ Вычисление связного дерева суффиксов T_x

Построение дерева T'_1

```

for  $i \leftarrow 1$  to  $n$  do           ▷ построение дерева  $T'_{i+1}$  на основе дерева  $T'_i$ 
     $u \leftarrow \text{parent}(\text{head}(i))$    ▷  $u \leftarrow \varepsilon$ , если  $\text{head}(i) = \varepsilon$ 
     $v \leftarrow \text{label}(u, \text{head}(i))$   ▷  $v \leftarrow \varepsilon$ , если  $\text{head}(i) = \varepsilon$ 
    if  $u \neq \varepsilon$  then
         $w \leftarrow \text{fastscan}(s(u), v)$ 
    else
         $w \leftarrow v[1..|v| - 1]$ 
    if  $w$  — новый узел then
         $\text{head}(i + 1) \leftarrow w$ 
        добавить  $\text{head}(i + 1)$  в дерево  $T'_i$ 
    else
         $\text{head}(i + 1) \leftarrow \text{slowscan}(w, \text{tail}(i))$ 
        if  $\text{head}(i + 1)$  — новый узел then
            добавить  $\text{head}(i + 1)$  в дерево  $T'_i$ 
     $s(\text{head}(i)) \leftarrow w$ 
    добавить конечный узел  $i + 1$  в дерево  $T'_i$ 
     $\text{label}(\text{head}(i + 1), i + 1) \leftarrow \text{tail}(i + 1)$ 
    
```

Теорема 5.2.2. Алгоритм 5.2.1 корректно вычисляет связное дерево суффиксов T_x .

Доказательство корректности основывается на следующих свойствах алгоритма:

- для непустой строки x в дереве T'_{n+1} существует невнутренний узел с менее чем двумя сыновьями;

- ребра, выходящие из внутренних узлов, ведут к разным префиксам;
- суффиксная связь для любого внутреннего узла устанавливается корректно;
- в дереве T'_{n+1} для любого $i = 1, 2, \dots, n + 1$ существует конечный узел i , соответствующий суффиксу $x[i..n]$. ■

Установить временную сложность алгоритма МК является более сложным делом.

Теорема 5.2.3. Алгоритм 5.2.1 вычисляет связное дерево суффиксов T_x для заданной строки $x[1..n]$, определенной на упорядоченном алфавите A размером α , за время $O(n \log \alpha)$ с использованием памяти объемом $O(n\alpha)$.

Доказательство. За исключением операций быстрого и медленного сканирования, каждый шаг алгоритма 5.2.1 выполняется за константное время с использованием стандартных методов обхода дерева и его представления. Поэтому, если не учитывать операции сканирования, для выполнения алгоритма необходимо время порядка $\Theta(n)$.

Рассмотрим операцию быстрого сканирования. Обозначим через B максимальное (по всем узлам дерева) время, необходимое для выбора правильного пути из текущего узла. Тогда суммарное время, затрачиваемое на быстрое сканирование при выполнении n шагов алгоритма, составляет $O(nB)$. Но необходимо также учесть возможный рекурсивный вызов функции быстрого сканирования вследствие существования промежуточных узлов на нисходящем пути из узла $s(\text{parent}(\text{head}(i + 1))) = s(u)$ в узел $uv = s(\text{head}(i + 1))$, которые не имеют “двойников” на пути от u до uv . Каждый такой узел может потребовать времени для операции быстрого сканирования порядка $O(B)$. При выполнении n шагов алгоритма может встретиться m таких узлов, поэтому на их обработку потребуется время порядка $O(mB)$. Но поскольку в дереве T_x нет пути, длиннее n , то $m < n$. Следовательно, общее время, затрачиваемое на операцию быстрого сканирования при выполнении алгоритма, составляет $O(nB)$.

Теперь рассмотрим медленное сканирование. Заметим, что поскольку $s(\text{head}(i))$ является префиксом строки $\text{head}(i + 1)$, то

$$|\text{head}(i + 1)| - |\text{head}(i)| + 1 \geq 0. \quad (5.3)$$

Фактически, последнее неравенство позволяет получить точное количество букв, просматриваемых в процессе медленного сканирования для локализации $\text{head}(i + 1)$, $i = 1, 2, \dots, n$. Для этого неравенство (5.3) надо просуммировать по всем n шагам алгоритма. Получим

$$|\text{head}(n + 1)| - |\text{head}(1)| + n = n.$$

Таким образом, общее время операции медленного сканирования составляет $O(nB)$.

Теперь оценим значение времени B . Поскольку алфавит A упорядочен, в каждом узле дерева суффиксов этот алфавит можно представить в виде упорядоченного массива или дерева поиска. Такая структура содержит не более $\Theta(\alpha)$ элементов, поиск по ней требует времени порядка $O(\log \alpha)$. Отсюда получаем оценки, приведенные в утверждении теоремы. ■

В упражнении 5.2.5 предлагается провести дальнейший анализ временной и пространственной сложности алгоритма МК.

5.2.3 Алгоритм Укконена

Подобно алгоритму МК, алгоритм Укконена (для краткости будем называть его алгоритмом Укк) выполняется итерационно, начиная с исходного дерева и заканчивая полным деревом суффиксов T_x . Но алгоритм Укк применяет другую стратегию, отличную от стратегии алгоритма МК, — в алгоритме МК в промежуточных деревьях последовательно представляются суффиксы $x[i..n] \$$, $i = 1, 2, \dots, n + 1$, тогда как в алгоритме Укк в промежуточных деревьях последовательно представляются префиксы $x[1..i] \$$, $i = 1, 2, \dots, n$, заканчивая префиксом $x[1..n] \$$ (т.е. заданной строкой). Таким образом, алгоритм Укк является онлайн-алгоритмом (см. раздел 4.1) — каждое построенное промежуточное дерево T_i фактически является законченным связным деревом суффиксов для строки $x[1..i]$.

Однако в этом алгоритме метка $\$$ конца строки добавляется в дерево только на последней итерации построения дерева T_x . В промежуточных деревьях T_i ($1 \leq i \leq n$) конечные узлы (т.е. те узлы, которые представляют суффиксы) не обязательно являются листьями этих деревьев. Если префикс $x[1..i]$ имеет два суффикса u и uv , тогда узел, соответствующий суффиксу u , может быть внутренним узлом на пути от корня дерева до узла, представляющего суффикс uv . Более того, если строка u является *только префиксом* строки uv , тогда строка u вообще может не иметь узла в дереве T_i , ее представляющего.

Чтобы понять, как алгоритм Укк “рисует” деревья, рассмотрим пример построения дерева суффиксов для строки $f_4 = abaab$, показанного на рис. 5.4 (здесь *только для пояснения* ребра помечены значениями подстрок, а не парой чисел (“ i_1 ”, “ i_2 ”), определяющих позиции этих подстрок). Отметим, что здесь, как и в алгоритме МК, суффиксные связи определяются только для внутренних узлов (узлов ветвления).

На рис. 5.4 видно, что в деревьях T_3 и T_5 некоторые суффиксы не представлены узлами дерева, и только дерево T_6 представляет дерево суффиксов в привычной форме. Однако важно отметить, что преобразование дерева T_5 в дерево T_6 (в общем случае дерева T_n в дерево T_{n+1}) происходит по тем же правилам алгоритма, что и преобразование дерева T_i в дерево T_{i+1} , $1 \leq i \leq n - 1$.

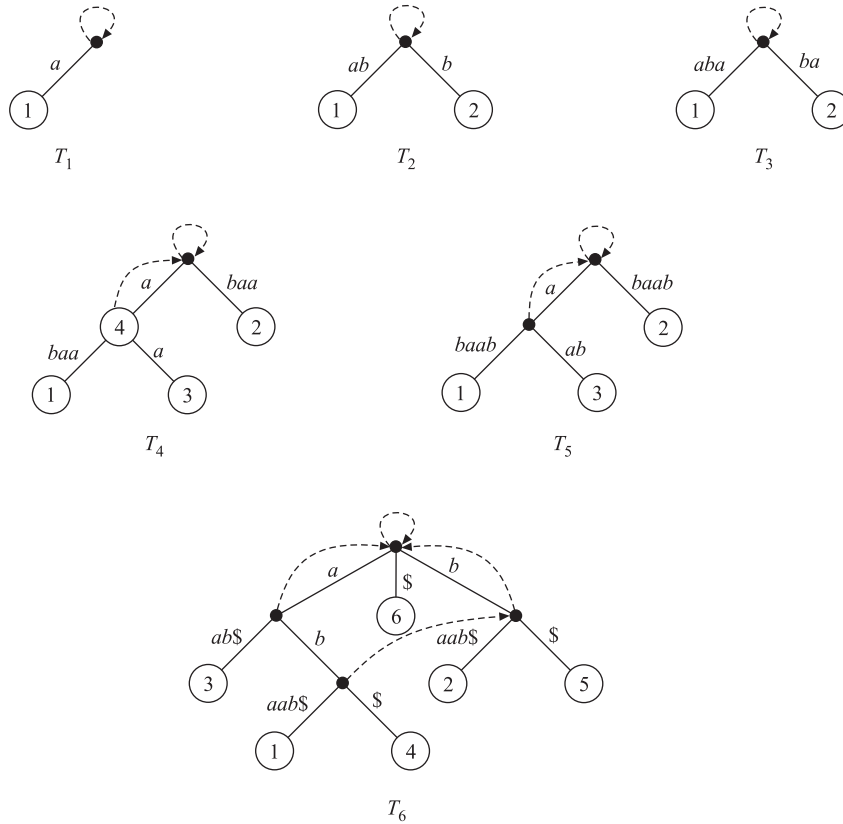


Рис. 5.4. Связное дерево суффиксов для f_4 , формируемое алгоритмом Укк

Поскольку алгоритм Укк является онлайн-овым, за один шаг он обрабатывает по одной букве строки x , поэтому при описании данного алгоритма основное внимание уделим преобразованию дерева T_i в дерево T_{i+1} , $1 \leq i \leq n$, при этом предполагая, что дерево T_1 всегда имеет такой вид, как показано на рис. 5.4. Чтобы немного упростить описание алгоритма, введем символ $\$$ в исходную строку, так что теперь $x[n+1] = \$$.

Если в дереве T_i определены i суффиксов $x[j..i]$ ($j \in 1..i$), то задача преобразования дерева T_i в дерево T_{i+1} будет сведена к определению $i+1$ суффиксов $x[j..i+1]$ в дереве T_{i+1} , при этом теперь $j \in 1..i+1$. При преобразовании дерева T_i в дерево T_{i+1} выполняются такие операции.

1. Новый j -й суффикс в дереве T_{i+1} получен из “старого” j -го суффикса дерева T_i добавлением к последнему одной буквы $x[i+1]$. Если “старый” суффикс уже представлен в дереве T_i листом, то все, что требуется сделать в этой ситуации, — это добавить к метке ребра, ведущего к этому листу, позицию

$i + 1$. Если же “старый” суффикс не представлен в дереве T_i листом, тогда необходимо “расширить” существующий в дереве T_i путь к суффиксу $x[j..i]$ до суффикса $x[j..i + 1]$. Такое расширение может потребовать или создание нового пути (и, следовательно, создание нового листа), или создание нового узла на уже существующем пути. Последний вариант возможен в случае, когда новый суффикс является префиксом уже существующего суффикса.

2. Новый $(i + 1)$ -й суффикс определяется как новый сын корня дерева T_i . Если буква $x[i + 1]$ ранее не встречалась в строке x , тогда $(i + 1)$ -й суффикс становится листом дерева T_{i+1} . В противном случае этот суффикс представлен на одном из существующих путей, исходящих от корня, а выполняемые при этом действия являются специфическим случаем операции, описанной в первом пункте, когда начальной точкой будет корень дерева.

Следующая лемма составляет основу операций, выполняемых при преобразовании дерева T_i в дерево T_{i+1} .

Лемма 5.2.4. Пусть число j определяет суффикс $x[j..n]$ строки x .

- а) Если $j > 1$ и представляет лист дерева T_x , тогда суффикс $j - 1$ также будет листом.
- б) Если в дереве T_x есть путь, начинающийся с некоторой буквы λ и проходящий через суффикс j ($j \leq n$), тогда существует путь, начинающийся с буквы λ и проходящий через суффикс $j + 1$.

Доказательство леммы предложено дать в упражнении 5.2.6. ■

Утверждения леммы справедливы для любого дерева суффиксов T_x , поэтому она применима к связным деревьям T_i , которые строятся в ходе выполнения алгоритма Укк. Поскольку в каждом дереве T_i ($i = 1, 2, \dots, n + 1$) суффиксу 1 всегда соответствует узел-лист, утверждение а) леммы говорит о том, что существует такое целое число $j_L \in 1..i$, что любой суффикс $j \in 1..j_L$ является листом, тогда как суффиксы $j \in j_L + 1..i$ таковыми не являются. Аналогично утверждение б) леммы говорит о том, что существует такое целое число $j_\lambda \in 2..i + 1$, что любому суффиксу $j \in j_\lambda..i$ в дереве T_i соответствует путь в этом дереве, который начинается в узле с буквой λ , тогда как суффиксам $j \in 1..j_\lambda - 1$ соответствуют пути, который начинаются в узлах, не совпадающих с буквой λ . Очевидно, что для любой буквы $\lambda j_L < j_\lambda$. На рис. 5.4 видно, что в дереве T_4 $j_L = 3$ и $j_a = j_b = 4$, тогда как в дереве T_6 $j_L = 3$ и $j_a = 4$, но $j_b = 6$ (поскольку нет суффикса, который следовал бы за последней буквой b).

Возвращаясь к преобразованию дерева T_i в дерево T_{i+1} , мы можем, основываясь на приведенных выше результатах и наблюдениях, разбить множество суффиксов на три непересекающихся класса (здесь $\lambda = x[i + 1]$, т.е. буква, добавляемая на i -м шаге алгоритма):

- $1 \leq j \leq j_L$ (листья дерева T_i);
- $j_L + 1 \leq j \leq j_{x_{[i+1]}} - 1$;
- $j_{x_{[i+1]}} \leq j \leq i$ (суффиксы, являющиеся префиксами подстроки $x[j..i + 1]$).

Суффиксы первого и третьего классов не требуют перестройки дерева T_i , а требуют только внесения изменений в метки ребер. Рассмотрим сначала первый класс суффиксов. Напомним замечание, сделанное выше, что для суффикса j , являющегося листом, необходимо только добавить букву $x[i + 1]$ к пути, ведущему к этому префиксу. Но поскольку данный лист дерева T_i останется листом в дереве T_{i+1} , добавление всей последовательности букв $x[i + 1], x[i + 2], \dots, x[n + 1]$ можно реализовать, просто добавив метку “ ∞ ” к ребру, ведущему к узлу j . Таким образом, если узел-лист j является сыном некоего узла $u = x[1..h]$, то ребро от узла u до узла j помечается как $(h + 1, \infty)$, что указывает на суффикс $x[h + 1..n + 1]$, присоединенный к этому ребру.

Для третьего класса суффиксов также нет необходимости изменять дерево T_i , поскольку уже существует путь от суффикса j , который начинается с буквы $\lambda = x[i + 1]$, — для преобразования дерева T_i в дерево T_{i+1} достаточно суффикс j переместить по дереву на одну позицию вниз для обозначения строки $x[j..i + 1]$, которая уже должна быть представлена в дереве T_i . В последнем случае надо эффективно локализовать суффикс j , для чего можно использовать метод быстрого сканирования, описанный в разделе 5.2.1.

Таким образом, в процессе преобразования дерева T_i в дерево T_{i+1} следует основное внимание уделить суффиксам “второго класса” $j \in j_L + 1..j_{x_{[i+1]}} - 1$, которые не являются листьями в дереве T_i и при этом не встречались ранее в строке $x[1..i]$ в качестве префиксов подстроки $x[j..i + 1]$. Рассмотрим пример строки, в которой присутствуют суффиксы всех трех классов.

1	2	3	4	5	6	7	8	9	10	
$x = a$	a	b	b	a	c	b	b	a	a	\dots

Для $i = 9$ имеем $j_L = 6$, поскольку

- буква c в строке x впервые встречается в позиции 6;
- подстрока $x[7..9]$ встречалась ранее в строках $x[1..i]$ как $x[3..5]$.

Также имеем $j_{x_{[i+1]}} = j_a = 9$, поскольку подстрока $x[9..10] = aa$ встречалась ранее, а подстрока $x[8..10] = baa$ не встречалась. Для $i = 9$ имеем следующие три класса суффиксов $j : \{1..6, 7..8, 9..9\}$.

При преобразовании дерева T_i в дерево T_{i+1} алгоритм Укк начинает с проверки первого суффикса $j_L + 1$, который не является листом в дереве T_i . Если суффикс $j_L + 1$ не встречается в строке $x[1..i]$ после буквы $x[i + 1]$, тогда в дереве T_{i+1} создается лист $j_L + 1$, значение j_L увеличивается на единицу и далее проверяется

следующий суффикс. Когда или $j_L > i$ или суффикс $j_L + 1$ следует за буквой $x[i + 1]$, преобразование дерева T_i заканчивается.

Поиск суффикса $j_L + 1$ начинается от префиксного узла u . Первоначально, в дереве T_1 , узел u является корнем дерева, в последующих деревьях узел u определяется как ближайший предок суффикса $j_L + 1$, для которого существует суффиксная связь.

Для манипуляций с узлом u используем измененный вариант быстрого сканирования, который назовем *интеллектуальным сканированием*, а соответствующую функцию, которая будет реализовать такое сканирование, назовем *smartscan*. Подобно быстрому сканированию, функция $smartscan(u, v)$ определяет в дереве T_i строку $w = uv$, а также заменяет префиксный узел u ближайшим предком строки w (если узел w уже существует, то происходит формальная замена w на w). Таким образом, при интеллектуальном сканировании вычисляется пара строк:

$$(u, w) \leftarrow smartscan(u, v).$$

Как и в алгоритме МК, здесь используется тот факт, что если u является префиксным узлом для суффикса j_L , тогда суффиксная связь $s(u)$ указывает на префиксный узел $j_L + 1$. Каждое дерево T_i создается как связное дерево суффиксов с суффиксными связями, вычисленными для каждого внутреннего узла. Эти связи облегчают начальный поиск для суффикса $j_L + 1$, и любой последующий поиск для того же суффикса использует в интеллектуальном сканировании измененный префиксный узел.

Приведем в формальной форме алгоритм Укк.

Алгоритм 5.2.2 (Алгоритм Укконена)

```

▷ Онлайновый алгоритм вычисления связного дерева суффиксов  $T_x$ 
построение дерева  $T_1$     ▷ с суффиксной связью  $s(\text{корень}) = \text{корень}$ 
 $j_L \leftarrow 1$            ▷ в соответствии с леммой 5.2.4 узлы-листья
                           ▷ добавляются в порядке  $1, 2, \dots$ 
 $u \leftarrow \varepsilon$       ▷ начальным префиксным узлом является корень
for  $i \leftarrow 1$  to  $n$  do    ▷ преобразование  $T_i$  в  $T_{i+1}$ 
                           ▷ в цикле repeat формируется последний узел ветвления,
                           ▷ его суффиксная связь изменяется в зависимости от
                           ▷ значения переменной exit

 $w_{\text{пред}} \leftarrow \varepsilon$ 
 $exit \leftarrow \text{FALSE}$ 
repeat    ▷ создание, если необходимо, новых листьев
           ▷ локализация в  $T_i$  суффикса  $w = uv = x[j_L + 1..i]$ 
           ▷  $u$ , если необходимо, замена префиксного узла  $u$ 
            $(u, w) \leftarrow smartscan(u, v)$ 
    
```

```

if не существует пути от  $w$ , помеченного  $x[i + 1]$  then
    ▷ формируется новый лист (лемма 5.2.4)
     $j_L \leftarrow j_L + 1$ ; if  $j_L > i$  then  $exit \leftarrow \text{TRUE}$ 
    if узел  $w$  в  $T_i$  не существует then
        создание узла  $w$ 
        помечается ребро, ведущее к узлу  $w$ 
        помечается ребро, выходящее из узла  $w$ 
        добавляется лист  $j_L$ , помечается ? ребро от  $w$  до  $j_L$ 
        ▷ вычисляется связь для предыдущего узла
    if  $w_{\text{пред}} \neq \varepsilon$  then  $s(w_{\text{пред}}) \leftarrow w$ 
        ▷ для следующего выполнения smartscan
        ▷ определяется префиксный узел
     $u \leftarrow s(u)$ 
else
    ▷ суффикс  $j_L + 1$  в  $T_i$  расширяется до суффикса  $j_L + 1$ 
    ▷ в  $T_{i+1}$ , аналогичные действия выполняются
    ▷ для всех последующих суффиксов (лемма 5.2.4)
    ▷ вычисляется суффиксная связь для последнего узла
    if  $w_{\text{пред}} \neq \varepsilon$  then
         $s(w_{\text{пред}}) \leftarrow w$     ▷ для  $w$  существует узел
     $exit \leftarrow \text{TRUE}$ 
        ▷ выход из цикла, если  $j_L > i$  или нет возможности
        ▷ добавить новые листья
until  $exit$ 
    
```

Осталось прояснить еще одно “темное место” данного алгоритма. Если в дереве T_i создана последовательность листьев, суффиксную связь $s(w)$ родителя последнего листа $j_L + h$ в этой последовательности можно установить только тогда, когда следующий суффикс $j_L + h + 1$ не является листом. В этом случае $s(w)$ будет указывать на строку $w' = x[j_L + h + 1..i]$, которой соответствует узел в дереве T_i . Это утверждение основывается на следующих рассуждениях: строка $w = x[j_L + h..i]$ порождает путь, помеченный, например, меткой π , который не начинается с буквы $x[i + 1]$. Однако узел w' имеет не менее двух нисходящих путей: путь π и еще один путь, который начинается с буквы $x[i + 1]$. Эти рассуждения доказывают следующую теорему.

Теорема 5.2.5. Алгоритм 5.2.2 корректно вычисляет связанные деревья суффиксов $T_i, i = 1, 2, \dots, n + 1$. ■

Доказательство оценки временной сложности алгоритма 5.2.2 повторяет схему доказательства теоремы 5.2.3.

Теорема 5.2.6. Алгоритм 5.2.2 вычисляет связное дерево суффиксов T_x для заданной строки $x[1..n]$, определенной на упорядоченном алфавите A размером α , за время $O(n \log \alpha)$ с использованием памяти объемом $O(n\alpha)$.

Доказательство. Сначала заметим, что любая операция в алгоритме 5.2.2 выполняется не более $2n - 1$ раз: один раз для каждого значения i , плюс не более $n - 1$ раз при создании листьев в цикле **repeat**.

За исключением, возможно, операций интеллектуального сканирования и определения пути от узла w , каждый шаг алгоритма можно выполнить за константное время, используя для этого стандартные методы обхода дерева и его представления. Путь от узла w определяется за время $O(B)$, где, как и в доказательстве теоремы 5.2.3, B обозначает максимальное (по всем узлам дерева) время, необходимое для выбора правильного пути из текущего узла. Таким образом, если не учитывать операцию сканирования, то алгоритм выполняется за время $O(nB)$.

Теперь рассмотрим операцию интеллектуального сканирования. Для каждого значения j_L функция *smartscan* вызывается несколько раз, но после ее первого вызова выполняется только повторное вычисление префиксного узла u , для чего необходимо фиксированное время. Первоначально строка u полагается равной пустой строке, затем в процессе сканирования она изменяет свое значение, при этом сканирование не уменьшает ее длину, тогда как вычисление суффиксной функции s может уменьшить ее длину, но не более чем на единицу. Поэтому после всех сканирований длина строки u не может уменьшиться более чем на n . Следовательно, общее количество промежуточных узлов, просмотренных в процессе сканирования, также не превышает n . Таким образом, функция *smartscan* обрабатывает не более $2n$ узлов, при этом на каждый из них тратит не более $O(B)$ времени. Теперь осталось привлечь оценку времени B из доказательства теоремы 5.2.3 и получить требуемое утверждение данной теоремы. ■

В упражнении 5.2.9 предложено провести детальное сравнение алгоритмов МК и Укк. На практике алгоритм МК выполняется немного быстрее, чем алгоритм Укк [102]. С другой стороны, онлайн-природа алгоритма Укк позволяет динамически строить дерево суффиксов $T_{x_1 x_2 \dots x_k}$ по мере поступления на обработку строк x_1, x_2, \dots, x_k .

5.2.4 Алгоритм Фараха

Алгоритм Фараха (будем его для краткости обозначать как алгоритм Φ) представляет собой значительное открытие: это первый алгоритм построения дерева суффиксов, которое действительно выполняется за время $\Theta(n)$, при этом даже не требуется выполнения условия конечности алфавита. Такая эффективность достигается за счет того, что строковые последовательности определяются на индексированном алфавите (см. раздел 4.1) или, что эквивалентно, на целочисленном

алфавите $A = \{1, 2, \dots, \alpha\}$, при этом накладывается дополнительное условие, что $\alpha \in O(n)$. Такие алфавиты часто встречаются на практике.

Основу алгоритма Φ составляет алгоритм, описанный в работе [85], он кардинально отличается от других алгоритмов построения деревьев суффиксов. Мы опишем алгоритм Φ в виде пяти выполняемых шагов, что несколько отличается от оригинального описания метода в статье [86]. Чтобы сделать описание алгоритма более понятным, используем в качестве примера строку $g = 121112212221$, определенную на алфавите $A = \{1, 2\}$ (в этом примере $n = 12$).

Шаг I. Построение нечетного дерева $T_x^{\text{нч}}$.

Нечетное дерево $T_x^{\text{нч}}$ является деревом суффиксов для строки x , узлы-листья которого ограничены нечетными позициями $1, 3, 5, \dots$ строки $x\$$. Дерево $T_g^{\text{нч}}$ показано на рис. 5.5.

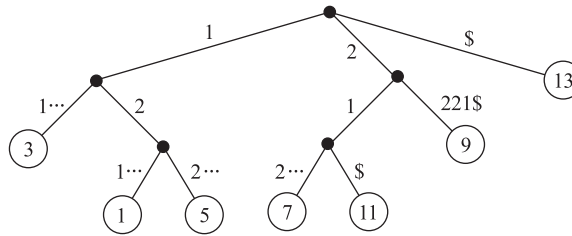


Рис. 5.5. Нечетное дерево для строки g

Основой эффективного построения дерева $T_x^{\text{нч}}$ за время $\Theta(n)$ служит поразрядная сортировка всех упорядоченных пар $\pi_i = (x[2i - 1], x[2i])$, $i = 1, 2, \dots, \lceil n/2 \rceil$, где для нечетного n последней парой будет $(x[n], \$)$ (символ $\$$ принимается как $(\alpha + 1)$ -я буква алфавита). Далее формируется новая строка x' из рангов $\rho(\pi_i)$ пар π_i . В нашем примере пары $(1, 2), (1, 1), (1, 2), (2, 1), (2, 2), (2, 1)$ имеют ранги соответственно 2, 1, 2, 3, 4, 3. Поэтому $g' = 212343$.

Для строки x' строится дерево суффиксов $T_{x'}$. Между этим деревом и деревом $T_x^{\text{нч}}$ существуют такие соответствия.

- Листьям i дерева $T_{x'}$ соответствуют листья $2i - 1$ в дереве $T_x^{\text{нч}}$.
- Внутренним узлам дерева $T_{x'}$, представляющим подстроки длиной k , соответствуют внутренние узлы дерева $T_x^{\text{нч}}$, представляющие подстроки длиной $2k$.

Дерево $T_{g'}$ показано на рис. 5.6.

На основе приведенных соответствий нетрудно показать, что дерево $T_{x'}$ можно преобразовать в дерево $T_x^{\text{нч}}$ за время $\Theta(n)$. Если обозначить через

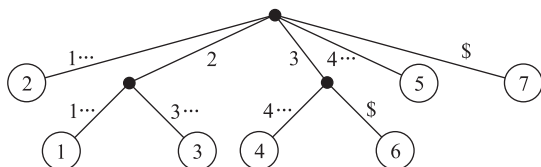


Рис. 5.6. Дерево суффиксов для строки g'

$t(n)$ максимальное (по всем возможным строкам x длиной n) время вычисления дерева T_x данным алгоритмом (выполнив все его пять шагов), тогда дерево $T_x^{\text{нч}}$ можно вычислить за время, не превышающее

$$t(n/2) + \Theta(n), \tag{5.4}$$

путем предварительного вычисления дерева $T_{x'}$ и последующего его преобразования в дерево $T_x^{\text{нч}}$. Если, кроме того, мы сможем показать, что время, требуемое для выполнения других четырех шагов алгоритма Φ , имеет порядок $O(n)$, тогда

$$t(n) - t(n/2) \in \Theta(n) \tag{5.5}$$

и, следовательно, $t(n) \in \Theta(n)$.

Шаг II. Построение четного дерева $T_x^{\text{чет}}$ на основе дерева $T_x^{\text{нч}}$.

Заметим следующее: если в строке x известен лексикографический порядок суффиксов, которые начинаются с нечетных позиций $2i + 1$, $i = 1, 2, \dots, \lceil n/2 \rceil - 1$, тогда можно определить лексикографический порядок суффиксов, которые начинаются с четных позиций $2i$, путем поразрядной сортировки первых элементов пар $(x[2i], 2i + 1)$, которую можно выполнить за $\Theta(n)$ время.

В нашем примере мы получим упорядочение 3, 1, 5, 7, 11, 9, 13 нечетных позиций в строке g путем обхода в прямом порядке дерева $T_x^{\text{нч}}$. Исключая позицию 1 и добавляя к этим числам значения $g[2i]$, $i = 1, 2, \dots, 6$, получим набор пар (2, 3), (1, 5), (2, 7), (2, 11), (1, 9), (1, 13), который надо отсортировать по первым элементам. После сортировки имеем (1, 5), (1, 9), (1, 13), (2, 3), (2, 7), (2, 11). Вычитая из вторых элементов пар первые, получим набор 4, 8, 12, 2, 6, 10, который является лексикографическим порядком четных суффиксов $g[2i..n + 1]$ и, следовательно, порядком листьев в дереве $T_x^{\text{чет}}$.

Как будет пояснено далее, теперь можно за время $\Theta(n)$ вычислить длины наименьших общих префиксов смежных листьев $2i$ и $2j$ в дереве $T_x^{\text{чет}}$:

$$\text{lcp}(2i, 2j) = |\text{LCP}(2i, 2j)|.$$

Теперь построить дерево $T_x^{\text{чет}}$ не представляет труда. Подобное дерево для строки g показано на рис. 5.7.

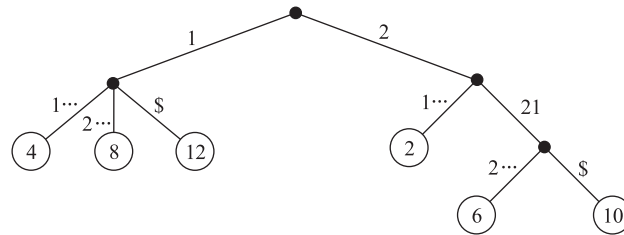


Рис. 5.7. Четное дерево для строки g

Шаг III. Слияние деревьев $T_x^{\text{нч}}$ и $T_x^{\text{чет}}$.

Далее необходимо найти эффективный способ слияния нечетного и четного деревьев в одно дерево T_x . Слияние будем производить начиная с корня. Предположим, что для каждого узла деревьев $T_x^{\text{нч}}$ и $T_x^{\text{чет}}$ выходящие из них ребра занесены в специальные списки, где они упорядочены в возрастающем лексикографическом порядке подстрок, которые представляют эти ребра. Алгоритм слияния деревьев просматривает только первые буквы подстрок, представленных ребрами деревьев $T_x^{\text{нч}}$ и $T_x^{\text{чет}}$, пусть это будут буквы $\lambda^{\text{нч}}$ и $\lambda^{\text{чет}}$. Тогда,

- если $\lambda^{\text{нч}} \neq \lambda^{\text{чет}}$, определяется поддереву, соответствующее меньшей из этих букв, и без изменений присоединяется к узлу-родителю;
- если $\lambda^{\text{нч}} = \lambda^{\text{чет}}$ и длины подстрок, представленных соответствующими ребрами, равны, в дерево слияния к текущему узлу добавляются два сына: один — из четного дерева, другой — из нечетного;
- если $\lambda^{\text{нч}} = \lambda^{\text{чет}}$ и длины подстрок, представленных соответствующими ребрами, различны, в дерево слияния к текущему узлу добавляются два узла, находящиеся на одном нисходящем пути, при этом ближайший узел будет соответствовать более короткой подстроке.

Если начать эту процедуру для корней нечетного и четного деревьев, далее она рекурсивно выполняется для корней всех поддеревьев, которые, возможно, уже содержат узлы из нечетного и четного деревьев, поскольку ранее мог быть реализован случай $\lambda^{\text{нч}} = \lambda^{\text{чет}}$. Так как время манипулирования с любым ребром этих деревьев фиксированно, то общее время слияния деревьев составит $\Theta(n)$.

В результате описанных действий получится дерево M_x , в котором будут присутствовать поддеревья, которые прошли процедуру слияния и которые ее избежали (т.е. были перенесены в дерево M_x без изменений). Теперь

встает задача преобразовать полученное дерево M_x в конечное дерево T_x . На рис. 5.8 показано полученное дерево слияния для строки g , где толстыми линиями выделены ребра, которые “прошли обработку” в процессе слияния. Отметим, что в результате слияния может получиться так, что некоторые суффиксы (как суффиксы 4 и 8 на рис. 5.8) будут представлены внутренними узлами дерева M_x .

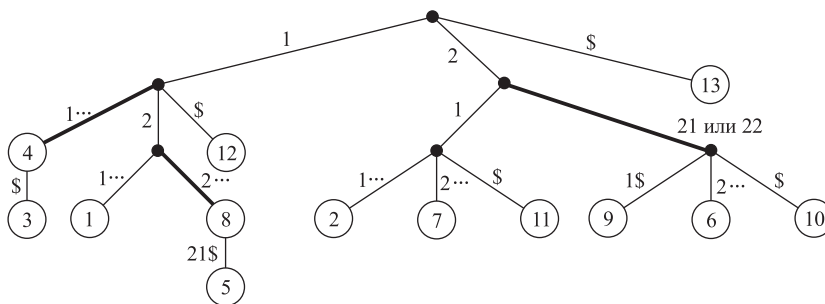


Рис. 5.8. Дерево слияния M_g для строки g

Шаг IV. Построение LCP-дерева.

В качестве прелюдии к преобразованию дерева M_x представим некоторые факты и простые утверждения, которые помогут нам в дальнейшем в построении дерева T_x . (Далее будем называть “слитой частью дерева M_x ” ту часть этого дерева, которая прошла обработку в процессе слияния нечетного и четного деревьев, и “неслитой частью” — необработанную часть этого дерева.)

- О1.** В процессе построения дерева M_x каждый его узел формируется на основе как дерева $T_x^{нч}$, так и дерева $T_x^{чет}$.
- О2.** Узел u принадлежит слитой части дерева M_x только тогда, когда он имеет потомков $2i$ и $2j - 1$ (возможно, один из них совпадает с самим узлом u), которые являются листьями в деревьях $T_x^{чет}$ и $T_x^{нч}$ соответственно.
- О3.** Пусть в дереве $M_x u = LCA(2i, 2j - 1)$, а в дереве T_x есть узел $u' = LCA(2i, 2j - 1)$, который в дереве M_x принадлежит поддереву с корнем u .⁵ Тогда u' является префиксом u . Следовательно, любая пара $(2i, 2j - 1)$ с наименьшим общим предком u в дереве M_x имеет наименьшим общим предком подстроку u' в дереве T_x .

⁵Напомним, что LCA обозначает “наименьший общий предок” (см. раздел 5.2.1). — Примеч. ред.

- О4.** Если узел $u = \text{LCA}(2i, 2j - 1)$ в дереве M_x не является корнем этого дерева, тогда $\text{LCA}(2i, 2j - 1) \neq \varepsilon$ и поэтому $u' = \text{LCA}(2i, 2j - 1)$ не является корнем дерева T_x .
- О5.** Пусть (i, i') — пара листьев дерева M_x , которые оба четные или нечетные. Тогда $\text{LCA}(i, i') = \text{LCP}(i, i')$, т.е. в этом случае выполняется равенство (5.1).
- О6.** С учетом фактов О2, О3 и О5 равенство $M_x = T_x$ будет выполняться только тогда, когда для любого внутреннего узла $u = \text{LCA}(2i, 2j - 1)$ дерева M_x выполняется условие $u = \text{LCP}(2i, 2j - 1)$.

Утверждение О6 говорит о том, что для корректного преобразования дерева M_x в дерево T_x необходимо проверить равенство (5.1) для каждой пары $(2i, 2j - 1)$ четных-нечетных листьев дерева M_x . Другими словами, надо вычислить $\text{LCP}(2i, 2j - 1)$ или, что то же самое, найти префикс строки $u = \text{LCA}(2i, 2j - 1)$ в дереве M_x длиной $\text{lcp}(2i, 2j - 1)$, который определит $u' = \text{LCA}(2i, 2j - 1)$ в дереве T_x . Заметим, что

$$\text{lcp}(2i, 2j - 1) = \text{lcp}(2i + 1, 2j) + 1, \quad \text{если } x[2i] = x[2j - 1], \quad (5.6)$$

и $\text{lcp}(2i, 2j - 1) = 0$ в противном случае. Но, согласно утверждению О4, $\text{lcp}(2i, 2j - 1) = 0$ только в том случае, когда $u = \text{LCA}(2i, 2j - 1)$ является корнем дерева M_x . Таким образом, для *всех* четных-нечетных пар $(2i, 2j - 1)$ в поддереве с корнем u (если u не является корнем дерева M_x) $\text{lcp}(2i, 2j - 1)$ можно вычислить на основе значения $\text{lcp}(2i + 1, 2j)$. Тогда, если в дереве M_x найдено $v = \text{LCA}(2i + 1, 2j)$, можно определить дерево, построенное на внутренних узлах дерева M_x , такое, что в этом дереве v будет родителем u . Если v не является корнем дерева M_x , то можно найти его родителя в дереве, подобном тому, в котором искался родитель u . Этот процесс можно продолжать до тех пор, пока родителем не окажется корень дерева M_x . Назовем дерево, которое определяется в процессе поиска родителей, **ЛСП-деревом**. Заметим, что $\text{lcp}(2i, 2j - 1)$ для всех четных-нечетных пар $(2i, 2j - 1)$ равняется глубине вершины $u = \text{LCA}(2i, 2j - 1)$ в ЛСП-дереве.

Ниже будет показано, что ЛСП-дерево можно вычислить за время $\Theta(n)$. ЛСП-дерево для строки $g = 121112212221$ приведено на рис. 5.9. Здесь, например, $\text{lcp}(5, 8) = 3$, поскольку узел 8 в этом дереве находится на глубине 3. Подобным образом вычисляются

$$\text{lcp}(3, 4) = \text{lcp}(1, 8) = \text{lcp}(2, 7) = \text{lcp}(2, 11) = \text{lcp}(6, 9) = \text{lcp}(9, 10) = 2.$$

Шаг V. Построение дерева T_x на основе дерева M_x и ЛСП-дерева.

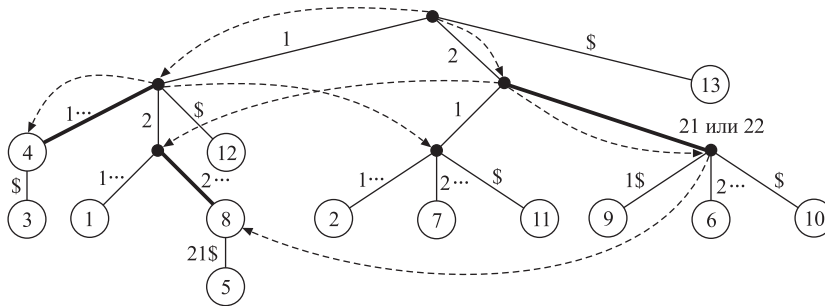


Рис. 5.9. Дерево M_g и LCP-дерево для строки g

Чтобы преобразовать дерево M_x в дерево T_x , сначала надо для каждого узла u в LCP-дереве проверить, будет ли его глубина $d(u)$ удовлетворять условию

$$d(u) = |u|, \tag{5.7}$$

где $|u|$ обозначает длину подстроки, определяемой путем в дереве M_x от корня до узла u . Глубину каждого узла можно определить во время обхода LCP-дерева, а $|u|$ для каждого узла можно вычислить в процессе формирования дерева M_x на основе нечетного и четного деревьев. Поскольку дерево M_x формируется путем слияния этих деревьев, то $d(u) \leq |u|$.

Предположим, что v — вершина, такая, что для всех ее предков u в LCP-дереве выполняется равенство (5.7). Если, кроме того, $d(v) = |v|$, то проверка закончена. Если же $d(v) < |v|$, то выполняются следующие действия.

P1. Обозначим через $p(v)$ родителя узла v в дереве M_x . Новый узел v' вставляется в дерево M_x как сын узла $p(v)$, если $|v'| = d(v)$. Фактически, $v' = \text{LCP}(2i, 2j - 1)$ для каждой четной-нечетной пары $(2i, 2j - 1)$, принадлежащей поддереву дерева M_x с корнем в узле v . Поэтому в дереве $T_x v'$ становится родителем поддеревьев из четного и нечетного деревьев с корнем v , которые сливаются на шаге III алгоритма, а узел v при этом удаляется.

P2. Узлы, подобные узлу v , представленному выше, обязательно появляются при слиянии четных и нечетных поддеревьев. Все четные и нечетные поддеревья, имеющие корнем узел v и слитые в дереве M_x , необходимо снова разделить. Такое разделение начинается с узла v' . Целесообразно применить итерационный процесс разделения, который несложно организовать, если в каждом узле дерева M_x хранить указатели на четные и нечетные поддеревья, которые слиты в этом узле. Тогда “правильное” поддерево с корнем в $p(v)$ можно сформировать так, как показано на

рис. 5.10. Здесь $\ell^{\text{чет}}$ (соответственно $\ell^{\text{нч}}$) — длина исходящего ребра⁶, ведущего к четному (нечетному) поддереву $T_x^{\text{чет}}$ (соответственно, $T_x^{\text{нч}}$), которое слито в дереве M_x с первоначальным корнем v .

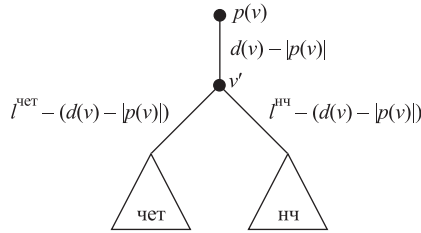


Рис. 5.10. Формирование поддерева с корнем $p(v)$

Чтобы лучше понять, как выполняются действия P1 и P2, рассмотрим для строки $g = 121112212221\$$ преобразование дерева слияния M_g , показанного на рис. 5.9, в дерево T_g , приведенное на рис. 5.11.

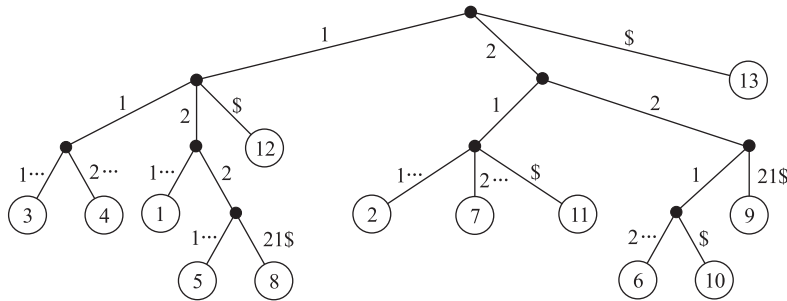


Рис. 5.11. Законченное дерево суффиксов T_g

Обозначим через v родителя узлов 9, 6 и 10 в дереве M_g ; как и выше, $p(v)$ обозначает родителя узла v . Поскольку в данном случае

$$2 = d(v) < |v| = 3,$$

вводится новый узел v' как сын узла $p(v)$ с ребром длиной

$$d(v) - |p(v)| = 1. \tag{5.8}$$

Последнее значение определяет длину префикса, который необходимо удалить из меток исходящих ребер сливаемых поддеревьев $T_x^{\text{чет}}$ и $T_x^{\text{нч}}$. Возвра-

⁶Здесь и далее длиной ребра называется длина метки этого ребра. — Примеч. ред.

щаясь к рис. 5.7 и 5.5, находим эти поддеревья (они отдельно показаны на рис. 5.12).

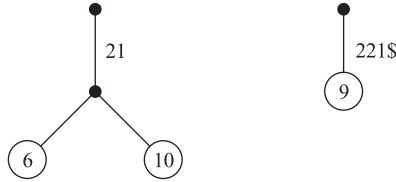


Рис. 5.12. Фрагменты четного и нечетного поддеревьев

В данном случае $\ell^{\text{чет}} = 2$, а $\ell^{\text{нч}} = 4$. Как видно на рис. 5.11, из меток исходящих ребер, согласно равенству (5.8), удаляется первая буква (буква 2), т.е. $\ell^{\text{чет}}$ и $\ell^{\text{нч}}$ уменьшаются на единицу. Поддеревья должны вставляться в лексикографическом порядке первых букв меток (подстрок) ребер нисходящего пути от узла v' . Поэтому на рис. 5.10 четное дерево расположено слева, а нечетное — справа.

Может случиться, что $\ell^{\text{чет}} = 1$ или $\ell^{\text{нч}} = 1$ (одновременно эти равенства выполняться не могут). В этом случае соответствующие исходящие ребра удаляются и корнем соответствующего поддерева становится узел v' .

Мы закончили поверхностный обзор алгоритма Φ — внимательный читатель обязательно найдет в этом обзоре много “темных мест”. Поэтому некоторые детали мы рассмотрим отдельно в следующем подразделе, но они будут, в основном, относиться к реализации алгоритма.

Реализация

Рассмотрим некоторые подробности выполнения алгоритма Φ , опущенные при общем его описании.

Шаг I. Построение нечетного дерева $T_x^{\text{нч}}$.

В описании этого шага остались неясными некоторые детали процесса сортировки, а также преобразования дерева $T_{x'}$ в дерево $T_x^{\text{нч}}$.

Чтобы выполнить сортировку целочисленных пар $\pi_i = (x[2i - 1], x[2i])$, необходимо дважды использовать поразрядную сортировку. Сначала происходит сортировка по вторым буквам $x[2i]$, а затем по первым буквам $x[2i - 1]$. Каждая из таких сортировок требует времени порядка $\Theta(\alpha + n/2)$. Как показано в упражнении 5.2.13, упорядочение рангов $\rho_i = \rho(\pi_i)$ этих отсортированных пар выполняется по дереву суффиксов строки $x' = \rho_1\rho_2 \dots \rho_{\lceil n/2 \rceil}$ и соответствует упорядочению листьев $2i - 1$ дерева $T_x^{\text{нч}}$. Аналогично каждый внутренний узел u дерева $T_{x'}$, определяющий подстроку x'

длиной $|u|$, соответствует узлу дерева $T_x^{нч}$, определяющему подстроку длиной $|2u|$, — отдельные буквы подстроки x' (эти буквы являются рангами) заменяются на соответствующие подстроки $x[2i - 1..2i]$ строки x . Реализация такой замены на примере дерева $T_{g'}$ (см. рис. 5.6) приводит к дереву $T_{g'_{3M}}$, показанному на рис. 5.13.

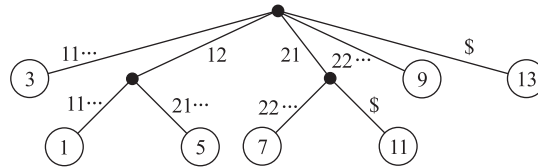


Рис. 5.13. Дерево $T_{g'_{3M}}$

Как видно из этого примера, дерево $T_{x'_{3M}}$ не обязательно будет синтаксическим деревом, поскольку в процессе замены может случиться, что выходящие из некоторого узла ребра будут заканчиваться на узлах, представляющих одинаковые буквы. В этом случае необходимо вернуться на предыдущий шаг. Но с другой стороны, одинаковыми могут быть только первые буквы на двух путях, начинающихся из одного узла, — если вторые буквы также совпадают, тогда в дереве $T_{x'}$ должно быть два пути с одинаковой первой буквой ρ , а это противоречит тому, что данное дерево является деревом суффиксов.

Обратное преобразование дерева $T_{x'_{3M}}$ в дерево $T_x^{нч}$ просто: в каждом внутреннем узле, предполагая, что исходящие из него ребра представлены в лексикографическом порядке, сравнивается первая буква метки ребра с первой буквой метки предшествующего ему ребра. И если эти буквы совпадают, вводится новое ребро с меткой длиной 1 (и, соответственно, новый узел на конце этого ребра). Как легко проверить, с помощью данной процедуры дерево $T_{g'_{3M}}$ (рис. 5.13) действительно преобразуется в дерево $T_{g^{нч}}$ (рис. 5.5). Построение дерева $T_{x'_{3M}}$ требует одного обхода дерева $T_{x'}$ с операцией замены во внутренних узлах. Замена выполняется за константное время, поэтому на такое преобразование необходимо время порядка $\Theta(n)$.

Шаг II. Построение четного дерева $T_x^{чет}$ на основе дерева $T_x^{нч}$.

Напомним, что выполнение второй поразрядной сортировки четных суффиксов строки x основано на известном лексикографическом порядке нечетных суффиксов этой строки. Для суффиксов четных номеров $2i$ и $2j$ можно получить аналог формулы (5.6):

$$\text{lcp}(2i, 2j) = \text{lcp}(2i + 1, 2j + 1) + 1, \quad \text{если } x[2i] = x[2j], \quad (5.9)$$

и $\text{lcp}(2i, 2j) = 0$ в противном случае. Тогда, зная $\text{lcp}(2i + 1, 2j + 1)$, значение $\text{lcp}(2i, 2j)$ можно вычислить за константное время. Поскольку дерево $T_x^{\text{нч}}$ является деревом суффиксов, для которого справедливо равенство (5.1), то $\text{LCA}(2i + 1, 2j + 1) = \text{LCP}(2i + 1, 2j + 1)$ для любой пары листьев. Известный алгоритм [112] позволяет найти наименьший общий предок для любых узлов в любом дереве за константное время, основываясь на предварительном линейном упорядочении номеров узлов дерева. Используя этот алгоритм, можно вычислить $\text{lcp}(2i, 2j)$ за константное время для каждой пары смежных суффиксов в лексикографически упорядоченном дереве $T_x^{\text{чет}}$. Но в этом дереве опять выполняется равенство (5.1), поэтому значение $\text{lcp}(2i, 2j)$ определяет $\text{LCA}(2i, 2j)$, т.е. имеем всю необходимую информацию для построения дерева $T_x^{\text{чет}}$.

Итак, используя $\Theta(n)$ времени для предварительной обработки дерева $T_x^{\text{нч}}$, дерево $T_x^{\text{чет}}$ можно построить за время $\Theta(n)$.

Шаг III. Слияние деревьев $T_x^{\text{нч}}$ и $T_x^{\text{чет}}$.

Этот шаг выполняется за время $\Theta(n)$.

Шаг IV. Построение LCP-дерева.

На этом шаге мы можем снова использовать алгоритм [112], вычисляющий наименьший общий предок для любой пары узлов за константное время. В качестве побочного продукта процесса построения дерева слияния M_x мы уже имеем соседние четные-нечетные (или нечетные-четные) пары $(2i, 2j - 1)$, которые являются листьями в деревьях $T_x^{\text{чет}}$ и $T_x^{\text{нч}}$ соответственно. Для каждой такой пары за константное время можно вычислить $u = \text{LCA}(2i, 2j - 1)$. Как указывалось ранее, родителем u в LCP-дереве является $v = \text{LCA}(2i + 1, 2j)$, что также вычисляется за константное время. Таких смежных неоднородных (четные-нечетные или нечетные-четные) пар, как показано в упражнении 5.2.15, всего n . Поэтому LCP-дерево строится за время, пропорциональное n .

В нашем примере в дереве M_x имеется $n + 1 = 13$ узлов, которые в дереве M_x появляются в следующем порядке: 4, 3, 1, 8, 5, 12, 2, 7, 11, 9, 6, 10, 13. Имеем 12 смежных неоднородных пар (4, 3), (4, 1), (1, 8), (8, 5), (5, 12), (5, 2), (2, 7), (2, 11), (2, 9), (9, 6), (9, 10), (10, 13), для которых надо вычислить наименьшие общие предки.

Шаг V. Построение дерева T_x на основе дерева M_x и LCP-дерева.

Как мы видели, обработка LCP-дерева сверху вниз позволяет определить узлы, исходящие ребра которых не удовлетворяют равенству (5.7). В таких узлах выполняются операции P1 и P2, которые требуют только ограниченного числа замены указателей на узлы. Таким образом, время выполнения шага V пропорционально количеству узлов в LCP-дереве, т.е. имеет порядок $O(n)$.

Корректность и эффективность

Итак, из предыдущего подраздела видно, что время $t(n)$, необходимое для выполнения первого шага алгоритма Φ , действительно имеет верхнюю границу (5.4), а другие шаги этого алгоритма требуют времени порядка $\Theta(n)$. Поэтому справедлива оценка (5.5), и мы имеем право сформулировать следующий формальный результат.

Теорема 5.2.7. Алгоритм Фарача вычисляет дерево суффиксов T_x для заданной строки $x[1..n]$, определенной на индексированном алфавите, за время порядка $\Theta(n)$. ■

Несмотря на асимптотическую оптимальность алгоритма Φ , остается неясным, когда наиболее целесообразно использовать его на практике. Хотя, как мы видели, основная идея алгоритма проста, в целом алгоритм весьма сложен для реализации, по сравнению с другими алгоритмами построения деревьев суффиксов, и, кроме того, требует значительного объема памяти. С другой стороны, как мы увидим на примере алгоритмов глав 7 и 8, хорошие алгоритмы могут порождать еще лучшие алгоритмы.

5.2.5 Применение и реализация

После того как мы рассмотрели в подробностях построение дерева суффиксов, будет более понятным, как это дерево можно использовать. Определить, является ли данная строка u подстрокой строки x , можно с помощью медленного сканирования (см. подраздел 5.2.1). Найти первое вхождение строки u в строку x несложно, если предварительно ребра дерева суффиксов пометить целочисленными метками и если все суффиксы строки x являются листьями в дереве T_x . Эффективно найти *все* вхождения строки u в строку x можно, используя стандартные способы обхода поддерева с корнем u (если u является узлом в дереве T_x) или обхода поддерева с корнями в сыновьях узлов, на которых оканчиваются ребра, определяющие префиксы строки u (если u не является узлом в дереве T_x). Если же заранее известно, что строка u является подстрокой строки x , то для нахождения первого вхождения строки u в строку x можно применить быстрое сканирование. В упражнениях 5.2.16 и 5.2.17 показаны детали подобного применения дерева суффиксов.

До сих пор мы не упоминали одно из основных свойств деревьев суффиксов, которое часто определяет возможность их использования: всем различным подстрокам строки x соответствуют свои пути в дереве T_x , начинающиеся в корне этого дерева. Поэтому реестр всех различных подстрок строки x можно получить как список всех подстрок, определенных на каждом таком пути. Даже если метки ребер имеют вид “ i_1, i_2 ”, что соответствует подстроке $x[i_1, i_2]$, в некоторых случаях получение реестра всех подстрок может потребовать времени порядка

$\Theta(n^2)$, поскольку, как показано в упражнении 1.2.10, строка x в общем случае может иметь C_n^2 различных подстрок. Но предположим, что требуется подсчитать только количество таких подстрок. Вследствие указанного свойства дерева суффиксов, количество различных подстрок, соответствующих любому пути в дереве T_x , равно количеству букв в метках ребер, составляющих этот путь, т.е. равно сумме величин $i_2 - i_1 + 1$, где сумма берется по всем меткам “ i_1, i_2 ” этого пути. Для каждого ребра, имеющего метку “ i_1, i_2 ”, назовем величину $i_2 - i_1 + 1$ *шириной* этого ребра, т.е. ширина ребра равна количеству букв, представленных его меткой. Тогда общее количество различных подстрок в строке x равно сумме значений ширины всех ребер дерева T_x . Используя стандартные способы обхода деревьев, эту сумму можно вычислить за время, пропорциональное количеству ребер в дереве T_x , т.е. пропорционально n . Таким образом, можем утверждать, что *если дерево суффиксов T_x построено*, то количество различных подстрок в строке x можно вычислить за время $\Theta(n)$.

Покажем еще одно применение дерева суффиксов. Пусть даны две строки x_1 и x_2 и необходимо вычислить наибольший общий фактор (подстроку) $\text{LCF}(x_1, x_2)$. (Эта задача поставлена в упражнении 2.2.17.) Используя свойства дерева суффиксов, которые показаны в процессе его построения, нетрудно показать, что данную задачу можно решить за время $\Theta(|x_1| + |x_2|)$, однако снова предполагая, что дерево суффиксов уже построено. Подробнее эта задача рассмотрена в упражнении 5.2.20.

Впечатляющая скорость выполнения алгоритмов, связанных с деревьями суффиксов, достигается не бесплатно, а требует определенных затрат, которые выражаются в необходимых объемах памяти. Как следует из теорем 5.2.3 и 5.2.6, в каждом узле дерева суффиксов время поиска составляет $O(\log \alpha)$, но только при условии, что в каждом узле хранятся данные, содержащие $O(\alpha)$ указателей и, возможно, другие дополнительные указатели и данные. Для индексированного алфавита поиск в каждом узле требует константного времени, но опять в каждом узле должна храниться информация объемом $O(\alpha)$. Кроме того, может потребоваться в каждом узле хранить указатель на родителя этого узла и указатели, необходимые для быстрого обхода дерева. Также необходимо $2 \log n$ бит для хранения меток каждого ребра. Таким образом, даже для малых α реализация дерева суффиксов невозможна без существенного, иногда недопустимо большого, объема памяти. Для очень малых алфавитов (например, $\alpha \leq 4$) дерево T_x можно эффективно реализовать как бинарное дерево, что исключает большие потребности в памяти. Как указывалось в обсуждении 2.1, использование деревьев суффиксов наиболее приемлемо в ситуациях, когда дерево T_x (или строка x) не изменяется, но используется часто, а объем алфавита α мал — этим условиям удовлетворяют, например, последовательности ДНК.

И последнее замечание: в принципе, для очень малых алфавитов на практике можно достичь теоретической эффективности алгоритмов, но такая эффек-

тивность все равно будет иллюзорной. Даже если для хранения каждой буквы тратится всего два бита (как в последовательности ДНК), то для хранения строки длиной n необходим массив объемом $n/4$ байт. Кроме того, для хранения дерева суффиксов, если в каждом узле сохраняются не менее двух указателей и дополнительная служебная информация, плюс метки ребер, потребуется не менее $10n$ байт. В действительности, “экономичная” реализация алгоритма Мак-Крейта [143] потребовала $20n$ байт памяти. На практике, особенно для больших значений n , время выполнения алгоритмов с деревьями суффиксов и требуемый для их реализации объем памяти могут значительно (в десятки раз) отличаться от *теоретических оценок* как в сторону увеличения, так и в сторону уменьшения, что обусловлено как обрабатываемыми данными, так и конкретной реализацией алгоритма под управлением конкретной операционной системы. Исследование этих вопросов проведено в работе [71].

Несколько пессимистические рассуждения этого раздела ведут к необходимости построения структур суффиксов, отличных от деревьев суффиксов. Две такие структуры описаны в следующем разделе.

Упражнения 5.2

1. Покажите, что для любого целого $i \in 1..n + 1$ $\text{tail}(i)$ реализуется в дереве T_x как конечный узел j , где $j = i + |\text{head}(i)|$.
2. Можно заметить, что для всех элементов табл. 5.1 выполняется неравенство

$$|\text{head}(i + 1)| \geq |\text{head}(i)| - 1.$$

Докажите это неравенство в общем случае (т.е. для любого положительного целого i). Докажите подобное неравенство для $\text{tail}(i)$.

3. Пусть u — узел дерева T_x (но не корень этого дерева) и пусть s обозначает суффиксную функцию. Покажите, что
 - а) $s(u)$ также является узлом дерева T_x ;
 - б) $s(u)$ не является узлом поддерева с корнем в узле u ;
 - в) если в дереве T_x узел v является родителем узла u , то или $s(v) = s(u) = \varepsilon$ или $s(v)$ является предком $s(u)$.
4. **Источником** называется узел v дерева T_x , такой, что для любого узла u этого дерева $v \neq s(u)$. Покажите, что
 - а) конечный узел 1 всегда будет источником;
 - б) если $|x| > 1$, тогда существует источник, который не является конечным узлом;

- в) если α' — количество различных букв в строке x (не считая символа \$), тогда существует ровно α' источников, которые не являются конечными узлами.
5. Приведите пример строки x , для которой алгоритм МК построит дерево T_x за время $\Theta(n^2)$, если в этом алгоритме, вместо быстрого сканирования, использовать медленное сканирование.
 6. Докажите лемму 5.2.4.
 7. Для строки $x = aabbacbbaa$ с помощью алгоритма Укк постройте дерево T_9 и опишите действия, необходимые для построения дерева T_{10} .
 8. В доказательствах теорем 5.2.3 и 5.2.6 предполагалось, что в каждом узле дерева используются для хранения букв упорядоченный массив или подходящее дерево поиска. Обсудите эти и другие структуры данных, которые можно применить для хранения информации в узлах дерева, и укажите ситуации, в которых их можно применить. Можно ли уменьшить требуемый объем памяти (соответственно, время выполнения) до величины порядка $\Theta(n)$, если алфавит неупорядочен, и если можно, то как это отразится на времени выполнения (соответственно, объеме памяти)?
 9. Очевидно, что алгоритм Укк более сложный в “интеллектуальном” смысле, чем алгоритм МК. Но каково соотношение их вычислительных сложностей? Предполагая, что оба алгоритма строят одинаковое дерево суффиксов и используют одинаковые структуры данных для представления информации в каждом узле дерева, что можно сказать о *реальном* времени их выполнения? Приведите примеры ситуаций, в которых свойство онлайнности алгоритма Укк является решающим преимуществом.
 10. Объясните, почему для доказательства эффективности алгоритма Φ необходимо условие, что $\alpha \in O(n)$.
Совет. Рассмотрите, как работает поразрядная сортировка.
 11. Используя соотношение (5.5), покажите, что для алгоритма Φ $t(n) \in O(n)$.
 12. Покажите, как две поразрядные сортировки можно использовать для сортировки в возрастающем порядке упорядоченных целочисленных пар $(x[2i - 1], x[2i])$.
 13. Покажите, что порядок суффиксов $2i - 1$, заданный деревом $T_x^{\text{нч}}$, совпадает с порядком рангов ρ_i , заданным деревом суффиксов строки $\rho_1\rho_2 \dots \rho_{\lceil n/2 \rceil}$.
 14. Покажите, что дерево суффиксов T_x можно построить, зная лексикографический порядок суффиксов строки x и порядок наименьших общих предков смежных пар.
 15. Покажите, что количество смежных четных-нечетных пар в последовательности из $n + 1$ целых чисел в точности равно n только тогда, когда эти числа не однородные (т.е. не все четные или не все нечетные).

16. Предложите способ определения реберных меток “ i_1, i_2 ” в процессе построения дерева суффиксов, который облегчал бы вычисление первого вхождения любой подстроки u в заданную строку x . Затем объясните, как быстрое сканирование находит первое вхождение подстроки u , если известно, что подстрока u действительно входит в строку x . Определите время, необходимое для выполнения быстрого сканирования.
17. Объясните, как, используя подходящие метки ребер и способ обхода дерева суффиксов, найти все вхождения заданной подстроки u в строку x . Определите необходимые для этого время и объем памяти.
18. Автор утверждал, что для алфавита размером в 4 буквы использование деревьев суффиксов может потребовать памяти объемом $20n$ байт (или даже $40n$ байт), где n — длина исходной строки. Поскольку он не доказал свое утверждение, оно может показаться пессимистическим или даже пораженческим. Для строк длиной $n = 100$ Кбайт, 1 Мбайт, 5 Мбайт и алфавита размером $\alpha = 4$ оцените минимальную требуемую память для подходящего представления дерева T_x .
19. *Деревом позиций* P_x [230] называется дерево суффиксов T_x , измененное таким образом, что в нем лист i представляет самый короткий префикс суффикса $x[i..n]$, который при этом не совпадает с другими подстроками строки x .
 - а) Покажите, что i является листом дерева P_x только в том случае, когда $i \leq j_L$, другими словами, только в том случае, когда i является листом дерева T_x .
 - б) Покажите, что дерево P_x имеет n листьев только тогда, когда буква $x[n]$ входит в строку x только один раз.
 - в) Предложите простой алгоритм построения дерева P_x на основе дерева T_x .
20. Предложите алгоритм на основе дерева суффиксов для вычисления наибольшего общего фактора для двух заданных строк x_1 и x_2 за время, пропорциональное $|x_1| + |x_2|$. Покажите, что ваш алгоритм можно расширить для обработки $k \geq 2$ строк.

Совет. Сформируйте строку $x = x_1\$1x_2\2 с использованием разных символов окончания строк $\$1$ и $\$2$. Заметьте, что $\text{LCF}(x_1, x_2)$ в дереве T_x соответствует узел u наибольшей длины, такой, что поддереву с корнем в узле u содержит символ $\$1$.

5.3 Альтернативные структуры для представления суффиксов

В этом разделе мы рассмотрим две структуры, альтернативные деревьям суффиксов. Обе эти структуры имеют перед деревьями суффиксов некоторые теоретические преимущества, которые должны проявиться на практике. Хотя описываемые здесь структуры различны между собой, они имеют общее свойство, которое заключается в том, что обе они строятся на основе деревьев суффиксов.

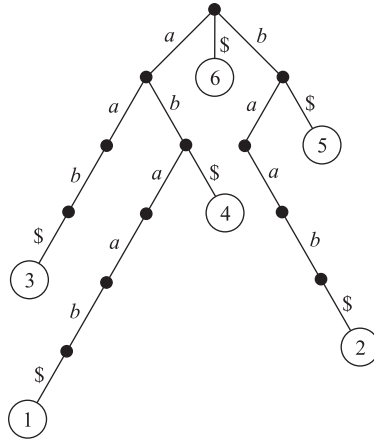
5.3.1 Ориентированные ациклические графы слов

Необходимость введения структур, отличных от деревьев суффиксов, вытекает из простого наблюдения, которое читатель уже должен был сделать: дерево суффиксов представляет по отдельности все одинаковые повторяющиеся подстроки, что кажется излишним, — узлам, представляющим повторяющиеся подстроки, соответствуют ребра с одинаковыми метками, и различие между ними проявляется только в метках узлов-листьев. Например, на рис. 5.2 поддерево T'_v , соответствующее подстроке $v = aababaabaab\$$, встречается дважды, присоединенное к узлам b и ab ; поддерево T'_w , соответствующее подстроке $w = abaabaab\$$, встречается три раза, присоединенное к узлам $aaba$, $baaba$ и $abaaba$. Подобным образом на рис. 5.4 в дереве суффиксов T_6 поддерево $T'_{aab\$}$ встречается дважды, присоединенное к узлам b и ab . Поэтому ориентированные ациклические графы слов (сокращенно, ОАГС) вводятся как средство для избежания таких излишеств [35, 70].

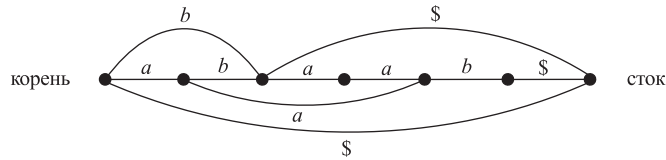
Но прежде чем приступить к непосредственному описанию ОАГС, укажем способ определения наличия в дереве суффиксов T_x одинаковых поддеревьев. В описании алгоритма Укконена для простоты мы предполагали, что $x = x[1..n + 1]$, где $x[n + 1] = \$$. Поскольку, как будет показано далее, ОАГС можно легко получить путем преобразования “неплотных” деревьев суффиксов, мы на время предположим, что каждая подстрока u строки x представляется в дереве T_x отдельным узлом, даже если такой узел имеет только одного сына. Пример неплотного дерева суффиксов для строки $x = abaab\$$ показан на рис. 5.14, *a*.

Для каждой непустой подстроки u строки x определим множество Π_u всех позиций в строке x , где оканчиваются подстроки u . Назовем такое множество **множеством окончаний**. По определению положим $\Pi_\epsilon = \{1, 2, \dots, n + 1\}$. С помощью множества окончаний можно отобразить структуру дерева суффиксов T_x . Пусть u — непустая подстрока строки x , которая входит в строку x k раз в позициях j_1, j_2, \dots, j_k . Эти номера позиций являются метками листьев в дереве T_x и представляют префиксы подстроки u . Множество окончаний для подстроки u вычисляется следующим образом:

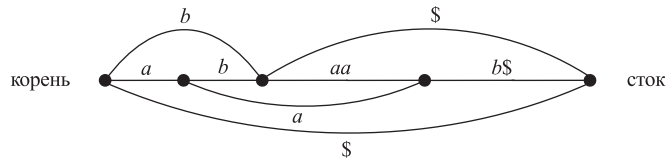
$$\Pi_u = \{j_1 + |u| - 1, j_2 + |u| - 1, \dots, j_k + |u| - 1\}.$$



а) Неплотное дерево суффиксов



б) Граф D_x (7 узлов, 10 дуг)



в) Граф D_x (5 узлов, 8 дуг)

Рис. 5.14. Преобразование дерева суффиксов в ОАГС для строки $x = abaab\$$

В случае, когда $u = \varepsilon$, множество Π_ε является множеством меток листьев, которые представляют префиксы ε , т.е. множеством меток всех листьев. Отметим, что в случае $u \neq \varepsilon$ множество Π_u является подмножеством множества Π_ε .

Рассмотрим несколько примеров. Для одиночной буквы λ Π_λ представляет множество всех позиций в строке x , в которых встречается эта буква (и, конечно, это — множество меток листьев, имеющих префикс λ). Если $u = x[j..n + 1]$ для всех j , $1 \leq j \leq n + 1$, тогда $\Pi_u = \{n + 1\}$, показывая тем самым, что подстроки

строки x могут быть префиксами только одного суффикса $x[j..n+1]$. В общем случае, если $x[j]$ — буква, которая входит в строку x только один раз, тогда для произвольной подстроки u , содержащей эту букву, $|\Pi_u| = 1$. В последнем примере подстроке u соответствует или лист дерева T_x , или узел, имеющий только одного сына.

Для построения ориентированных ациклических графов слов практический интерес представляют подстроки u_1 и u_2 , имеющие одинаковые множества окончаний. В этом случае эти подстроки в дереве T_x представимы узлами, которые являются корнями одинаковых поддеревьев. Если $\Pi_{u_1} = \Pi_{u_2}$, то в этом случае будем говорить, что подстроки u_1 и u_2 *Π-эквивалентны*. В упражнении 5.3.1 доказывается, что Π-эквивалентность действительно является отношением эквивалентности.

Лемма 5.3.1. Пусть u_1 и u_2 — две непустые подстроки строки x , такие, что $|u_1| \leq |u_2|$. Тогда подстроки u_1 и u_2 будут Π-эквивалентными только тогда, когда подстрока u_1 встречается в строке x только как суффикс подстроки u_2 .

Доказательство. Если $\Pi_{u_1} = \Pi_{u_2}$, тогда по определению множеств окончаний подстроки u_1 и u_2 оканчиваются в одинаковом множестве позиций. И поскольку $|u_1| \leq |u_2|$, то u_1 является суффиксом u_2 . Обратное: если u_1 встречается в строке x только как суффикс u_2 , тогда множество окончаний для u_1 должно совпадать с аналогичным множеством подстроки u_2 . ■

Можно интерпретировать этот результат в терминах знакомой нам суффиксной функции, введенной в разделе 5.2.1. Пусть s^k ($k \geq 1$) обозначает последовательное применение k раз функции s , т.е. если $s(u_3) = u_2$ и $s(u_2) = u_1$, тогда $s^2(u_3) = u_1$. Лемма 5.3.1 говорит о том, что если подстроки u_1 и u_2 Π-эквивалентны и $|u_1| \leq |u_2|$, тогда существует такое целое положительное число k , что $s^k(u_2) = u_1$. Таким образом, с помощью суффиксной функции можно определить число вхождений Π-эквивалентных подстрок в строку x .

Лемма 5.3.2. Пусть u_1 и u_2 — две непустые Π-эквивалентные подстроки строки x , такие, что $|u_1| \leq |u_2|$. Тогда для некоторого целого положительного числа k $s^k(u_2) = u_1$ и подстрока u_1 будет Π-эквивалентной всем подстрокам $s^i(u_2)$, $i = 1, 2, \dots, k$.

Доказательство. Пусть $u = s^i(u_2)$ для некоторого $i \in 1..k$. Вследствие леммы 5.3.1, подстрока u_1 может входить в строку x только как суффикс подстроки u_2 . Тогда подстрока u также должна входить в строку x только как суффикс подстроки u_2 . Но тогда, вследствие леммы 5.3.1, подстрока u будет Π-эквивалентной подстроке u_2 и, следовательно, Π-эквивалентной подстроке u_1 . ■

Далее покажем, что множества окончаний не могут *частично* перекрываться. Отсюда будет следовать, что эти множества соответствуют структуре дерева.

Лемма 5.3.3. Пусть u_1 и u_2 — две непустые подстроки строки x , такие, что $|u_1| \leq |u_2|$. Тогда может выполняться только одно из следующих утверждений:

- а) $\Pi_{u_2} \cap \Pi_{u_1} = \emptyset$;
- б) $\Pi_{u_2} \subset \Pi_{u_1}$.

Доказательство тривиально, когда $u_1 = \varepsilon$. Поэтому предположим, что подстрока u_1 не пустая. Пусть множества Π_{u_1} и Π_{u_2} имеют общий элемент. Этот элемент определяет конечную позицию как подстроки u_1 , так и подстроки u_2 . Поэтому подстрока u_1 является суффиксом подстроки u_2 . Но тогда каждый элемент множества Π_{u_2} должен входить в множество Π_{u_1} . ■

Теперь можно “связать” понятие Π -эквивалентности с наличием одинаковых (изоморфных) поддеревьев в дереве T_x .

Лемма 5.3.4. Пусть u_1 и u_2 — узлы дерева суффиксов T_x . Поддеревья T'_{u_1} и T'_{u_2} дерева T_x будут изоморфны тогда и только тогда, когда подстроки u_1 и u_2 будут Π -эквивалентны.

Доказательство тривиально, когда $u_1 = u_2$. Поэтому полагаем, что $u_1 \neq u_2$. Поскольку пустая строка ε Π -эквивалентна сама себе, также полагаем, что подстроки u_1 и u_2 не пусты.

Пусть $\Pi_{u_2} = \Pi_{u_1}$. Без потери общности можно считать, что $|u_1| \leq |u_2|$. Тогда, вследствие леммы 5.3.1, подстрока u_1 может входить в строку x только как суффикс подстроки u_2 . В таком случае поддерево T'_{u_2} с корнем в узле u_2 является поддеревом для суффикса u_1 подстроки u_2 и поэтому должно совпадать с поддеревом T'_{u_1} .

Теперь пусть совпадают поддеревья T'_{u_1} и T'_{u_2} . Тогда в строке x существует суффикс $u_1x[j..n+1]$ только в том случае, когда существует суффикс $u_2x[j..n+1]$. Для каждого такого целого j подстроки u_1 и u_2 заканчиваются в одной и той же позиции $j-1$. Поскольку эти подстроки не имеют окончаний в других позициях, отличных от позиций j , то отсюда следует, что $\Pi_{u_2} = \Pi_{u_1}$. ■

Теперь мы можем дать четкое определение ориентированному ациклическому графу слов для данной строки x (будем обозначать такой граф как D_x) и показать, как он строится на основании “неплотного” дерева суффиксов T_x . Иллюстрацией к процессу построения графа D_x послужит рис. 5.14. Сначала опишем этапы построения графа D_x .

1. Для всех Π -эквивалентных узлов u_1, u_2, \dots, u_k произведем слияние поддеревьев $T'_{u_i}, i = 1, 2, \dots, k$ в единый подграф. Эту операцию можно выполнить путем слияния узлов в каждом поддереве.

Поскольку, как мы видели, $\Pi_{x[j..n+1]} = \{n+1\}$ для любого листа $j \in 1..n+1$, поэтому все листья Π -эквивалентны. Вследствие леммы 5.3.4, на этом шаге все листья сольются в один узел, который назовем *стоком*.

На рис. 5.14 можно найти четыре множества Π -эквивалентных узлов: множество всех листьев, $\{aa, baa, abaa\}$, $\{ba, aba\}$ и $\{b, ab\}$.

2. Дуги графа ориентированы от корня в направлении, соответствующем отношению “родитель-сын” в дереве T_x . На рис. 5.14 ориентация дуг не показана: она совпадает с направлением слева направо.

Дополнительный третий шаг, показанный на рис. 5.14, *в*, производит “зачистку” графа D_x — удаляются узлы, имеющие единичные полустепени входа и исхода.⁷ Отметим, что для хранения даже первоначального графа D_x хватает линейной (относительно n) памяти.

Отметим, что построенный ОАГС по буквенно представляет те же $n+1$ суффиксы, что и соответствующее дерево суффиксов. И так же, как в дереве суффиксов, в графе метки на исходящих из какого-либо узла дугах попарно различны. ОАГС имеет в точности $n+1$ путей от корня до стока длиной $1, 2, \dots, n+1$, которые можно проверить теми же способами, что и пути в дереве суффиксов. Очевидно, что граф D_x ацикличен, поскольку все пути от корня до стока ориентированы согласно ориентации “родитель-сын” в дереве T_x .

Это определение ОАГС, данное в полуалгоритмической форме, на самом деле не является алгоритмом. Алгоритм построения ОАГС обязательно будет зависеть от метода “разуплотнения” дерева T_x , который может потребовать памяти порядка $\Theta(n^2)$, что может быть нежелательным. Мы оставим “на потом” вопрос об эффективном построении ОАГС, сначала разберемся с требованиями, предъявляемыми к объему памяти, необходимому для хранения ОАГС.

Теорема 5.3.5. Для строки $x = x[1..n]$ длиной $n+1$ граф D_x имеет $N \leq 2n+1$ узлов и не более $N+n-1 \leq 3n$ дуг.

Доказательство. Каждый узел графа D_x соответствует одному или нескольким множествам окончаний подстрок, определяемых узлами дерева T_x . Лемма 5.3.3 говорит о том, что эти множества не могут частично пересекаться. Также очевидно, что для корня $\Pi_\varepsilon = \{1, 2, \dots, n+1\}$. Узлы графа D_x должны соответствовать узлам поддеревьев с корнями в узлах $\{1, 2, \dots, n+1\}$, у которых потомки соответствуют непустым непересекающимся множествам, объединение которых включается в узел родителя. В упражнении 5.3.4 предложено показать, что такое дерево содержит не более $2n+1$ узлов.

⁷Полустепенью входа (исхода) узла (вершины) ориентированного графа называется количество дуг, входящих в данный узел (выходящих из данного узла). — *Примеч. ред.*

Теперь рассмотрим остовное дерево⁸ ST_x графа D_x , который показывает путь от корня до стока наибольшей длины, помеченный буквами строки x , как показано на рис. 5.15. Остовное дерево ST_x имеет N узлов и $N - 1$ ребер. Теперь в это остовное дерево будем по одной добавлять дуги графа D_x . Поскольку добавление каждой такой дуги также добавляет не менее одного нового пути, а всего в графа D_x путей ровно n , то отсюда следует, что общее количество дуг в графе D_x не более $(N - 1) + n$. ■

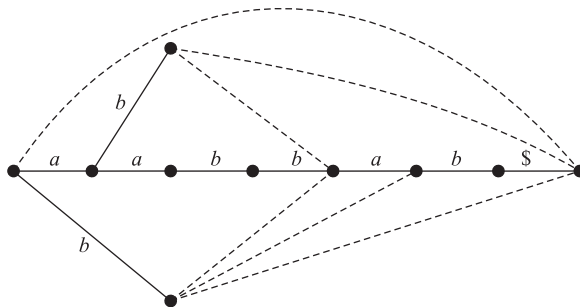


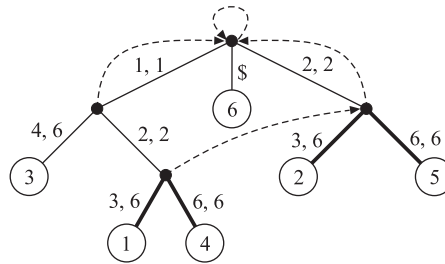
Рис. 5.15. Остовное дерево графа D_x для $x = aabbab$ ($N = 10, n + 1 = 7$)

Из последней теоремы вытекает искомый результат, что для хранения ОАГС D_x требуется память, объем которой пропорционален длине строки x . Труднее разработать алгоритм построения графа D_x , который выполнялся бы за время, пропорциональное $|x|$. В решении последней проблемы может помочь суффиксная функция s , обладающая полезными свойствами, установленными в лемме 5.3.2. Эти свойства позволяют установить изоморфизмы поддеревьев в дереве суффиксов T_x и, следовательно, П-эквивалентность соответствующих подстрок.

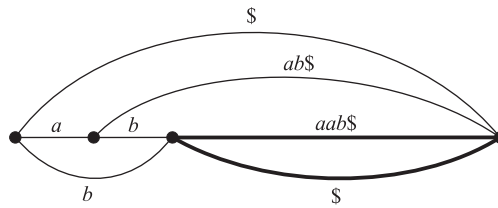
Предложим следующий алгоритм преобразования дерева T_x в граф D_x , основанный на использовании суффиксной функции. Кратко этот алгоритм можно представить в виде двух последовательных шагов.

1. На первом шаге на основе суффиксных связей дерева T_x строится **ориентированный ациклический граф** (сокращенно, — ОАГ; этот граф может не совпадать с ОАГС). Процесс создания ОАГ для строки $x = abaab$ показан на рис. 5.16, *a* и *б*.
2. На втором шаге минимизируется ширина (см. определение в разделе 5.2.5) входящих дуг для каждого узла ОАГ. Этот процесс показан на рис. 5.16, *a* и *б*.

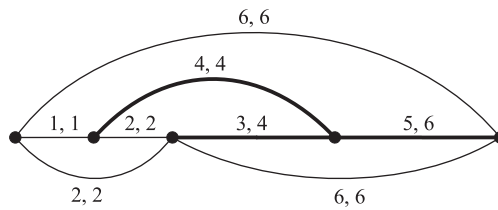
⁸Остовным деревом графа называется минимальный (по количеству ребер) связный ациклический граф, содержащий все вершины исходного графа. — Примеч. ред.



а) Суффиксные связи в дереве суффиксов
(толстыми линиями обозначены совпадающие поддеревья)



б) ОАГ (толстыми линиями
обозначены слитые поддеревья)



в) ОАГС (толстыми линиями
обозначены измененные дуги)

Рис. 5.16. Преобразование дерева суффиксов
в ОАГС (см. рис. 5.14)

Следующая лемма обосновывает первый шаг алгоритма.

Лемма 5.3.6. Пусть u_1 и u_2 — два узла дерева суффиксов T_x , такие, что $|u_1| \leq |u_2|$. Тогда поддеревья T'_{u_1} и T'_{u_2} изоморфны (будут совпадать), если выполняются условия:

- а) эти поддеревья содержат одинаковое количество листьев;
- б) $s^k(u_2) = u_1$ для некоторого целого положительного числа k .

Доказательство тривиально, когда $u_1 = \varepsilon$, поскольку в этом случае $T'_\varepsilon = T_x$. Поэтому будем считать, что $u_1 \neq \varepsilon$.

Предположим, что поддеревья T'_{u_1} и T'_{u_2} изоморфны. Тогда, вследствие леммы 5.3.4, подстроки u_1 и u_2 П-эквивалентны и, как результат леммы 5.3.2, существует целое число k , такое, что $s^k(u_2) = u_1$. Тогда, как показано в упражнении 5.3.3, поддеревья T'_{u_1} и T'_{u_2} имеют одинаковое количество листьев. Это доказывает достаточность утверждения леммы.

Теперь предположим, что выполняются условия *а)* и *б)* леммы. Обозначим через v_u количество листьев в поддереве T'_u . Заметим, что v_u равно количеству вхождений подстроки u в строку x . По условию *б)* леммы u_1 является суффиксом подстроки u_2 , т.е. каждый лист в поддереве T'_{u_2} соответствует листу в поддереве T'_{u_1} . По условию *а)* леммы $v_{u_2} = v_{u_1}$, поэтому подстрока u_1 может входить в строку x только как суффикс подстроки u_2 . Тогда, вследствие леммы 5.3.1, $\Pi_{u_2} = \Pi_{u_1}$, теперь из леммы 5.3.4 следует, что поддеревья T'_{u_1} и T'_{u_2} совпадают. ■

Подчеркнем, что лемма 5.3.6 справедлива только в случае, когда u_1 и u_2 являются листьями дерева суффиксов. Поэтому процесс слияния одинаковых поддеревьев включает слияние всех листьев дерева T_x в узел-сток графа D_x . Все другие изоморфные поддеревья можно локализовать, если следовать суффиксным связям в дереве T_x . Такие одинаковые поддеревья будут сливаться только в том случае, когда они содержат одинаковое количество листьев. Процесс слияния можно реализовать за время, пропорциональное величине n , если выполнить предварительную обработку дерева T_x , как показано ниже в алгоритме 5.3.1. (Эту процедуру предварительной обработки можно выполнить как “побочный продукт” процесса построения самого дерева T_x .)

Алгоритм 5.3.1

▷ Преобразование дерева суффиксов T_x в $OAGCD_x$

▷ Предварительная обработка дерева T_x

обход дерева T_x :

а) сохранение в каждом узле u количества листьев в поддереве T'_u

б) сохранение узла-источника для каждой цепочки суффиксных связей

▷ Шаг 1: построение OAG из дерева T_x

for каждый сохраненный узел-источник u **do**

▷ Проверка цепочки суффиксных связей до корня

while $s(u) \neq \varepsilon$ **do**

if T'_u и $T'_{s(u)}$ имеют одинаковое количество листьев **then**

слияние T'_u и $T'_{s(u)}$

$u \leftarrow s(u)$

▷ Шаг 2: построение D_x на основе OAG

for каждая вершина v **do**

определение входящей дуги (u, v) наибольшей ширины $|w|$
for каждая входящая дуга (u', v) ширины $1 < |w'| < |w|$ **do**
 удаление дуги (u', v)
 создание дуги (u', v') \triangleright как показано на рис. 5.17

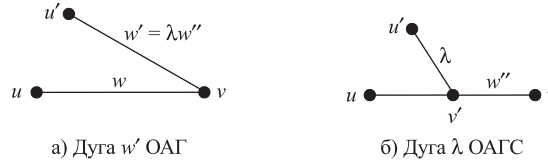


Рис. 5.17. Преобразование входящих дуг ОАГ в дуги ОАГС

Как видно на рис. 5.16, б и в, граф ОАГ может не быть графом ОАГС: ОАГ может иметь дуги (подобные дуге с меткой ab на рис. 5.16, б), метки которых можно определить как суффиксы меток на других дугах (на том же рисунке метка aab имеет суффикс ab). Эта проблема порождена тем, что используется дерево T_x , представленное в компактной форме, поэтому некоторые изоморфные поддеревья (такие, как ребро b на рис. 5.14, а) могут быть не найдены.

На втором шаге алгоритма 5.3.1 находятся “пропавшие” суффиксы путем просмотра входящих дуг для каждого узла v ОАГ. Сначала определяется дуга (u, v) с меткой w наибольшей длины и затем одна за другой просматриваются другие дуги с метками w' длиной $2 \leq |w'| \leq |w|$. Если w' является суффиксом w , то такая дуга удаляется, а вместо нее добавляется дуга (u', v') с меткой единичной длины, как показано на рис. 5.17.

Приведенные выше рассуждения позволяют сформулировать следующую теорему.

Теорема 5.3.7. Алгоритм 5.3.1 корректно вычисляет ОАГС D_x для данной строки x на основе ранее построенного дерева суффиксов T_x . ■

Очевидно, что время выполнения алгоритма 5.3.1 должно каким-то образом зависеть от размера алфавита, поскольку этот размер влияет на построение дерева суффиксов. Но поскольку мы предполагали, что дерево суффиксов уже построено, преобразование дерева T_x в граф D_x можно выполнить за время, линейное относительно $|x|$.

Теорема 5.3.8. Алгоритм 5.3.1 преобразует дерево суффиксов T_x для данной строки x в ОАГС D_x за время порядка $\Theta(n)$.

Доказательство. В алгоритме 5.3.1 предварительная обработка дерева суффиксов требует посетить каждый узел дерева T_x для вычисления количества листьев

в поддереве T'_u , для чего можно использовать стандартный способ обхода дерева в обратном порядке. Во время обхода дерева T_x можно сформировать список узлов, которые или являются листьями или указывают на суффиксную связь. Поскольку дерево T_x имеет $\Theta(n)$ узлов (метки ребер не исследуются на этом этапе), то при использовании стандартного способа обхода дерева предварительная обработка потребует времени порядка $\Theta(n)$.

На первом шаге алгоритма проверяется каждый узел каждой цепочки суффиксных связей от источника до корня. Если на этапе предварительной обработки получен список всех узлов, от которых начинается цепочка суффиксных связей, то “слияние” поддеревьев требует только константного времени для переустановки указателей. Поскольку количество узлов с суффиксными связями имеет порядок $\Theta(n)$, выполнение первого шага алгоритма также требует времени порядка $\Theta(n)$.

На втором шаге обработка каждой дуги графа ОАГ также требует только константного времени. Количество дуг в ОАГ равно количеству дуг в ОАГС, по теореме 5.3.5 дуг в ОАГС всего $\Theta(n)$. Поэтому общее время выполнения второго шага также имеет порядок $\Theta(n)$. ■

Отметим, что в недавних работах [68, 127] описаны алгоритмы, которые вычисляют ОАГС непосредственно, без предварительного вычисления дерева суффиксов.

Недавние исследования показали [143, 127, 62], что при применении соответствующих методов ОАГС требует значительно меньше памяти для его хранения, чем соответствующие деревья суффиксов. Но основное преимущество ОАГС заключается не в этом, а в том, что они предлагают другой взгляд на структуру суффиксов строковых последовательностей, что показывают леммы 5.3.1–5.3.4 и 5.3.6. В следующем разделе рассмотрим еще одну структуру суффиксов, которая также более экономична по сравнению с деревьями суффиксов и, кроме того, не зависит от размера алфавита.

5.3.2 Массивы суффиксов

Как обычно, предполагаем, что задана строковая последовательность $x = x[1..n]$. Как и выше во всех многочисленных обсуждениях деревьев суффиксов, полагаем, что каждый суффикс $x[j..n]$ строки x ($j = 1, 2, \dots, n$) определяется просто номером j листа в дереве T_x . Тогда **массив суффиксов** (или **суффиксная строка!**) $\sigma = \sigma_x$ — это просто массив n целых чисел $j = \sigma[i]$, определяющих суффиксы, упорядоченные в возрастающем лексикографическом порядке (т.е. i — номер суффикса $x[j..n]$ в этом упорядоченном множестве всех суффиксов строки x) [171, 172]. Например, строка

$$g = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ a & b & a & a & b & a & a & b, \end{matrix}$$

введенная в разделе 2.1, имеет следующие суффиксы, упорядоченные в лексикографическом порядке.

6	<i>aab</i>	1	<i>abaabaab</i>
3	<i>aabaab</i>	8	<i>b</i>
7	<i>ab</i>	5	<i>baab</i>
4	<i>abaab</i>	2	<i>baabaab</i>

Поэтому для данной строки массив суффиксов имеет вид $\sigma_g = 63741852$.

Отметим одно почти очевидное, но очень важное свойство массива суффиксов: все суффиксы строки x , имеющие одинаковый префикс, располагаются последовательно в массиве σ . Это свойство незаменимо для поиска всех вхождений подстроки u в строку x — всем этим вхождениям в массиве σ соответствует ряд последовательных элементов. Если в нашем примере положить $u = ab$, то этой подстроке будут соответствовать элементы массива $\sigma[3..5]$. Если $u = baa$, тогда этой подстроке соответствуют элементы $\sigma[7..8]$.

Имея дерево суффиксов T_x , массив σ_x можно вычислить путем обхода дерева T_x в глубину, когда для каждого узла просматриваются в лексикографическом порядке первые буквы меток ребер, исходящих из данного узла. Если символу \$ окончания строки присвоить наименьшее лексикографическое значение, тогда в процессе такого обхода дерева T_x все листья этого дерева будут посещаться в лексикографическом порядке. Таким образом, мы почти доказали следующую теорему.

Теорема 5.3.9. Массив суффиксов σ_x для данной строки $x = x[1..n]$, определенной на алфавите A размером α , можно вычислить на основе ранее построенного дерева суффиксов T_x за время порядка $O(n \log \alpha)$. ■

В зависимости от вида структуры данных, реализованной в каждом узле T_x , время вычисления массива σ_x в определенных ситуациях действительно может достигать верхней границы $\Theta(n)$. Более подробно эта тема рассмотрена в упражнении 5.3.7. В работе [172] показано, что среднее (по всем возможным строкам x) время вычисления массива σ_x составляет $\Theta(n)$, причем без использования дерева суффиксов. Этот прямой метод вычисления массива суффиксов требует значительно меньшего объема памяти, чем необходимо для построения дерева суффиксов, но за счет небольшого увеличения времени вычислений [172, 151, 31, 200].

Теперь, зная, что такое массив суффиксов и как его вычислить, мы основное внимание в оставшейся части данного раздела уделим эффективному использованию этих массивов для решения основной задачи поиска: найти в строке $x = x[1..n]$ все $k \geq 0$ вхождений подстроки $u = u[1..m]$. Мы представим три алгоритма решения этой задачи, назовем их условно **наивным**, **простым** и **сложным** алгоритмами соответственно.

Наивный и простой алгоритмы используют только массив суффиксов σ , n элементов которого требуют $\log n$ бит для хранения каждого значения-элемента. Они выполняются за время, не превышающее теоретической верхней границы. На практике простой алгоритм выполняется быстрее наивного, а также, согласно работе [172], быстрее сложного алгоритма. Сложный алгоритм, кроме массива суффиксов, использует также дополнительный целочисленный массив объемом, не превышающим $3n$ (этот дополнительный массив можно получить как побочный продукт в процессе вычисления массива σ). Этот алгоритм *теоретически* в самом худшем случае выполняется за время, близкое к времени решения данной задачи с помощью деревьев суффиксов и ОАГС, но *на практике* он выполняется значительно быстрее. В табл. 5.2 приведены основные временные и пространственные характеристики этих алгоритмов в предположении, что массив суффиксов уже вычислен (на этапе предварительной обработки), а дерево суффиксов, если оно использовалось для вычисления массива σ , удалено из памяти.

Таблица 5.2. Характеристики алгоритмов (на основе массивов суффиксов) нахождения всех вхождений подстроки u в строку x

Алгоритм	Требуемый объем памяти (байты)	Теоретическая верхняя граница времени выполнения
Наивный	$n \log(n/8)$	$O(m(\log n + k))$
Простой	$n \log(n/8)$	$O(m(\log n + k))$
Сложный	$4n \log(n/8)$	$O(m + \log n + k)$

Основная идея всех трех алгоритмов элементарна и известна любому студенту, изучавшему курс компьютерных наук, — это бинарный поиск! Наивный алгоритм не использует ничего, кроме бинарного поиска, который использует σ как индексный массив для выбора M -го лексикографически наибольшего суффикса. Как показано в листинге алгоритма 5.3.2, он (алгоритм) возвращает ненулевое значение, если подстрока $u = u[1..m]$ входит в строку $x = x[1..n]$, при этом возвращаемое значение равно целому числу j , такому, что $u = x[j..j + m - 1]$. Алгоритм выполняет $O(\log n)$ сравнений, каждое из которых требует $O(m)$ времени. Поэтому все сравнения выполняются за время $O(m \log n)$. Как упоминалось выше, все k вхождений подстроки u в строку x индексированы соседними элементами массива σ . Таким образом, чтобы найти остальные $k - 1$ вхождений, надо просмотреть в массиве σ соседние позиции из интервала $[L, R]$, который содержит M -ю позицию первого найденного вхождения. Это требует дополнительно $O(mk)$ времени, что доказывает оценку времени выполнения, приведенную в табл. 5.2.

Алгоритм 5.3.2. (Наивный Алгоритм)

▷ *Наивное использование массива суффиксов для локализации u в x*
 $j \leftarrow 0; L \leftarrow 1; R \leftarrow n$
repeat
 $M \leftarrow \lceil (R + L)/2 \rceil$
 if $u = x[\sigma[M].. \sigma[M] + m - 1]$ **then**
 $j \leftarrow \sigma[M]$
 else if $u > x[\sigma[M].. \sigma[M] + m - 1]$ **then**
 $L \leftarrow M$
 else
 $R \leftarrow M$
until $L = M$ или $j \neq 0$

Теперь представим простой и сложный алгоритмы. Простой алгоритм можно рассматривать как усовершенствование наивного алгоритма, где очевидна необходимость повысить эффективность обоих операторов **if**. Один способ такого усовершенствования показан в упражнении 5.3.8. Основой более значительных изменений может служить ранее описанное свойство массивов суффиксов: если суффиксы $\sigma[L]$ и $\sigma[R]$ имеют одинаковый префикс, тогда такой же префикс имеют суффиксы $\sigma[K]$, где $K = L, L + 1, \dots, R$. Назовем это свойство массивов суффиксов *кластерным*. Используя это свойство мы с неизбежностью возвращаемся к еще одному понятию, а именно к понятию *наибольшего общего префикса* (LCP), введенному в разделе 5.2.1. Напомним, что $\text{lcp} = |\text{LCP}|$ обозначает длину такого префикса. Для модификации алгоритма 5.3.2 введем величины $P_L = \text{lcp}(u, \sigma[L])$, $P_R = \text{lcp}(u, \sigma[R])$ и $P = \min\{P_L, P_R\}$. В соответствии с кластерным свойством сравнение u с префиксами $\sigma[K]$ необходимо выполнять только для $u[P + 1..m]$, поскольку префикс $u[1..P]$ одинаковый у всех этих суффиксов. Приведем листинг простого алгоритма.

Алгоритм 5.3.3. (Простой Алгоритм)

▷ *Простое использование массива суффиксов для локализации u в x*
 $j \leftarrow 0; L \leftarrow 1; R \leftarrow n$
 $P_L \leftarrow 0; P_R \leftarrow 0$
repeat
 $P \leftarrow \min\{P_L, P_R\}$.
 $M \leftarrow \lceil (R + L)/2 \rceil$
 ▷ *Вычисление $\text{lcp}(u, \sigma[M])$*
 $P_M \leftarrow P + \text{lcp}(u[P + 1..m], x[\sigma[M] + P..n])$
 if $P_M = m$ **then**
 $j \leftarrow \sigma[M]$
 else if $u[P_M + 1] > x[\sigma[M] + P_M]$ **then**

```

     $L \leftarrow M; P_L \leftarrow P_M$ 
else
     $R \leftarrow M; P_R \leftarrow P_M$ 
until  $L = M$  или  $j \neq 0$ 
    
```

Для выполнения каждого оператора этого алгоритма, за исключением вычисления lcr , требуется константное время, поэтому справедлива верхняя оценка временной сложности алгоритма, приведенная в табл. 5.2. Нетрудно привести примеры строк u и x , для которых алгоритм выполняется за время, равное верхней оценке. Поэтому простой алгоритм не является улучшением наивного алгоритма в теоретическом плане. Однако на практике, поскольку P монотонно не убывает, операция сравнения букв в среднем требует значительно меньше времени, чем в наивном алгоритме. По-видимому, простой алгоритм является наиболее быстрым среди всех известных последовательных алгоритмов нахождения вхождений заданной подстроки u в заданную строку x .

Теперь рассмотрим сложный алгоритм. Ключевая идея этого алгоритма — одновременно определить все возможные интервалы поиска $[L, R]$ ($L < R$), которые могут возникнуть в процессе бинарного поиска, а также предварительно вычислить $\text{lcr}(\sigma[L], \sigma[R])$ для этих интервалов. Как мы увидим далее, общее количество таких интервалов равно $2n - 3$. Все предварительно вычисленные значения lcr поместим в новый массив (или строку) $\pi = \pi_x$, длина которого не превысит $3n - 8$ при $n \geq 5$. Мы также покажем, что массив π можно получить в процессе первоначального вычисления массива суффиксов σ на основе дерева суффиксов T_x . Для этого дополнительно потребуются время и объем памяти, линейные относительно $|x|$. После создания массива π вычисление величины $\text{lcr}(\sigma[L], \sigma[R])$ можно будет выполнить за константное время просто путем выбора соответствующей позиции в массиве π_x . В этом случае общее время бинарного поиска можно уменьшить до величины порядка $O(m + \log n)$.

Из всего, что сказано выше, наиболее просто доказать, что количество интервалов поиска $[L, R]$ действительно равно $2n - 3$. Для примера на рис. 5.18 показано бинарное дерево для $n = 7$, где каждому из 11 узлов соответствует свой интервал поиска. Общий случай построения бинарного дерева поиска и соответствующих интервалов предложено рассмотреть в упражнении 5.3.12. На рис. 5.18 номера, расположенные внутри узлов дерева, представляют позиции в массиве π . Эти номера определяются согласно “правилу кучи”, поэтому интервал, значение lcr которого хранится в $\pi[i]$, разбивается на два подынтервала, значения lcr которых хранятся в $\pi[2i]$ и $\pi[2i + 1]$. Читателю не составит труда проверить, что при такой схеме хранения максимальная длина массива π при $n \geq 5$ не превышает $3n - 8$.

Далее, предполагая, что массив π уже вычислен, опишем выполнение сложного алгоритма, основное внимание уделив шагу обработки интервала $[L, R]$. Пред-

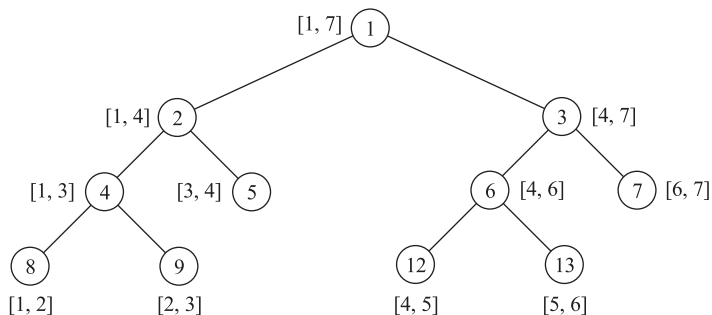


Рис. 5.18. Бинарное дерево интервалов $[L, R]$

положим, как и в простом алгоритме, что доступны значения P_L и P_R и каким-либо способом вычислено M . Также предположим, что вычислено значение i , такое, что

$$P_{LR} \equiv \text{lcp}(\sigma[L], \sigma[R]) = \pi[i].$$

Для выполнения первого шага алгоритма должны быть известны все значения (L, R, P_L, P_R, i) . Далее мы покажем, как эти величины вычисляются на последующих шагах алгоритма за время, пропорциональное величине изменения $\max\{P_L, P_R\}$, плюс константа.

Сначала рассмотрим случай $P_L = P_R$. Можно вычислить P_M таким же способом, как и в простом алгоритме, используя равенство $P = P_L$. Вычисление значения lcp требует $P_M - P_L + 1$ сравнений букв и выполняется, как показано ранее, за константное время.

Теперь пусть $P_L \neq P_R$. Для определенности положим $P_L > P_R$. (Случай $P_L < P_R$ является симметричным рассматриваемому сейчас.) Возможны три ситуации, в зависимости от значения $P_{LR} = \pi[2i]$.

1. $P_{LM} > P_L$. Для каждого $K \in L..M$ $\sigma[K]$ имеет префикс длиной $P_{LM} > P_L = \text{lcp}(u, \sigma[L])$. Поэтому, согласно кластерному свойству, u не может входить в интервал $[L, R]$, но, возможно, u входит в интервал $[M, R]$. Однако проверять последнюю возможность нет необходимости: поскольку уже $\sigma[M]$ имеет префикс длиной P_L , подстановки $L \leftarrow M, i \leftarrow 2i + 1$ позволяют определить все параметры, необходимые для следующего шага, без изменения величины $\max\{P_L, P_R\}$.
2. $P_{LM} < P_L$. В этом случае выполняется неравенство $\text{lcp}(u, \sigma[L]) > \text{lcp}(u, \sigma[M])$. Предположим, что u входит в строку x в некоторой позиции K , принадлежащей интервалу $[M, R]$. Тогда $\sigma[K]$ имеет префикс длиной P_L (как и $\sigma[L]$), но в промежуточных позициях $M \in L..K$ $\sigma[M]$, согласно кластерному свойству, не имеет префиксов длиной P_L . Поэтому u не может

входить в строку x в позициях из интервала $[M, R]$, а если все же входит в строку x , то в позициях из интервала $[L, M]$. Тогда подстановки $R \leftarrow M$, $P_R \leftarrow P_{LM}$, $i \leftarrow 2i$ позволяют определить все параметры, необходимые для следующего шага, и опять без изменения величины $\max\{P_L, P_R\}$.

3. $P_{LM} = P_L$. В этом случае $\text{lcp}(u, \sigma[L]) = \text{lcp}(u, \sigma[M])$. Снова, используя равенство $P = P_L$, вычисляем P_M тем же способом, что и в простом алгоритме. Затем вычисляем (L, P_L, i) или (R, P_R, i) за константное время.

Теперь покажем, что на каждом шаге сложного алгоритма пять параметров (L, R, P_L, P_R, i) вычисляются за время, пропорциональное величине изменения $\max\{P_L, P_R\}$, плюс константа. Поскольку величина $\max\{P_L, P_R\}$ монотонно не возрастает, а сам интервал $[L, R]$ монотонно сходится к конечному интервалу $[L, L + 1]$, то алгоритм определяет вхождение u (если оно существует) либо в некоторой позиции M , либо в конечном интервале $[L, L + 1]$. Тогда, используя немного модифицированное условие завершения поиска для конечного интервала, алгоритм найдет вхождение подстроки u в строку x за время $O(m + \log n)$, как указано в табл. 5.2. Мы оставим для упражнения 5.3.16 представление сложного алгоритма в алгоритмической форме. В упражнении 5.3.17 предложено показать, как с помощью массива π можно найти все k вхождений подстроки u в строку x за дополнительное время $\Theta(k)$.

Нам осталось показать, как вычислить массив π в процессе вычисления массива σ . Для вычисления массива π применяется следующий двухшаговый процесс.

1. При лексикографическом обходе в глубину дерева T_x путь от суффикса $\sigma[i]$ до суффикса $\sigma[i + 1]$ проходит через узел v , который является корнем поддеревя минимального размера, содержащего как $\sigma[i]$, так и $\sigma[i + 1]$. Как отмечалось в разделе 5.2.1, $\text{lcp}(\sigma[i], \sigma[i + 1]) = |v|$. Тогда для всех $i = 1, 2, \dots, n - 1$ значения lcp для интервала $[i, i + 1]$ можно вычислить при обходе в глубину дерева T_x . В упражнении 5.3.14 показано, как вычислить позицию j в массиве π , соответствующую интервалу $[i, i + 1]$. Значения lcp для этого интервала можно сохранить путем присвоения $\pi[j] \leftarrow |v|$.
2. Как показано в упражнении 5.3.10, для любого интервала $[L, R]$, $R \geq L$, выполняется равенство

$$\text{lcp}(\sigma[L], \sigma[R]) = \min\{\text{lcp}(\sigma[L], \sigma[M]), \text{lcp}(\sigma[M], \sigma[R])\} \quad (5.10)$$

для всех $M \in L..R$. Создавая дерево интервалов от листьев до корня (т.е. в порядке заполнения массива π), это равенство позволяет вычислить все элементы массива π за время $O(n)$.

Этим заканчивается наше исследование массивов суффиксов, возможно, наиболее привлекательной суффиксной структуры. Они вычисляются с относительно небольшими временными и пространственными затратами и не зависят напрямую

от размера алфавита. Но, конечно, опосредованно массивы суффиксов зависят от размера алфавита, поскольку вычисляются на основе деревьев суффиксов. Однако, как упоминалось ранее, можно исключить деревья суффиксов из процесса вычисления массивов за счет некоторого увеличения времени вычислений [172]. Фактически, большинство современных алгоритмов непосредственного вычисления массивов суффиксов [151, 31, 200] появилось в результате конкуренции по критерию времени выполнения с алгоритмом Мак-Крейта построения дерева суффиксов [200].

В разделе 13.2 мы используем массивы суффиксов для построения эффективных решений некоторых задач вычисления повторяющихся подстрок строки x .

Упражнения 5.3

1. Докажите, что отношение Π -эквивалентности, определенное на подстроках строки x , является отношением эквивалентности.
2. Если v — подстрока строки x , то обязано ли поддереву T'_v с корнем в узле v быть деревом суффиксов?
3. Пусть подстроки u_1 и u_2 Π -эквивалентны и выполняется неравенство $|u_1| < |u_2|$. Покажите, что j будет меткой листа в поддереве T'_{u_1} только в том случае, когда в поддереве T'_{u_2} есть лист с меткой $j + |u_2| - |u_1|$.
4. Пусть T — дерево, узлы которого соответствуют непустым множествам положительных целых чисел. Предположим, что корню дерева соответствует множество $\{1, 2, \dots, n\}$, и пусть сынам каждого узла соответствуют непересекающиеся множества, объединение которых входит в множество родителя. Покажите, что такое дерево может иметь не более $2n - 1$ узлов.
5. Основываясь на исходном определении ОАГС, докажите, что для любого узла только одна входящая в него дуга имеет ширину, превышающую 1.
6. На этапе предварительной обработки алгоритма 5.3.1 для выполнения вычислений совершался обход дерева T_x в определенном порядке. Объясните, как организовать обход дерева так, чтобы исключить зависимость алгоритма от размера алфавита, и как при этом выполнить вычисления за время, пропорциональное только $|x|$.
7. В тексте сделано утверждение, что поиск в глубину по дереву T_x для вычисления σ_x можно “в определенных случаях” выполнить за время, пропорциональное только n . Приведите эти “определенные случаи” и опишите соответствующий процесс поиска.
8. Наивный алгоритм выполняет повторные проверки позиций подстроки u в случае, если она имеет общий префикс с листом $\sigma[M]$. Измените этот алгоритм таким образом, чтобы исключить повторные проверки.

9. Объясните пессимистические замечания автора об алгоритме 5.3.2 и покажите изменения, которые необходимо сделать в алгоритме для усовершенствования бинарного поиска.
10. Докажите, что для произвольных позиций L и R в строке x и произвольной строке u справедливо равенство

$$\min\{|u|, \text{lcp}(\sigma[L], \sigma[R])\} = \min\{P_L, P_R\}.$$

Приведите примеры классов строк x и u , для которых простой алгоритм будет выполняться за время $\Omega(m \log n)$.

11. Пусть алгоритм 5.3.2 выполняется над массивом длиной n . Покажите, что
 - а) для $n \geq 2$ количество возможных интервалов $[L, R]$ ($L < R$) равно в точности $2n - 3$;
 - б) для $n \geq 5$ максимальный размер массива π , если он вычисляется так, как описано в тексте, не превышает $3n - 8$.

Охарактеризуйте те значения n , для которых достигается верхняя оценка утверждения б).

12. Вычислите массив π для строки $g = abaabaab$.
13. Запишем число n как $n = 2^k + h$ для $k \geq 1$ и $h \in 0..2^k - 1$. Получите формулы для вычисления размера массива π как функции от k и h . Затем покажите, как найти в массиве π позицию, соответствующую интервалу $[L, L + 1]$, $L \in 1..n - 1$.
14. Вместо массива π с $3n - 8$ элементами можно использовать бинарное дерево с $2n - 3$ узлами. Обсудите эту возможность и объясните, почему первый вариант (массив π) предпочтительнее.
15. Основываясь на описании в тексте сложного алгоритма, запишите его в алгоритмической форме.
16. Пусть с помощью сложного алгоритма найдено первое вхождение подстроки u в строку x . Покажите, что все другие вхождения подстроки u можно локализовать за время, пропорциональное количеству этих вхождений.
Совет. Используйте кластерное свойство и тот факт, что массив π содержит позиции, соответствующие всем интервалам $[L, L + 1]$, $L = 1, 2, \dots, n - 1$.
17. Докажите равенство (5.10).
18. Запишите второй шаг алгоритма вычисления массива π и затем покажите, что он выполняется за время $\Theta(n)$.
19. В упражнении 5.2.18 предлагалось оценить необходимый объем памяти для деревьев суффиксов для строк длиной $n = 100$ Кбайт, 1 Мбайт и 5 Мбайт

в случае $\alpha = 4$. На основе табл. 5.2 получите аналогичные оценки для массивов суффиксов и выведите коэффициент отношения этих оценок к объему памяти, необходимому для хранения только строки x . Какие заключения можно сделать на основании этого коэффициента? Какой эффект произведет на дерево суффиксов и массив суффиксов увеличение или удвоение размера алфавита?

ГЛАВА 6

Декомпозиция строковых последовательностей

“Всегда” и “никогда” — вот два слова, которые вы всегда должны помнить, чтобы никогда их не произносить.

— У. А. Л. Джонсон (1906–1965)

В разделе 1.4 мы изучили специальные строковые последовательности, которые назвали линдонскими словами, и некоторые их свойства. Там же показали, что линдонские слова могут стать основой внутреннего паттерна, который называется линдонской декомпозицией, причем такой паттерн единственный для любой строковой последовательности. Однако в главе 1 мы не показали, как для данной строки x эффективно вычислить линдонскую декомпозицию $w_1^{q_1} w_2^{q_2} \dots w_k^{q_k}$. Поэтому в разделе 6.1 опишем два разработанных Дювалем (Duval, [80]) элегантных алгоритма, которые выполняются за линейное время. В следующем разделе покажем два применения этих алгоритмов. Одно из этих применений уже описано в разделе 1.4 — линдонская декомпозиция строки x непосредственно ведет к канонической форме соответственной строковой петли $C(x)$; таким способом получаем эффективное решение задачи, которая часто возникает в вычислительной геометрии.

В разделе 6.3 описан другой тип декомпозиции, полностью отличной от линдонской, которая называется s -факторизацией [156] и которая также выполняется за линейное время. Этот тип декомпозиции первоначально разрабатывался как

алгоритм сжатия данных [235], однако s -факторизация имеет также важные приложения для вычисления кратных подстрок строки x (см. раздел 12.2).

В этой главе, как и в предыдущей, предполагается, что алфавит A упорядочен. Для линдонской декомпозиции мы покажем, что все алгоритмы, ее вычисляющие, не зависят от размера алфавита — вычислить линдонскую декомпозицию на неограниченных алфавитах не труднее, чем на бинарном алфавите. Однако таким замечательным свойством не обладает s -факторизация, для эффективного вычисления которой используется дерево суффиксов.

6.1 Линдонская декомпозиция: алгоритм Дюваля

Сначала рассмотрим первый алгоритм Дюваля, который за линейное время обрабатывает строку $x = x[1..n]$ по одной букве слева направо, используя при этом некоторое количество возвратов к предыдущим позициям строки x . Затем опишем второй алгоритм Дюваля, который является модификацией первого алгоритма и уменьшает количество обратных просмотров за счет использования дополнительной памяти объемом $\Theta(n)$. Эти алгоритмы основаны на интересной взаимосвязи двух внешне совершенно различных понятий: нормальной формы и линдонских слов.

Перед описанием алгоритмов Дюваля введем следующие обозначения.

- Пусть L обозначает множество всех линдонских слов.
- В честь Чена, Фокса и Линдона (Chen, Fox, Lyndon, [48]) линдонскую декомпозицию строки x будем обозначать как $\text{CFL}(x)$.
- Через $\text{NF}(x)$ обозначим нормальную форму $u^r u'$ строки x , где u — образующая подстрока строки x минимальной длины, u' — собственный префикс подстроки u и $r \geq 1$. В этом случае будем говорить, что подстрока u **порождает** строку x . Будем обозначать как L^* множество всех строк, порожденных линдонскими словами.

Напомним, что лемма 1.4.5 утверждает, что любое линдонское слово примитивно. Поскольку линдонское слово порождает само себя при $r = 1$ и $u' = \varepsilon$, то, очевидно, $L \subset L^*$. Но, конечно, не каждая строка с примитивной порождающей подстрокой принадлежит множеству L^* . Например, $\text{NF}(baaba) = (baa)(ba)$ с примитивной порождающей подстрокой baa , но поскольку $baa \notin L$, то и $baaba \notin L^*$. С другой стороны, из теоремы 1.4.2 следует, что любая строка, не являющаяся кратной, имеет единственный сдвиг, который будет линдонским словом. Так $aab = R_1(baa) \in L$, и поэтому, например, $(aab)^r aa \in L^*$ для любого $r \geq 1$.

На каждом шаге первого алгоритма Дюваля исходная строка x рассматривается как конкатенация (сцепление) трех факторов (подстрок) $x_1 x_2 x_3$, где

- подстрока x_1 является префиксом строки x , для которой $\text{CFL}(x_1)$ уже вычислена;
- подстрока x_2 уже просмотрена алгоритмом, но CFL для нее еще не вычислена;
- подстрока x_3 является суффиксом строки x , которая еще не обрабатывалась алгоритмом.

На каждом шаге алгоритма подстрока x_1x_2 увеличивается справа на одну букву (пусть это будет буква μ), тогда как подстрока x_3 уменьшается слева на эту же букву. В алгоритме на любом шаге должно выполняться условие $x_2 \in L^*$. Тогда на следующем шаге это условие выполнится, если $x_2\mu \in L^*$. Если же $x_2\mu \notin L^*$, тогда $x_2\mu$ разбивается на два фактора $x'_1x'_2$, таких, что можно вычислить $\text{CFL}(x'_1)$, а $x'_2 \in L^*$. В последнем случае новым префиксом строки x , для которого необходимо вычислить CFL , становится подстрока $x_1x'_1$.

Теперь рассмотрим детали первого алгоритма Дювала совместно с их обоснованием. Основой этих обоснований послужат следующие три леммы, которые, в свою очередь, основываются на леммах из раздела 1.4. Первая из этих лемм показывает, как можно классифицировать подстроки $x_2\mu$.

Лемма 6.1.1. Пусть $x_2 \in L^*$ и $\text{NF}(x_2) = w^r w'$ для некоторого $w \in L$ и пусть $w = w' \lambda w''$, где $\lambda \in A$. Тогда для $\mu \in A$

- а) $\mu > \lambda \Rightarrow x_2\mu \in L$;
- б) $\mu = \lambda \Rightarrow x_2\mu \in L^* - L$;
- в) $\mu < \lambda \Rightarrow x_2\mu \notin L^*$.

Доказательство. Чтобы доказать утверждение а), сначала покажем, что $w'\mu \in L$. Если $w' = \varepsilon$, то это выполняется по определению. Поэтому пусть w' является непустой строкой. Поскольку $w \in L$, из леммы 1.4.6 следует, что

$$w = w' \lambda w'' < \text{suff}(w') \lambda w''$$

для любого непустого собственного суффикса $\text{suff}(w') \lambda$ подстроки $w' \lambda$. Так как $\lambda w'' < \mu$, то

$$w' \lambda w'' < \text{suff}(w') \lambda w'' < \text{suff}(w') \mu,$$

и поэтому $w' < \text{suff}(w') \mu$. Следовательно, для любого непустого собственного суффикса $\text{suff}(w') \mu$ подстроки $w' \mu$

$$w' \mu < \text{suff}(w') \mu.$$

Тогда из леммы 1.4.6 вытекает, что $w' \mu \in L$.

По индукции, основанной на лемме 1.4.7, доказываем утверждение а): поскольку $w \in L$ и $w' \mu \in L$ и при этом $w < w' \mu$, то отсюда следует, что $ww' \mu \in L$.

В общем случае, если $w \in L$ и $w^s w' \mu \in L$ для некоторого неотрицательного целого s и при этом выполняется неравенство $w < w^s w' \mu$, то отсюда вытекает, что $w^{s+1} w' \mu \in L$.

Для доказательства утверждения б) заметим, что подстрока $x_2 \mu = x_2 \lambda$ имеет непустую грань $w' \lambda$. Поэтому, вследствие леммы 1.4.5, $x_2 \mu \notin L$. Однако подстрока $x_2 \mu = w^r w' \lambda$ порождена подстрокой w , и поэтому $x_2 \mu \in L^*$.

Для доказательства утверждения в) имеем

$$\lambda > \mu \Rightarrow w' \lambda > w' \mu \Rightarrow x_2 \mu > w' \mu,$$

поэтому строка $x_2 \mu$ лексикографически больше, чем ее собственный префикс. Тогда, вследствие леммы 1.4.6, $x_2 \mu \notin L$. Предположим, что строка $x_2 \mu$ порождена подстрокой $\bar{w} \in L$. Поскольку $x_2 \mu \notin L$, подстрока \bar{w} должна быть собственным префиксом строки $x_2 \mu$ и, следовательно, префиксом x_2 . Но любой префикс строки x_2 длиной больше длины подстроки w не может быть примитивной строкой. Поэтому на основании леммы 1.4.5 заключаем, что $|\bar{w}| \leq |w|$. Но строгого неравенства $|\bar{w}| < |w|$ быть не может, поскольку тогда x_2 должна порождаться \bar{w} , что противоречит требованию нормальной формы о том, что w является порождающей x_2 подстрокой минимальной длины. Поэтому остается только возможность, когда $\bar{w} = w$. Но эта ситуация также не имеет места, поскольку $\mu \neq \lambda$ и $w' \mu$ не является префиксом w . Отсюда заключаем, что порождающая $x_2 \mu$ подстрока не принадлежит множеству L , и поэтому $x_2 \mu \notin L^*$. ■

Чтобы пояснить доказанную лемму, рассмотрим пример, где $x_2 = (abb)^2 aab$ и $w = abb \in L$, $r = 2$, $w' = ab$, $\lambda = b$. Если следующая буква $\mu = c > \lambda$, тогда строка $x_2 c$, очевидно, будет лексикографически меньше любого собственного сдвига и поэтому $x_2 c \in L$, как предсказывает утверждение а) леммы. Для случая $\mu = b = \lambda$ строка $x_2 b = (abb)^3$ принадлежит множеству $L^* - L$, что соответствует утверждению б) леммы. Если же $\mu = a < \lambda$, тогда $NF(x_2 a) = (abbbaabbaab)a$. Поскольку порождающая подстрока $(abb)^2 aab$ не является примитивной, то $x_2 a \notin L^*$, как предсказывает утверждение в) леммы.

Следующие две леммы дают основу для построения линдонской декомпозиции.

Лемма 6.1.2. Пусть $x = wx'$, где w — наибольший префикс строки x , который является линдонским словом. Тогда

$$\text{CFL}(x) = w \text{CFL}(x').$$

Доказательство. Пусть $\text{CFL}(x') = w'_1 w'_2 \dots w'_k$. Поскольку w — наибольший префикс, то $ww' \notin L$. Поэтому, вследствие леммы 1.4.7, $w \geq w'_1$. Тогда $ww'_1 w'_2 \dots w'_k$ по определению является линдонской декомпозицией строки x , причем единственной (теорема 1.4.9). Следовательно, $\text{CFL}(x) = w \text{CFL}(x')$. ■

Следующая лемма показывает, как вычисляются компоненты декомпозиции $\text{CFL}(x)$ в случае, когда очередная буква μ “мала”, т.е. когда выполняется утверждение в) леммы 6.1.1.

Лемма 6.1.3. Пусть $x = x_2\mu x'$, где $x_2 \in L^*$ и $\mu \in A$. Пусть также $\text{NF}(x_2) = w^r w'$, где $w = w'\lambda w''$ и $\lambda \in A$. Если $\mu < \lambda$, тогда

$$\text{CFL}(x) = w^r \text{CFL}(x'\mu x').$$

Доказательство. Отметим, что поскольку $x_2 \in L^*$, то $w \in L$. Покажем, что w — наибольший префикс строки x , являющийся линдонским словом. Предположим противное, что существует префикс $\bar{w} \in L$ строки x , такой, что $|\bar{w}| > |w|$. Но если $|\bar{w}| \leq |x_2|$, тогда \bar{w} имеет непустую грань и, следовательно, не является примитивной строкой. Но, вследствие леммы 1.4.5, \bar{w} не может быть линдонским словом. Поэтому возможен только случай, когда $|\bar{w}| > |x_2|$. Тогда $\bar{w} = x_2\lambda\bar{x}$ для некоторого, возможно пустого, префикса \bar{x} строки x' . Следовательно, \bar{w} имеет суффикс $w'\mu\bar{x} < w'\lambda$, где $w'\lambda$ — префикс подстроки \bar{w} . Поэтому

$$\bar{w} \geq w'\lambda > w'\mu\bar{x} = \text{suff}(w),$$

что противоречит лемме 1.4.6. Отсюда заключаем, что не существует префикса $\bar{w} \in L$ строки x , такого, что $|\bar{w}| > |w|$.

Теперь можно применить лемму 6.1.2 к строке $x = w^r w'\mu x'$, откуда получаем, что

$$\text{CFL}(x) = w \text{CFL}(w^{r-1} w'\mu x'),$$

где или $r - 1 = 0$ или w является наибольшим префиксом строки $w^{r-1} w'\mu x'$ и линдонским словом. Это доказывает утверждение леммы. ■

Возвращаясь к примеру $x\mu = (aabb)^2(aab)a$, мы видим, что в случае $\mu = a$

$$\text{CFL}(x\mu) = (aabb)^2 \text{CFL}(aaba).$$

Отметим, что лемму 6.1.3 можно применить и в случае $\mu = \$$, где $\$$ — символ окончания строки, который лексикографически меньше любой буквы алфавита. Этот факт используется в алгоритме Дюваля.

Теперь можно описать работу алгоритма Дюваля: на каждом шаге этого алгоритма в процессе просмотра букв слева направо строка x представляется в виде $x = x_1 x_2 x_3$, где подстрока x_1 уже представима в форме линдонской декомпозиции, подстрока $x_2 \in L^*$ представима в нормальной форме, а начальный символ μ подстроки x_3 используется для управления нормальной формой $\text{NF}(x_2)$ и для определения того, когда префикс w^r подстроки x_2 присоединится к $\text{CFL}(x_1)$. Подробно алгоритм показан в приведенном ниже листинге, в котором переменные h, i, j указывают позиции правых окончаний x_1, λ и μ соответственно. Отметим, что присвоение $i \leftarrow h + 1$ соответствует случаю а) леммы 6.1.1, а присвоение $i \leftarrow i + 1$ — случаю б) этой леммы.

Алгоритм 6.1.1. (Первый алгоритм Дюваля)

▷ Для входной строки x выводятся упорядоченные позиции
 ▷ концов всех линдонских слов в $\text{CFL}(x)$
 $h \leftarrow 0 - h = |x_1|$, текущая длина уже вычисленной части CFL
while $h < n$ **do** ▷ цикл выполняется до тех пор, пока $|\text{CFL}(x_1)| < n$
 $i \leftarrow h + 1 - x[i] = \lambda$
 $j \leftarrow h + 2 - x[j] = \mu$
 ▷ следующий цикл выполняется до тех пор, пока $x_2 \in L^*$
 while $x[j] \geq x[i]$ **do** ▷ выполняется проверка $\mu \geq \lambda$
 if $x[j] > x[i]$ **then**
 $i \leftarrow h + 1$
 else
 $i \leftarrow i + 1$
 $j \leftarrow j + 1$
 ▷ для $\mu < \lambda$ (на основе леммы 6.1.3) этот цикл выводит все r
 ▷ вхождений текущей порождающей строки w
 ▷ длиной $j - i + 1$
repeat
 $h \leftarrow h + (j - i)$; **output** h
until $h \geq i$

Отметим, что в случае $\mu < \lambda$ алгоритм выполняет возврат назад, когда после вывода r копий линдонского слова w алгоритм возвращается к началу строки $w'\mu$ с помощью операторов $i \leftarrow h + 1$ и $j \leftarrow h + 2$. Поскольку w' является префиксом w , то для w' повторяются уже ранее выполненные действия всякий раз, когда $|w'| > 1$. В примере $x\mu = (aabb)^2(aab)a$, приведенном выше, строка $w' = aab$ будет обрабатываться еще раз, хотя она уже обрабатывалась как префикс строки $w = aabb$ в начале выполнения внешнего цикла **while**. Как мы покажем далее, второй алгоритм Дюваля исключает повторную обработку путем введения записей о периодичности строки w .

Сформулируем основной результат данного раздела.

Теорема 6.1.4. Алгоритм 6.1.1 корректно вычисляет линдонскую декомпозицию за время порядка $\Theta(n)$ с использованием памяти фиксированного объема.

Доказательство. Корректность алгоритма вытекает из лемм 6.1.1 и 6.1.3, как показано в приведенных выше обсуждениях. Ограниченность используемой памяти следует из того факта, что для выполнения алгоритма не требуется никакой дополнительной памяти, кроме памяти для хранения переменных h , i и j .

Для анализа времени выполнения алгоритма заметим, что каждый вывод состоит из $r \geq 1$ копий строки $w \in L$ и что повторные вычисления выполняются

только с собственным префиксом w' строки w . Поскольку w является порождающей строкой для текущей подстроки $x_2 = w^r w'$, то ее длина $p = |w|$ является периодом подстроки x_2 . Поэтому требуется $rp + |w'| - 1$ итераций внутреннего цикла **while**, чтобы присоединить r строк общей длиной rp к текущей (уже вычисленной) декомпозиции. В лучшем случае ($|w'| = 0$) необходимо $rp - 1$ итераций внутреннего цикла **while** для вывода каждого множества, состоящего из rp букв. В самом худшем случае ($r = 1$ и $|w'| = p - 1$) требуется $2p - 2$ таких итераций для вывода каждого множества из p букв. Поскольку вычисления на каждой итерации выполняются за константное время, общее время вывода результата составит величину порядка $O(n)$. Отсюда следует, что алгоритм 6.1.1 выполняется за время $\Theta(n)$. ■

Прежде чем перейти к рассмотрению второго алгоритма Дюваля, обсудим взаимосвязь между нормальной формой и линдонскими словами, которая проявилась в первом алгоритме Дюваля. Этот алгоритм можно рассматривать как алгоритм, который вычисляет нормальную форму максимальных элементов множества L^* , которые встречаются в данной строке x . При этом кажется, что нормальная форма вычисляется без использования массива граней. Однако это не совсем так, как мы сейчас покажем.

Для переменных h , i и j из первого алгоритма Дюваля введем обозначения $i' = i - h$ и $j' = j - h$. Поскольку на любом шаге алгоритма $h = |x_1|$, тогда $\lambda = x[i] = x_2[i']$ и $\mu = x[j] = x_2[j']$. Заметим, что на каждом шаге в начале внешнего цикла **while** для переменных i' и j' справедливо такое утверждение: i' — наибольшее целое, меньшее j' , и при этом выполняется равенство

$$x_2[1..i' - 1] = x_2[j' - i' + 1..j' - 1]. \quad (6.1)$$

Другими словами, $x_2[1..i' - 1]$ является наибольшей гранью подстроки $x_2[1..j' - 1]$ или, в обозначениях раздела 1.3, $\beta_2[1..j' - 1] = i' - 1$, где β_2 — массив граней строки x_2 . Значения переменных i и j изменяются только во внутреннем цикле **while**, при этом переменная h там не изменяется. Поэтому равенство (6.1) останется неизменным во внутреннем цикле **while** в следующих ситуациях.

- Если $x_2[j'] = x_2[i']$. Тогда значения i' и j' возрастут на единицу, поэтому равенство (6.1) сохранит силу.
- Если $x_2[j'] > x_2[i']$. Тогда $i' \leftarrow 1$, а значение j' увеличится на единицу. Из утверждения *a*) леммы 6.1.1 вытекает, что в этом случае новая строка x_2 является линдонским словом и поэтому примитивна. Следовательно, ее наибольшей гранью будет ε . Поэтому равенство (6.1) также будет иметь место.

Таким образом, во время выполнения первого алгоритма Дюваля равенство (6.1) будет выполняться всегда, что доказывает следующую лемму.

Лемма 6.1.5. Во внутреннем цикле **while** алгоритма 6.1.1 выполняется $\beta_2[j' - 1] = i' - 1$, где $i' = i - h$, $j' = j - h$ и β_2 — массив граней строки $x_2[1..j']$. ■

Основываясь на этом результате, мы можем иметь массив β_2 для каждого нового обрабатываемого значения j' , содержащего записи соответствующих значений i' , для которых уже проведено сравнение букв $x_2[j']$ и $x_2[i']$. На этом основан второй алгоритм Дюваля. Здесь для уменьшения вычислений сохраняется $n/2$ позиций массива β_2 в несколько измененном виде, как показано ниже.

Алгоритм 6.1.2. (Второй Алгоритм Дюваля)

▷ Для входной строки x выводятся упорядоченные позиции
 ▷ концов всех линдонских слов в $\text{CFL}(x)$
 $h \leftarrow 0$ ▷ $h = |x_1|$, текущая длина уже вычисленной части CFL
 $j \leftarrow 2$; $\beta'[2] \leftarrow 1$ ▷ инициализация j и $\beta'[j]$
while $h < n$ **do** ▷ цикл выполняется до тех пор, пока $|\text{CFL}(x_1)| < n$
 $i \leftarrow h + \beta'[j - h]$ ▷ $x[i] = \lambda$
 ▷ следующий цикл выполняется до тех пор, пока $x_2 \in L^*$
 while $x[j] \geq x[i]$ **do** ▷ выполняется проверка $\mu \geq \lambda$
 if $x[j] > x[i]$ **then**
 $i \leftarrow h + 1$
 else
 $i \leftarrow i + 1$
 $j \leftarrow j + 1$
 ▷ поскольку i и j изменены, то
 ▷ $\beta'[j - h] = i - h \Leftrightarrow +\beta_2[j - h - 1] = i - h - 1$
 if $j - h \leq n/2$ **then**
 $\beta'[j - h] \leftarrow i - h$
 ▷ для $\mu < \lambda$ (на основании леммы 6.1.3) этот цикл выводит все r
 ▷ вхождений текущей порождающей строки w
 ▷ длиной $j - i + 1$
 repeat
 $h \leftarrow h + (j - i)$; **output** h
 until $h \geq i$
 ▷ если $w' = \varepsilon$, то μ перемещается на одну позицию вправо
 if $h = j - 1$ **then**
 $j \leftarrow j + 1$

Сравнивая первый и второй алгоритмы Дюваля, мы видим, что в процессе вычислений переменные h , i и j принимают одинаковые значения, да и структура циклов совпадает. Различия между этими алгоритмами связаны только с повтор-

ным вычислением i после того, как выведено множество r линдонских слов w . Во втором алгоритме для этих вычислений используется массив β' , изменения в который вносятся во внутреннем цикле **while**. Изменение массива β' эквивалентно созданию массива граней β_2 для строки x_2 . Так, записи в массиве β' для любого $j' - 1$ соответствуют граням $i' - 1 = \beta_2[j' - 1] = \beta'[j']$. Поэтому значения массива β' можно использовать при обработке суффикса w' строки x_2 после вывода w^r . Осталось прояснить еще две небольшие неясности во втором алгоритме Дюваля.

1. Изменения в массив β' вносятся только в случае, когда $j - h \leq n/2$. Напомним, что $j - h = |x_2|$, где $x_2 = w^r w'$. Элементы массива β' используются только для определения позиций в строке w' после вывода w^r . Но

$$|w'| < |x_2|/2 = (j - h)/2 \leq (n + 1)/2.$$

Поэтому могут потребоваться только элементы массива β' , определяющие позиции, не превосходящие $n/2$.

2. Если $j = h + 1$, тогда вывод w^r завершается и значение j увеличивается на единицу. В этом случае $w' = \varepsilon$. Для следующего шага устанавливаются значения $j = h + 2$, $i = h + \beta'[2] = h + 1$.

Очевидно, что второй алгоритм Дюваля имеет такое же асимптотическое время выполнения, как и первый алгоритм. Поэтому справедлива следующая теорема.

Теорема 6.1.6. Алгоритм 6.1.2 корректно вычисляет линдонскую декомпозицию за время порядка $\Theta(n)$ с использованием памяти объемом $\Theta(n)$. ■

Мы закончим изучение представленных алгоритмов вопросом о более точной оценке времени их выполнения в самом худшем случае. Второй алгоритм Дюваля уменьшает объем обработки непустых суффиксов w' , но остается вопрос о том, насколько велико это уменьшение и “сколько оно стоит”. Критерием, используемым для измерения эффективности строковых алгоритмов, часто выступает количество сравнений букв в строке. Существует обширная литература [52, 96, 97, 50], посвященная определению верхних и нижних границ количества буквенных сравнений, необходимых в самом худшем случае для нахождения заданного паттерна u в данной строке x . Чтобы поддержать эту традицию, рассмотрим число буквенных сравнений в алгоритмах Дюваля.

Чтобы сделать такой анализ осмысленным, сначала надо определить, что мы понимаем под “сравнением букв”. В этой книге мы будем использовать этот термин в смысле *двоичного буквенного сравнения*, т.е. как сравнение букв, результатом которого могут быть два логических значения TRUE (Истина) или FALSE (Ложь). Например, требуется только одно сравнение букв μ и λ , чтобы определить, будет ли $\mu = \lambda$ или $\mu \geq \lambda$. Но чтобы различить три возможных случая $\mu > \lambda$, $\mu = \lambda$, $\mu < \lambda$, требуется уже два сравнения этих букв. Это определение, очевидно, согласуется с обычной двоичной логикой цифровых компьютеров.

Напомним, что в первом алгоритме Дюваля, как показано в доказательстве теоремы 6.1.4, для вывода p букв требуется не более $2p - 2$ итераций внутреннего цикла **while**. Чтобы “перевести” это число в количество буквенных сравнений, заметим, что на каждой итерации выполняется два сравнения: одно в операторе **while** и другое в операторе **if**. Поэтому во внутреннем цикле **while** всего выполняется $4p - 4$ буквенных сравнений. Кроме того, еще одно сравнение используется для проверки $\mu < \lambda$, завершающей цикл. Таким образом, в первом алгоритме Дюваля в самом худшем случае требуется не более $4n - 3$ буквенных сравнений для вывода всех линдонских слов строки x .

Во втором алгоритме Дюваля во внутреннем цикле **while** выполняется не более одной итерации, требующей не более двух буквенных сравнений для каждого значения $j = 2, 3, \dots, n$. Поэтому всего выполняется не более $2n - 2$ буквенных сравнений. Еще одно сравнение (во внутреннем операторе **while**) требуется для каждого вывода подстроки w' — это случай, когда префикс w' линдонского слова содержит одно или несколько линдонских слов, формирующих часть линдонской декомпозиции строки x . Если $w' = \varepsilon$, то значение j увеличивается на единицу. Поэтому такое дополнительное сравнение может потребоваться не более чем для половины всех значений j . Таким образом, второй алгоритм в самом худшем случае требует не более $2,5n - 2$ буквенных сравнений, что значительно меньше, чем в первом алгоритме.

Поэтому второй алгоритм лучше, чем первый. Но так ли это на самом деле?

Обсуждение 6.1.1. Конечно, и с практической, и с теоретической точки зрения особый интерес представляют наиболее точные границы времени выполнения строковых алгоритмов. Но на практике алгоритмы, требующие малые или минимальные количества буквенных сравнений, могут быть не самыми эффективными. Может случиться так, что дополнительные вычисления, позволяющие исключить или уменьшить количество буквенных сравнений, потребуют больших временных затрат, чем сами сравнения. Например, как отмечалось в разделе 4.1, простой алгоритм 2.2.1 для нахождения паттернов в среднем может выполняться быстрее, чем алгоритмы (такие, как алгоритм 7.1.1 КМП), где количество буквенных сравнений для самого худшего случая сведено к минимуму; подобно другим “примитивным созданиям”, основной добродетелью простого алгоритма является его простота.

Подобные суждения применимы и для алгоритмов Дюваля. Хотя мы привели оценки количеств буквенных сравнений для одного и другого алгоритмов, остается все-таки не ясным вопрос, когда *на практике* второй алгоритм предпочтительнее первого? Приведем некоторые суждения по этому поводу.

- Второй алгоритм требует дополнительной памяти объемом $n/2$ и соответствующей дополнительной обработки каждой новой буквы μ и каждого множества одинаковых линдонских слов. Будут ли эти затраты давать преимущества времени выполнения в среднем?

- Как показано в упражнении 6.1.3, время вывода любого линдонского слова, найденного в строке x , можно уменьшить до константного времени. Будут ли эти простые изменения в алгоритме более эффективными, чем попытка уменьшить количество буквенных сравнений?
- Верхние границы для количества буквенных сравнений в алгоритмах Дюваля, приведенные выше, не очень точны (см. упражнение 6.1.5). Даже без рассмотрения этих оценок “в среднем”, может оказаться, что в самом худшем случае различия между этими алгоритмами не определяются полученными верхними оценками. Поэтому без более точного анализа опрометчиво делать заключения об эффективности алгоритмов Дюваля только на основании имеющихся верхних оценок числа буквенных сравнений.
- Второй алгоритм Дюваля имеет преимущества перед первым алгоритмом только в случае, когда $|w'| > 1$. Для случайно выбранных строк, даже определенных на малых алфавитах, такое событие крайне редкое. Получается, что второй алгоритм придуман для решения проблемы, которой практически не существует!
- С общей точки зрения: что более важно для алгоритмов с временем выполнения $O(n)$ — поведение этих алгоритмов в среднем или их выполнение в самом худшем случае? Различие между двумя алгоритмами с временем выполнения $\Theta(n)$ в самом худшем случае может быть относительно невелико, поэтому повышается значимость их поведения в среднем.
- Необходимо также учитывать влияние размера алфавита. Даже если алфавит очень мал, например $\alpha = 2$ или 4, то для хранения строки x требуется $n/8$ или $n/4$ байт соответственно. С другой стороны, во втором алгоритме Дюваля в новом массиве β' содержится $n/2$ элементов, являющихся целыми числами из интервала от 0 до $n/2$. Поэтому всего требуется памяти объемом не менее $\lceil n' \log n' \rceil / 8$ байт, где $n' = n/2 + 1$. Если n велико, то такой большой объем памяти может быть неприемлемым.
- Необходимо также учитывать характеристики компьютера, на котором реализуется алгоритм. Например, возможность компьютера производить относительно эффективно сравнение битовых строк или сравнение байтов (или даже полных машинных слов) может значительно повлиять на то, какой из алгоритмов Дюваля будет выполняться быстрее, поскольку даже для бинарного алфавита надо производить как сравнение одиночных битов, так и обработку длинных строк, состоящих из нескольких байтов (элементы массива β').

Все эти соображения высказаны не с целью приуменьшить роль теоретических компьютерных наук, предлагающих новые и все более элегантные и эффективные алгоритмы. Они высказаны с целью показать многочисленные факторы, которые

вливают на поведение алгоритмов, когда они *реализованы* в виде компьютерных программ, при этом влияние таких факторов зачастую можно оценить только экспериментальным путем. ■

Упражнения 6.1

1. Покажите, что строка x будет линдонским словом только тогда, когда в массиве суффиксов $\sigma_x[1] = 1$.
2. Покажите, что на стандартном алфавите (4.1) любое линдонское слово будет p -канонической строкой (терминология раздела 4.1). Справедливо ли обратное утверждение?
3. Пусть в лемме 1.4.7 целочисленная переменная r равна количеству вхождений выходного линдонского слова w . Выразите r через переменные h , i и j , используемые в алгоритме 6.1.1. (Заметьте, что при $r > 1$ выводятся кратные строки.) Далее примените кодирование, которое уже многократно использовалось в этой книге, для того чтобы представить все r вхождений строки w в виде числового кортежа. Сделайте соответствующую модификацию алгоритма 6.1.1. Какова асимптотическая временная сложность модифицированного алгоритма?
4. Для алфавитов размером $\alpha = 2, 4, 8$ и строк длиной $n = 100$ Кбайт, 1 Мбайт и 5 Мбайт оцените объемы памяти, необходимые для первого и второго алгоритмов Дюваля. Далее без проведения обширных экспериментов проанализируйте соотношение “стоимость–эффективность использования” каждого алгоритма Дюваля для приведенных 9 вариантов “размер алфавита–длина строки”.
5. Приведите примеры самых худших (или хотя бы “плохих”) строк для обоих алгоритмов Дюваля. На основе этих примеров сделайте заключение о приведенных в тексте границах количеств буквенных сравнений для этих алгоритмов.
6. Во втором алгоритме Дюваля массив β' изменяется таким образом, чтобы гарантировать выполнение неравенства $j - h \leq n/2$. Поэтому данный массив содержит только $n/2$ элементов. Имеется возможность “освободить” массив β' от гарантии выполнения этого неравенства. Но в этом случае он будет содержать n элементов. Обсудите все “за” и “против” этих двух вариантов массива β' .
7. Предположим, что при сравнении букв μ и λ равновероятны три исхода $>$, $=$ и $<$. Определите нижнюю и верхнюю границы *ожидаемого* количества буквенных сравнений в каждом алгоритме Дюваля.
8. Является ли второй алгоритм Дюваля онлайн-овым?

9. В статье [80] введено множество P как множество всех префиксов линдонских слов, определенных на данном алфавите A . Если алфавит A содержит максимальную букву ω , то множество P' — это множество всех кратных строк с буквой ω . (Если в алфавите A нет буквы ω , то $P' = \emptyset$.) Докажите, что $L^* = P \cup P'$.

Совет. Используйте лемму 6.1.1, в частности утверждение *a*) этой леммы.

6.2 Применения линдонской декомпозиции

В этом разделе покажем, как первый алгоритм Дюваля можно изменить так, чтобы он за время порядка $\Theta(n)$ вычислял следующие внутренние паттерны для заданной строки x :

- лексикографически минимальный непустой суффикс $s_{\min}(u)$,
- лексикографически максимальный суффикс $s_{\max}(u)$,

определяемые для любого непустого префикса u строки x . Конечно же, подобные изменения можно внести и во второй алгоритм Дюваля.

Основой вычисления обоих паттернов послужат приведенные ниже две леммы. В первой лемме рассматривается разбиение линдонской декомпозиции строки x , подобное применяемому в алгоритме Дюваля:

$$x[1..j] = x_1 x_2 [1..j - h],$$

где $|x_1| = h$ и декомпозиция $CFL(x_1)$ считается уже вычисленной. В этой лемме утверждается, что

$$s_{\min}(x[1..j]) = s_{\min}(x_2[1..j - h]).$$

Это означает, что вычисление s_{\min} для позиций подстроки x_2 не зависит от префикса x_1 .

Лемма 6.2.1. Пусть $x = x_1 w$, где x_1 и w — непустые строки. Тогда равенство $CFL(x) = CFL(x_1)$ выполняется только в том случае, когда $s_{\min}(x) = w$.

Доказательство. Если $CFL(x) = CFL(x_1)w$, тогда w является линдонским словом и мы можем записать

$$CFL(x) = w_1 w_2 \dots w_k,$$

где $w_k = w$ и $w_1 \geq w_2 \geq \dots \geq w_k$. Без потери общности можно положить, что $s_{\min}(x) = w'_i w_{i+1} \dots w_k$ для некоторого целого $i \in 1..k$, где w'_i — непустой суффикс подстроки w_i . Вследствие леммы 1.4.6, $w'_i \geq w_i$, и поскольку $w_i \geq w_k$, то отсюда следует, что $i = k$ и $s_{\min}(x) = w_k = w$.

Теперь докажем в обратную сторону: предположим, что $s_{\min}(x) = w$, но $CFL(x) = CFL(x_1)\bar{w}$, где $\bar{w} \neq w$. Тогда, повторяя приведенные выше аргументы, получим $s_{\min}(x) = \bar{w} \neq w$, т.е. пришли к противоречию. ■

Вторая лемма структурно похожа на лемму 6.1.1 и поэтому также ведет к алгоритму, который структурно похож на первый алгоритм Дюваля.

Лемма 6.2.2. Пусть $x_2 \in L^*$ и $NF(x_2) = w^r w'$ для некоторого $w \in L$ и пусть $w = w' \lambda w''$, где $\lambda \in A$. Пусть также $s' \lambda = s_{\min}(w' \lambda)$. Тогда для $\mu \in A$

а) $\mu > \lambda \Rightarrow s_{\min}(x_2 \mu) = x_2 \mu$;

б) $\mu = \lambda \Rightarrow s_{\min}(x_2 \mu) = s' \mu$;

в) $\mu < \lambda \Rightarrow$ для любой строки $x' s_{\min}(x_2 \mu x') = s_{\min}(s' \mu x')$.

Доказательство. утверждения а): из леммы 6.1.1 следует $x_2 \mu \in L$, поэтому на основании леммы 1.4.6 можем утверждать, что $x_2 \mu$ меньше любого собственного суффикса.

Для доказательства утверждения б) достаточно заметить, что из леммы 6.1.2 следует $CFL(x_2 \lambda) = w^r CFL(w' \lambda)$, где $w \geq w' \lambda$ и $w < R_j(w)$ для всех $j \in 1..n - 1$.

Доказательство утверждения в): из леммы 1.4.7 следует, что

$$CFL(x_2 \mu x') = w^r CFL(w' \mu x').$$

Тогда (как и в утверждении б)) $s_{\min}(x_2 \mu x') = s_{\min}(w' \mu x')$.

Обозначим через S' суффикс строки w' , такой, что $|S'| > |s'|$. Так как $s' \lambda = s_{\min}(w' \lambda)$, то $s' \lambda < S' \lambda$. Поскольку $|s' \lambda| \leq |S'|$, тогда $s' \lambda \leq S'$ и $s' \lambda < S' \mu x'$. Поэтому $s' \mu x' < s' \lambda < S' \mu x'$ и, следовательно, $S' \mu x'$ не может быть минимальным суффиксом строки $w' \mu x'$. Отсюда заключаем, что любой минимальный суффикс S' должен удовлетворять неравенству $|S'| \leq |s'|$ и поэтому должен быть суффиксом строки s' . ■

Леммы 6.2.1 и 6.2.2 содержат всю информацию, необходимую для вычисления функции s_{\min} для любого префикса строки x . Как указывалось выше, лемма 6.2.1 позволяет не рассматривать строку x_1 , т.е. префикс, для которого линдонская декомпозиция уже вычислена (и для которого уже вычислена функция s_{\min}). Утверждение а) леммы 6.2.2 является просто частным случаем леммы 6.2.1: линдонское слово является собственным минимальным суффиксом. Из утверждений б) и в) леммы 6.2.2 следует, что $s_{\min}(x_2 \mu)$ полностью определяется предварительно вычисленной величиной $s_{\min}(w' \lambda)$, где $w' \lambda$ — префикс строки w и собственный префикс строки $x_2 \mu$.

Алгоритм вычисления минимального суффикса показан в приведенном ниже листинге как алгоритм 6.2.1. Чтобы упростить вывод, выводятся не сами найденные минимальные суффиксы, а целочисленный массив $\sigma' = \sigma'[1..n]$, где для каждого $j \in 1..n$ $s_{\min}(x[1..j]) = x[\sigma'[j]..j]$. Приведем пример такого массива для строки Фибоначчи f_6 (пробелами отделены линдонские слова).

	1	2	3	4	5	6	7	8	9	10	11	12	13
$f_6 =$	a	b	a	a	b	a	b	a	a	b	a	a	b
$\sigma' =$	1	1	3	4	3	6	3	8	9	8	11	12	11

Алгоритм 6.2.1 (Алгоритм вычисления минимальных суффиксов)

▷ Для входной строки $x\$$ выводится массив $\sigma'[1..n]$, такой, что

▷ подстрока $x[1..j]$ имеет непустой суффикс $x[i..j]$ только тогда,

▷ когда $\sigma'[j] = i$

$h \leftarrow 0; \sigma'[1] \leftarrow 1$

while $h < 0$ **do**

$i \leftarrow h + 1$ ▷ $x[i] = \lambda$

$j \leftarrow h + 2$ ▷ $x[j] = \mu$

while $x[j] \geq x[i]$ **do** ▷ выполняется проверка $\mu \geq \lambda$

if $x[j] > x[i]$ **then**

▷ случай $\mu > \lambda$, тогда $s_{\min}(x_2\mu) = x_2\mu$

▷ $\sigma'[j] \leftarrow h + 1; i \leftarrow h + 1$

else

▷ случай $\mu = \lambda$, тогда $s_{\min}(x_2\mu) = s_{\min}(w'\mu)$

$\sigma'[j] \leftarrow \sigma'[i] + (j - i); i \leftarrow i + 1$

$j \leftarrow j + 1$

▷ случай $\mu < \lambda$, вычисление h

$r \leftarrow \lceil (i - h)/(j - i) \rceil; h \leftarrow h + r \times (j - i)$

if $h = j - 1$ **then**

$\sigma'[j] \leftarrow h + 1$

Теорема 6.2.3. Алгоритм 6.2.1 корректно вычисляет массив $\sigma'[1..n]$ (минимальные непустые суффиксы для всех непустых префиксов) для заданной строки $x\$$ за время порядка $\Theta(n)$.

Доказательство. Заметим, что ход вычислений и использование переменных h , i и j в алгоритме 6.2.1 такие же, как и в первом алгоритме Дюваля, за исключением вычисления h для случая $\mu < \lambda$. Но в последнем случае для вычисления h используется переменная r , которая дает то же самое значение, что и вычисления в цикле **repeat** в алгоритме Дюваля (см. упражнение 6.2.1). Увеличение значения h (если $h = j - 1$) просто указывает на особый случай: когда одна оставшаяся буква $x[j]$ оказывается наименьшей буквой подстроки $x[1..j]$, тогда $s_{\min}(x[1..j]) = x[j]$ (т.е. это тот самый случай, когда $h = 0$ и $j = 1$). Таким образом, алгоритм 6.2.3 выполняется практически так же, как первый алгоритм Дюваля, и поэтому время его выполнения составляет $\Theta(n)$.

Теперь рассмотрим четыре изменения, производимые в массиве σ' . Мы уже рассмотрели случай, когда $h = j - 1$. Также тривиально изменение $\sigma'[1] \leftarrow 1$. Два

других изменения, которые вызваны ситуациями, когда $\mu > \lambda$ и $\mu = \lambda$, соответствуют утверждениям а) и б) леммы 6.2.2. Следовательно, алгоритм корректно вычисляет массив σ' . ■

Теперь рассмотрим вторую задачу, сформулированную в начале этого раздела: вычисление лексикографически максимального суффикса $s_{\max}(u)$ любого непустого префикса u строки x . Очевидный подход к решению данной задачи — заменить отношение лексикографического порядка $<$ на обратный лексикографический порядок $<_{\mathcal{R}}$ в надежде, что использование этого нового порядка в алгоритме 6.2.1 позволит получить максимальные, а не минимальные суффиксы. Итак, сначала определим новый лексикографический порядок для пар букв λ и μ алфавита A :

$$\mu <_{\mathcal{R}} \lambda \Leftrightarrow \lambda < \mu. \quad (6.2)$$

Далее следует распространить новый порядок на все строки, определенные на алфавите A . Но здесь нас подстерегают сложности: например, для строк $u = a$ и $v = ab$ выполняется неравенство $u < v$, но $v \not<_{\mathcal{R}} u$. Однако имеет место следующее утверждение.

Лемма 6.2.4. Для любых строк u и v , определенных на упорядоченном алфавите A , одновременное выполнение неравенств $u < v$ и $v <_{\mathcal{R}} u$ возможно только в том случае, когда u является префиксом v .

Доказательство этой леммы предложено дать в упражнении 6.2.6. ■

Таким образом, отношение порядка $<_{\mathcal{R}}$ распространяется на строки только в том случае, если одна из них является префиксом другой. Рассмотрим применение леммы 6.2.4 к строкам $x = cab$ и $y = bab$, у обеих этих строк максимальными суффиксами являются они сами. Для строки x нет проблем: эта строка является линдонским словом, она примитивна и лексикографически меньше, в соответствии с отношением порядка $<_{\mathcal{R}}$, чем любой из ее суффиксов, а в соответствии с отношением порядка $<$ она больше, чем любой ее суффикс. Для строки y ситуация иная: эта строка не является линдонским словом, она не примитивна и лексикографически не меньше, в соответствии с отношением порядка $<_{\mathcal{R}}$, чем любой из ее суффиксов (поскольку $bab \not<_{\mathcal{R}} b$), но в соответствии с отношением порядка $<$ она больше, чем любой ее суффикс (так как $bab > b > ab$). Поэтому применение алгоритма 6.2.1 с заменой отношения $<$ на отношение $<_{\mathcal{R}}$ приведет к желаемому результату для строки $x = cab$, но не для строки $y = bab$. В общем случае лемма 6.2.4 показывает, что такое применение алгоритма 6.2.1 будет некорректным для любых строк с непустой гранью.

Таким образом, такой наивный подход, как замена отношения $<$ на отношение $<_{\mathcal{R}}$, задачу нахождения максимального суффикса не решает. Однако в введении

обратного отношения порядка $<_{\mathcal{R}}$ есть рациональное зерно (например, для строки $x = cab$, как показано выше, получается правильный результат). Поэтому рассмотрим данное отношение подробнее. Напомним, что L^* обозначает множество всех строк, порожденных линдонскими словами. Обозначим через $L^{\mathcal{R}}$ множество всех строк x , таких, что, в соответствии с отношением порядка $<_{\mathcal{R}}$, $s_{\max}(x) = x$. Следующая цепочка лемм покажет, что $L^{\mathcal{R}} = L^*$.

Лемма 6.2.5. Пусть $x \in L^{\mathcal{R}}$ и x' — непустой префикс строки x . Тогда $x' \in L^{\mathcal{R}}$.

Доказательство. Пусть $x = x'y$ и пусть u' — собственный префикс строки x' . Поскольку $x \in L^{\mathcal{R}}$, то $u' <_{\mathcal{R}} x = x'y$. Но так как $|u'| < |x'|$, то $u' <_{\mathcal{R}} x'$ и, в соответствии с отношением порядка $<_{\mathcal{R}}$, $s_{\max}(x') = x'$. Следовательно, $x' \in L^{\mathcal{R}}$, что и требовалось доказать. ■

Лемма 6.2.6. Пусть $x \in L^{\mathcal{R}}$ и x' — префикс строки x . Тогда $xx' \in L^{\mathcal{R}}$.

Доказательство. Утверждение тривиально, если $x' = \varepsilon$, поэтому положим, что $x' \neq \varepsilon$. Пусть u — собственный суффикс строки xx' . Возможны два случая.

1. $u = x''x'$, где x'' — непустой собственный суффикс строки x .

Поскольку $x \in L^{\mathcal{R}}$, то $x'' <_{\mathcal{R}} x$. Из неравенства $|x''| < |x|$ следует, что $u = x''x' <_{\mathcal{R}} xx'$.

2. u является суффиксом x' .

Если $u = x'$, тогда u — собственный префикс строки xx' и, следовательно, $u <_{\mathcal{R}} xx'$. Если u — собственный суффикс строки x' , то, на основании леммы 6.2.5, $x' \in L^{\mathcal{R}}$. Тогда $u <_{\mathcal{R}} x' <_{\mathcal{R}} xx'$.

В обоих случаях $xx' \in L^{\mathcal{R}}$. ■

Лемма 6.2.7. Если $x \in L$, тогда $x \in L^{\mathcal{R}}$.

Доказательство. Вследствие леммы 1.4.5 строка x примитивна, поэтому любой ее непустой собственный суффикс u гарантированно не будет ее префиксом. Из леммы 1.4.6 вытекает, что $x < u$, и, вследствие леммы 6.2.4, $u <_{\mathcal{R}} x$. Но тогда, в соответствии с отношением $<_{\mathcal{R}}$, $s_{\max}(x) = x$, что и требовалось доказать. ■

Лемма 6.2.8. Если $x \in L^*$, тогда $x \in L^{\mathcal{R}}$.

Доказательство непосредственно следует из лемм 6.2.6 и 6.2.7. ■

Лемма 6.2.9. $L^{\mathcal{R}} = L^*$.

Доказательство. Из леммы 6.2.8 следует, что $L^* \subset L^{\mathcal{R}}$. Докажем включение в обратную сторону. Пусть $x \in L^{\mathcal{R}}$, u — ее наибольший префикс, который принадлежит L^* . Тогда $\text{NF}(u) = w^r w'$ для некоторой строки $w = w' \lambda w'' \in L$.

Предположим, что $u \neq x$. Тогда $x = u\mu x'$ для некоторого непустого суффикса x' . Поскольку $u\mu \notin L^*$, из леммы 6.1.1 следует, что $\mu < \lambda$. Поэтому $\lambda <_{\mathcal{R}} \mu$, тогда $w'\lambda <_{\mathcal{R}} w'\mu$ и $x = u\mu x' <_{\mathcal{R}} u\lambda x'$, что противоречит предположению о том, что $x \in L^{\mathcal{R}}$. Отсюда заключаем, что $u = x \in L^*$. ■

Чтобы иметь возможность применить полученные результаты к алгоритмам Дюваля, надо сначала разобраться, что же следует из этих результатов. Лемма 6.2.9 говорит, что строка u принадлежит множеству L^* только в том случае, когда строка u сама является своим максимальным суффиксом в соответствии с отношением $<_{\mathcal{R}}$. Поскольку $(<_{\mathcal{R}})_{\mathcal{R}} \equiv <$, обратное утверждение также имеет место: строка $u \in L^*$ в соответствии с отношением $<_{\mathcal{R}}$ только в том случае, когда u является своим максимальным суффиксом в соответствии с отношением $<$. Поэтому, выполняя алгоритм Дюваля с отношением $<_{\mathcal{R}}$, мы не потеряли надежду каким-то образом определить строки u , такие, что $u = s_{\max}(u)$.

Эти рассуждения подводят нас к критическому замечанию: строка u будет максимальным суффиксом строки x и линдонским словом только в том случае, когда строка u будет *наибольшим* суффиксом строки x , таким, что $u = s_{\max}(u)$. Но вследствие леммы 6.2.9 каждая строка u , являющаяся своим максимальным суффиксом и вычисленная на основании одного отношения порядка, принадлежит множеству L^* , вычисленному с помощью обратного отношения порядка. Таким образом, приходим к алгоритму, где сначала на основании отношения порядка $<_{\mathcal{R}}$ вычисляются элементы множества L^* (т.е. суффиксы строки x), затем, в соответствии с отношением порядка $<$, в множестве L^* выбирается максимальная строка — это и будет максимальный суффикс u строки x . Но это в точности описание алгоритма Дюваля: для каждого префикса $x[1..j]$ определяются элементы множества L^* , которые заканчиваются в позиции j , т.е. вычисляются суффиксы префикса $x[1..j]$. Выбрать для каждого j наибольший элемент из этого множества просто, поскольку алгоритм Дюваля просматривает строку x слева направо — надо взять самый левый элемент. Таким образом, приходим к следующему алгоритму вычисления максимального суффикса.

Алгоритм 6.2.2 (Алгоритм вычисления максимальных суффиксов)

▷ Для входной строки $x\$$ выводится массив $\sigma''[1..n]$, такой, что
 ▷ подстрока $x[1..j]$ имеет максимальный суффикс $x[i..j]$ только тогда,
 ▷ когда $\sigma''[j] = i$
for $j \in 1..n$ **do** $\sigma''[j] \leftarrow 0$
 $h \leftarrow 0$; $\sigma''[1] \leftarrow 1$
while $h < 0$ **do**
 $i \leftarrow h + 1 - x[i] = \lambda$
 $j \leftarrow h + 2 - x[j] = \mu$
 while $x[j] >_{\mathcal{R}} x[i]$ **do** ▷ выполняется проверка $\mu \leq \lambda$

if $x[j] > \mathcal{R}x[i]$ **then**

$i \leftarrow h + 1$

else

$i \leftarrow i + 1$

if $\sigma''[j] = 0$ **then**

$\sigma''[j] \leftarrow h + 1$

$j \leftarrow j + 1$

▷ случай $\mu > \lambda$, вычисление h

$r \leftarrow \lceil (i - h)/(j - i) \rceil$; $h \leftarrow h + r \times (j - i)$

if $h = j - 1$ **then**

if $\sigma''[j] = 0$ **then**

$\sigma''[j] \leftarrow h + 1$

Как и в случае алгоритма 6.2.1, здесь для компактного представления вычисленных данных используется целочисленный массив $\sigma''[1..n]$, где для любого целого $j \in 1..n$ $\sigma''[j] = i$ только в том случае, если $s_{\max}(x[1..j]) = x[i..j]$. Другими словами, $\sigma''[j] = i$ только тогда, когда $x[i..j]$ является наибольшим суффиксом строки $x[1..j]$ и при этом $s_{\max}(x[i..j]) = x[i..j]$. Таким образом, для каждой позиции j , т.е. для каждого префикса $x[1..j]$, находится наименьшее значение i , такое, что $x[i..j] \in L^*$ в соответствии с отношением $<_{\mathcal{R}}$. Поскольку позиция $h + 1$ всегда обозначает начало элемента множества L^* , который обрабатывается в данный момент, то очевидно предназначение оператора $\sigma''[j] \leftarrow h + 1$. Но так как позиция j может встретиться в нескольких перекрывающихся словах множества L^* , то для того, чтобы использовалось наименьшее подходящее значение переменной h , данный оператор выполняется только после проверки условия $\sigma''[j] = 0$, как показано в листинге.

Покажем пример вычисления максимальных суффиксов для строки f_6 , которая использовалась выше в примере вычисления минимальных суффиксов. Отметим, что в этом примере пробелы отделяют линдонские слова в соответствии с отношением порядка $<_{\mathcal{R}}$. Эти же пробелы отделяют слова в линдонской декомпозиции (в соответствии с отношением $<$) сопряженной строки f'_6 , полученной из строки f_6 путем взаимной замены букв a и b .

	1	2	3	4	5	6	7	8	9	10	11	12	13
$f_6 = a$	b	a	a	b	a	b	a	a	b	a	a	b	
$f'_6 = b$	a	b	b	a	b	a	b	b	a	b	b	a	
$\sigma' = 1$	2	2	2	2	2	5	5	5	5	5	5	5	

Теорема 6.2.10. Алгоритм 6.2.2 корректно вычисляет массив $\sigma''[1..n]$ (максимальные суффиксы для всех непустых префиксов) для заданной строки x за время порядка $\Theta(n)$.

Доказательство. Алгоритмы 6.2.1 и 6.2.2 отличаются только вычислением массивов σ' и σ'' , поэтому асимптотические оценки времени их выполнения совпадают. Корректность вычисления массива σ'' следует из леммы 6.2.9 и суждений, приведших к данному алгоритму. ■

В этом разделе мы рассмотрели несколько вариантов алгоритмов Дюваля, но, как часто бывает с интересными строковыми задачами, существуют и другие алгоритмы, которые вычисляют линдонскую декомпозицию для заданной строки x или решают близкую к этой задаче задачу построения канонической формы для петли $C(x)$ с определением минимальной начальной точки $\text{MSP}(x)$ (см. раздел 1.4). Алгоритмы [36, 208] вычисляют $\text{MSP}(x)$, используя для этого обобщение массива граней. Другой алгоритм [126] вычисления линдонской декомпозиции конкурирует с алгоритмами Дюваля по времени выполнения в среднем, но в самом худшем случае выполняется за время порядка $\Theta(n \log n)$. Специальный алгоритм предложен в работе [207] для сравнения петель $C(x_1)$ и $C(x_2)$, который можно непосредственно применить для решения задачи о конгруэнтности многоугольников, которая рассматривалась в разделе 1.4.

Возможно, наиболее интересный алгоритм вычисления канонических форм предложен в работе [13]. В этой работе также показано, что первый алгоритм Дюваля можно применить для вычисления $\text{MSP}(x)$, причем без увеличения числа буквенных сравнений. Второй алгоритм Дюваля можно сделать онлайн-алгоритмом вычисления минимальных начальных точек для всех префиксов строки x за время $\Theta(n)$ (и при этом потребуется всего в полтора раза больше буквенных сравнений, чем в первом алгоритме Дюваля). Этот алгоритм может вычислить минимальные начальные точки для всех *подстрок* строки x за время $\Theta(n^2)$ при дополнительной памяти $\Theta(n)$.

Упражнения 6.2

1. Покажите, что в алгоритме 6.2.1 операторы $r \leftarrow \lceil (i - h)/(j - i) \rceil$; $h \leftarrow h + r \times (j - i)$ можно заменить одним оператором $h \leftarrow \sigma'[i] + (j - i + 1)$.
2. Необходим ли для алгоритма 6.2.1 во входной строке символ окончания строки \$? Покажите, как следует изменить этот алгоритм, чтобы можно было не использовать этот символ.
3. Покажите, как надо изменить алгоритм 6.1.2, чтобы он вычислял массив σ' .
4. Охарактеризуйте линдонскую декомпозицию строк Фибоначчи f_n .
5. Укажите необходимые и достаточные условия для справедливости следующих выражений.
 - а) $\sigma'[j] = 1$;
 - б) $\sigma'[j] = j$;

- в) $\sigma'[j + 1] = \sigma'[j] + 1$.
6. Докажите лемму 6.2.4.
 7. Докажите лемму 6.2.8.
 8. Покажите, как можно изменить алгоритм 6.1.2, чтобы он вычислял массив σ'' .
 9. Докажите, что $\sigma''[j + 1] \geq \sigma''[j]$ для любого $j \in 1..n - 1$.
 10. В упражнении 5.2.19 введено дерево позиций. Покажите, как дерево позиций можно использовать для вычисления канонических форм петель $C(x)$.

6.3 s -факторизация

В этом разделе рассмотрим другую, отличную от линдонской, декомпозицию строк $x = x[1..n]$, которая зависит от повторяемости подстрок внутри строки x . В отличие от линдонской, данная декомпозиция не зависит от отношения лексикографического порядка и поэтому может применяться в случае неупорядоченных алфавитов.

Определение 6.3.1. Декомпозиция $x = w_1 w_2 \dots w_k$ называется **s -факторизацией**, если каждая подстрока w_j ($j = 1, 2, \dots, k$) является или

- а) буквой, которая не встречается в подстроке $w_1 w_2 \dots w_{j-1}$, или
- б) подстрокой наибольшей длины, которая по крайней мере один раз встречалась в $w_1 w_2 \dots w_{j-1}$. ■

Для строки x s -факторизация строится при просмотре этой строки слева направо, последовательно вычисляя текущий фактор w_j на основе уже вычисленных факторов $w_1 w_2 \dots w_{j-1}$. По определению $w_1 = x[1]$. На j -м шаге, если буква, стоящая в позиции

$$p_j = |w_1| + |w_2| + \dots + |w_{j-1}| + 1,$$

еще не встречалась в подстроке $x[1..p_j - 1]$, то в этом случае применяется часть а) определения 6.3.1. Если же эта буква уже встречалась в подстроке $x[1..p_j - 1]$, тогда, согласно части б) определения, надо найти наибольшую подстроку, которая начинается в позиции p_j и которая уже имеется в подстроке $x[1..p_j - 1]$.

Рассмотрим пример s -факторизации для строки Фибоначчи

$$f_6 = abaababaabaab.$$

Эта факторизация состоит из таких факторов: $w_1 = a$, $w_2 = b$, $w_3 = a$, $w_4 = aba$, $w_5 = baaba$, $w_6 = ab$. Кратко s -факторизацию можно записать как $a/b/a/aba/baaba/ab$.

Заметим, что в случае б) определения 6.3.1 допускается перекрытие предыдущего вхождения фактора w_j с суффиксом этого же фактора. Например, s -факторизацией строки a^n будет a/a^{n-1} .

Интересно отметить, что хотя определение s -факторизации не зависит от природы алфавита, но известные эффективные алгоритмы вычисления s -факторизации используют деревья суффиксов и поэтому зависят от алфавита, который предполагается упорядоченным и, желательно, не слишком большим (см. раздел 2.1).

Чтобы показать, как вычисляется s -факторизация для заданной строки x , предположим, что дерево суффиксов T_x уже построено (см. раздел 5.2). Далее предположим, что каждый внутренний узел i дерева T_x помечен числовой меткой, совпадающей с наименьшей числовой меткой конечных узлов, являющихся потомками данного узла i . Присвоение меток внутренним узлам можно выполнить с помощью стандартной процедуры обхода дерева T_x в обратном порядке, когда каждый из $\Theta(n)$ узлов посещается только один раз [135, 5]. Такие метки внутренних узлов являются номерами самых левых позиций в строке x подстрок, представленных путем от корня дерева до данного внутреннего узла. По соглашению корень дерева помечается меткой 0. Помеченное дерево суффиксов для строки Фибоначчи $f_6 = abaababaabaab\$$ показано на рис. 6.1.

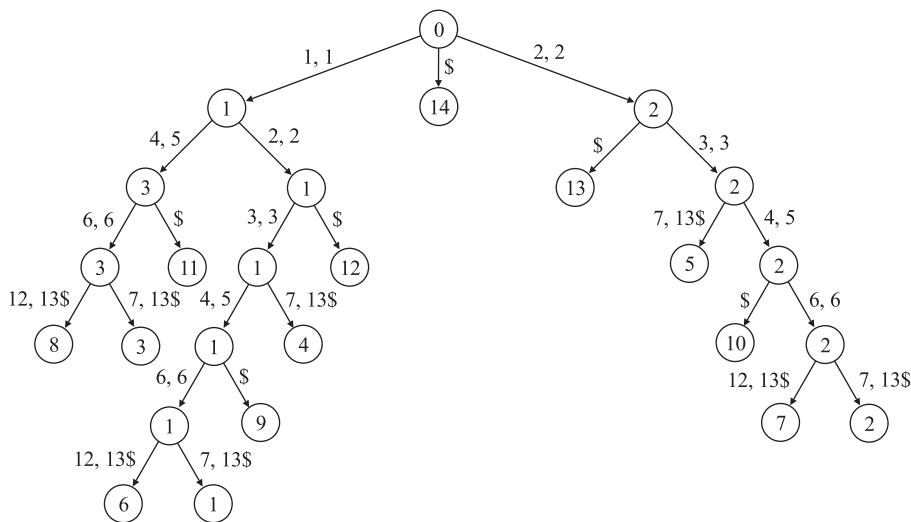


Рис. 6.1. Помеченное дерево суффиксов

Чтобы более детально показать процесс построения s -факторизации, обозначим каждый фактор w_j парой чисел (i_L, ℓ) , где i_L — начальная позиция в строке x первого вхождения повторяющейся подстроки w_j , $\ell = |w_j|$ — длина подстро-

ки w_j . Пусть i_0 — текущая позиция строки x , т.е. позиция, следующая за префиксом $w_1 w_2 \dots w_{j-1}$. Для значения i_0 корректное значение ℓ можно вычислить путем просмотра дерева T_x от корня до конечного узла i_0 . Для этого сначала надо положить $i \leftarrow i_0$, затем проверяется буква $x[i]$, при этом i увеличивается при прохождении нисходящих ребер дерева, ведущих к узлу с меткой i_0 . Если в конечной точке пути узел имеет ненулевую метку, например v , то тогда $i_L \leftarrow v$ и ℓ полагается равной длине подстроки, представленной узлом v . Если же последним узлом пути будет корень дерева (он помечен меткой 0), то полагаем $(i_L, \ell) \leftarrow (i_0, 0)$, что указывает на то, что в позиции i_L обнаружена новая буква. Алгоритм s -факторизации представлен в приведенном ниже листинге и основан на функции сравнения $match$, подробно описать которую предложено в упражнении 6.3.2. Этот алгоритм для строки f_6 должен дать следующую последовательность пар (i_L, ℓ) :

$$(1, 0), (2, 0), (1, 1), (1, 3), (2, 5), (1, 2). \tag{6.3}$$

Алгоритм 6.3.1 (Алгоритм Зива-Лемпела)

```

▷ Вычисление  $s$ -факторизации заданной строки  $x$ 
▷ как последовательности пар  $(i_L, \ell)$ 
▷ на основе помеченного дерева суффиксов  $T_x$ 
 $i_0 \leftarrow 1$ 
while  $i_0 \leq n$  do
     $(i_L, \ell) \leftarrow match(i_0, T_x)$ 
    output  $(i_L, \ell)$ 
     $i_0 \leftarrow i_0 + \ell$ 
    
```

Используя алгоритм Укконена (см. раздел 5.2.3) для построения дерева суффиксов, алгоритм 6.3.1 можно реализовать как онлайн-овый, т.е. s -факторизацию строки x можно вычислять в процессе построения дерева T_x [108, 235].

Чтобы использовать алгоритм 6.3.1 для сжатия данных, его следует немного изменить. Например, чтобы облегчить восстановление данных, первое вхождение каждой буквы необходимо указать точно, а для отделения факторов удобно использовать разделители. С учетом этих “пожеланий”, s -факторизацию строки Фибоначчи $f_6 = abaababaabaab$, полученную в виде последовательности (6.3), следует переписать как

$$a/b/1,1/1,3/2,5/1,2 \tag{6.4}$$

Но в таком случае возрастает количество используемых символов от 13 (количество букв в исходной строке f_6) до 19 в представлении (6.4) (если считать разделители и запятые). В общем случае чем длиннее исходная строка, тем больше проявляется эффект сжатия данных. Например, для следующей строки Фибоначчи

$$f_7 = abaababaabaababaababa$$

длиной 21 s -факторизация

$$a/b/1,1/1,3/2,5/4,8/2,2 \quad (6.5)$$

имеет длину 23, что не на много больше длины факторизации (6.4).

Метод сжатия данных Зива-Лемпена [235] и его модификации основаны на s -факторизации. Один из вариантов этого метода реализован в системе Unix в виде процедур *compress* (сжатие) и *uncompress* (восстановление). Обычно процедура *compress* сжимает длинный текст на английском языке в два или три раза по сравнению с исходным объемом текста. В упражнении 6.3.5 предложено детализировать алгоритм восстановления исходной строковой последовательности на основе имеющейся s -факторизации, подобной (6.5).

В разделе 12.2 мы рассмотрим применение s -факторизации для нахождения кратных подстрок в данной строке x .

Упражнения 6.3

1. Докажите, что любая строка x имеет только одну s -факторизацию.
2. Запишите функцию *match*, которая используется в алгоритме 6.3.1. Затем докажите корректность этого алгоритма.
3. Постройте дерево суффиксов T_x для строки $x = abaababcabab$ и затем вычислите s -факторизацию для этой строки.
4. Запишите s -факторизацию для строки $x = (ab)^n$, $n \geq 2$, в форме (6.4). Каков коэффициент сжатия в данном случае, т.е. чему равно отношение количества символов в s -факторизации к длине строки $2n$?

Совет. Предположите, что любой символ как в строке x , так и в ее s -факторизации требует для хранения полбайта (четыре бита).

5. Разработайте алгоритм процедуры *uncompress*, выполняющий за линейное время восстановление строки x на основании ее s -факторизации.
6. Иногда используется другая форма s -факторизации, где в определении 6.3.1 условие б) заменено следующим условием б):

б) лексикографически максимальной подстрокой, которая по крайней мере один раз встречалась в $w_1 w_2 \dots w_{j-1}$.

Перепишите функцию *match* и алгоритм процедуры *uncompress* в соответствии с новым определением s -факторизации. Какой коэффициент сжатия в этом случае будет для строки $x = (ab)^n$?

ЧАСТЬ III

Вычисление частных паттернов

Не хлебом одним будет жить человек, но всяким словом,
исходящим из уст Божиих

— Матфей 4.4

Компьютерная наука, интеллектуальная научная дисциплина, без сомнения, восходит к Алану Тьюрингу (1935 год), но распознать эту науку в общем случае несложно: как заметил Винни-Пух, время определяют вещи. С такой точки зрения вполне уверенно можно утверждать, что наша дисциплина сформировалась примерно 40 лет назад. Прилагательное “классический” (и его производные), которое часто используется в таких науках и искусствах, как физика, музыка или математика, и которое указывает на достижения, достигнутые в этих науках и искусствах известными и прославленными первопроходцами сотни, а иногда и тысячи лет назад, кажется совсем неуместным применительно к компьютерной науке. В этом случае фразы “классический подход” или “классический алгоритм” будут относиться к работам, выполненным менее 25 лет назад людьми, которые до сих пор живы и продолжают плодотворно работать! То же самое можно сказать о той области компьютерной науки, которая изучается в данной книге: самые ранние работы по вычисляемым паттернам строковых последовательностей относятся к 60-м годам, интерес к вычисляемым частным паттернам (в виде “классического подхода” сравнения с паттерном) проявился в 70-х годах, в частности, после публикаций в 1977 году “классических” алгоритмов Бойера–Мура и Кнута–Морриса–Пратта, авторы которых продолжают творить среди нас.

После такого разъяснения представим материал данной части как “классический”: задачи сравнения с паттернами были первыми изученными задачами этой области компьютерной науки и сегодня продолжают оставаться в поле пристального внимания исследователей. То, что эта область продолжает интенсивно развиваться, является ее преимуществом, а не недостатком, поскольку расширяется сфера применения данной области. В главах 7 и 8 рассмотрим алгоритмы локализации точных единичных паттернов в строковых последовательностях. В главах 9 и 10 эти алгоритмы будут “расширены” для выполнения сравнения с “приближенными” паттернами (как это определено в разделе 2.2). Наконец, в главе 11 рассмотрим применение конечных автоматов для сравнения с паттернами, в частности для сравнения регулярных выражений и для одновременной локализации в строке нескольких различных паттернов, причем как для точного, так и приближенного сравнения.

ГЛАВА 7

Базовые алгоритмы

Слова, без сомнения, наиболее мощный наркотик, который использует человечество.

— Редьярд Киплинг (1865–1936).
The Times, 15 февраля 1923

В этой главе мы рассмотрим четыре различных “классических” подхода к вычислению всех вхождений заданного непустого паттерна $p = p[1..m]$ в заданную непустую строку $x = x[1..n]$. Первые два подхода (правда, второй после небольшой модификации) гарантируют выполнение соответствующих алгоритмов за время порядка $O(n + m)$. Другие два не гарантируют такой скорости вычислений, однако обладают другими полезными свойствами, которые делают их конкурентоспособными в определенных ситуациях. В следующей главе мы систематизируем алгоритмы точного сравнения с паттернами, где будет представлено несколько десятков разработанных в последние 20 лет вариантов базовых алгоритмов, описанных в данной главе. Превосходный обзор алгоритмов сравнения с паттернами имеется на Web-узле [47], который включает краткое их описание, программы на языке C и апплеты Java.

7.1 Алгоритм Кнута–Морриса–Пратта

Основная идея алгоритма КМП (так для краткости будем называть алгоритм Кнута–Морриса–Пратта, Knuth–Morris–Pratt, [137]) заключается в последовательном сдвиге паттерна p вдоль строковой последовательности x и сравнении его

с просматриваемым участком последовательности x так, как это делается в простом алгоритме 2.2.1 (см. раздел 2.2). При этом паттерн p сдвигается на следующую букву вдоль строки x только в том случае, когда p не совпадает с текущим участком последовательности x .

Сначала рассмотрим этот алгоритм на примере, который уже использовался в разделе 2.2: $x = a^n$, $p = a^{m-1}b$, $2 \leq m \leq n$. В этом примере после сравнения $a^{m-1}b$ с подстрокой $x[1..m]$ мы будем располагать информацией, что

$$x[1..m-1] = p[1..m-1] = a^{m-1}.$$

Чтобы проверить, будет ли паттерн p совпадать с подстрокой $x[2..m+1]$, нет необходимости проводить сравнения в подстроке $x[2..m-1]$, поскольку уже известно, что $x[2..m-1] = a^{m-2}$. Поэтому следует провести только одно сравнение $x[m] : p[m-1]$, и в случае равенства этих букв, мы фиксируем несовпадение

$$x[m+1] \neq b = p[m].$$

Применяя эту стратегию в каждой позиции $i \leq n - m + 1$, можно показать (см. упражнение 7.1.4), что выполняется всего $2n - m$ буквенных сравнений, а не $m(n - m + 1)$, как можно было ожидать.

На этом примере показана основная особенность алгоритма КМП [184, 137]: сравнения

$$x[i..i+h-1] = p[1..h], 1 \leq h \leq m,$$

можно (частично или полностью) использовать для избежания дальнейшей проверки подстроки $u = x[i..i+h-1]$. Если паттерн p входит в строку x в некоторой позиции $i' \in i..i+h-1$, тогда строка $u = p[1..h]$ имеет непустой собственный префикс $p[1..h']$, $h' = i+h-i'$, а также суффикс $x[i'..i+h-1]$. Другими словами, подстрока $p[1..h']$ должна быть *гранью* строки $u = p[1..h]$!. Эта ситуация проиллюстрирована на рис. 7.1.

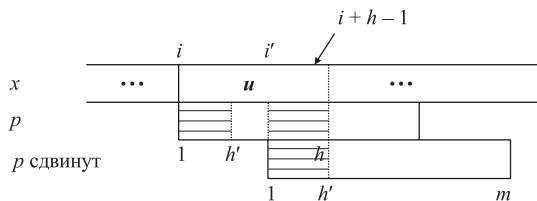


Рис. 7.1. Грань строки $p[1..h]$

Таким образом, с помощью граней строки u надо определить только те позиции i' , в которых может встретиться p . Но поскольку u — это подстрока $p[1..h]$ паттерна p для некоторого $h \in 1..m$, то нет необходимости вычислять грани строки

u в процессе выполнения алгоритма — можно вычислить их *заранее* и сохранить в массиве граней $\beta = \beta[1..m]$. Поэтому алгоритм КМП имеет предварительный этап, на котором вычисляется массив β . Такие вычисления можно выполнить за время $\Theta(m)$, если использовать для этого алгоритм 1.3.1.

Фактически надо предварительно вычислить не сам массив граней, а его более слабый вариант. Обозначим через j позицию в паттерне p первого несовпадения. Тогда префикс паттерна p , который не требует проверки при следующем сравнении, — это наибольшая грань подстроки $p[1..j-1]$. Поэтому следующей буквой для сравнения будет $p[\beta[j-1]+1]$, где $\beta[j-1]$ — длина наибольшей грани подстроки $p[1..j-1]$, при этом, конечно, $\beta[j-1] < j-1$. Следовательно, в действительности для всех $j \in 2..m+1$ надо вычислить $\beta'[j] = \beta[j-1] + 1$, положив при этом $\beta'[1] = 0$. Отметим, что для $j \geq 2$ $\beta[j-1] + 1 < (j-1) + 1 = j$, поэтому $\beta'[j] < j$. Если $j = 1$, тогда сравнивать нечего, и в этом случае необходимо изменить значение i , т.е. сдвинуть паттерн вдоль строки x . После этих пояснений мы можем записать алгоритм КМП (см. следующий листинг). Заметим, что оператор присвоения $j \leftarrow \beta'[j]$, выполняемый в результате несовпадения $p[j] \neq x[i]$, соответствует сдвигу вправо вдоль строки x паттерна p на $j - \beta'[j]$ позиций.

Алгоритм 7.1.1 (Алгоритм Кнута–Морриса–Пратта)

```

▷ Нахождение всех вхождений паттерна  $p$  в строку  $x$ 
 $i \leftarrow 1; j \leftarrow 1$ 
while  $i \leq n - m + j$  do
  ▷  $(i, j) \leftarrow match(i, j, m)$ 
  if  $j = m + 1$  then output  $i - m$ 
  if  $j = 1$  then  $i \leftarrow i + 1$ 
  else  $j \leftarrow \beta'[j]$ 
    
```

Отметим, что в этом алгоритме процедура *match* (сравнить) выполняется так же, как и ее двойник в алгоритме 2.2.1, т.е. она выполняет сравнение “буква за буквой” слева направо. Однако есть между ними и небольшое отличие: в алгоритме 2.2.1 процедура *match* необязательно начинает сравнение с первой буквы $p[1]$, тогда как в алгоритме КМП она начинает сравнение с буквы $x[i]$ с заданной позиции j , $1 \leq j \leq m$, т.е. с буквы $p[j]$. Кроме того, здесь она возвращает как значение i , так и значение j , показывающее позицию несовпадения букв в x и p .

Обсуждение алгоритма КМП мы подытоживаем в следующей теореме.

Теорема 7.1.1. Алгоритм 7.1.1 корректно вычисляет все вхождения заданного непустого паттерна p в заданную строку x . ■

Чтобы оценить сложность алгоритма КМП, подсчитаем количество буквенных сравнений, выполняемых в процедуре *match*. Если результат сравнения показывает равенство $x[i] = p[j]$, тогда значения переменных i и j увеличиваются на еди-

ницу, при этом паттерн p не сдвигается относительно x . Если результат сравнения показывает неравенство $x[i] \neq p[j]$, тогда для случая $j = 1$ значение i увеличивается на единицу и паттерн p сдвигается вправо на одну позицию, а в случае $j > 1$ значение i не меняется, паттерн p сдвигается вправо на $j - \beta'[j] \geq 1$ позиций. Заметим, что переменная i может возрасти не более n раз, тогда как паттерн p сдвигается вправо не более $n - m + 1$ раз. Однако эти две верхние границы не могут достигаться одновременно, поскольку цикл **while** закончится как только выполнится хотя бы одно из этих условий. Следовательно, общее количество буквенных сравнений не может превысить величины

$$n + (n - m + 1) - 1 = 2n - m.$$

Поскольку все другие операции в цикле **while** выполняются за константное время, мы приходим к следующей теореме.

Теорема 7.1.2. Алгоритм 7.1.1 выполняется за время порядка $O(n)$ с использованием дополнительной памяти объемом $\Theta(m)$, при этом выполняется не более $2n - m$ сравнений с элементами строки x . ■

Пример с $x = a^n$ и $p = a^{m-1}b$ показывает, что верхняя граница числа буквенных сравнений наименьшая из возможных. Если к времени выполнения алгоритма 7.1.1 добавить время $\Theta(m)$ выполнения предварительного этапа, то получим полное время порядка $\Theta(n+m)$, необходимое для поиска всех вхождений паттерна p в строку x . В контексте обсуждения 1.2.1 заметим, что эта верхняя временная граница зависит от возможности вывода каждого вхождения p за константное время, используя для локализации p в строке x только целые числа.

Рассмотрим также поучительный пример с $x = a^n$ и $p = a^m$, где паттерн $n - m + 1$ раз входит в строку x . Поскольку паттерн p имеет грань a^{m-1} , то $\beta'[m+1] = m$, и поэтому, для того чтобы определить каждое вхождение паттерна p после первого, требуется всего одно буквенное сравнение $p[m] : x[i]$. Этот пример показывает замечательное свойство алгоритма КМП: необходимо только n буквенных сравнений для вывода $n - m + 1$ локализаций вхождений паттерна p в строку x — другими словами, если значение n растет значительно быстрее m , то имеем константное время нахождения одного вхождения паттерна. Как мы увидим в следующем разделе, этот пример (периодический паттерн с малым периодом входит много раз в строковую последовательность) доставит большие трудности алгоритму Бойера–Мура, основному конкуренту алгоритма КМП.

Описав общий алгоритм КМП, можно приступить к его усовершенствованию, точнее, к усовершенствованию вычисления вспомогательного массива β' . Напомним, что в алгоритме 7.1.1 ситуация, когда $x[i] \neq p[j]$, $j > 1$, ведет к выполнению оператора $j \leftarrow \beta'[j]$. Но выполнение этого оператора бесполезно, если $p[j] = p[\beta'[j]]$, поскольку тогда на следующей итерации по-прежнему будет выполняться неравенство $x[i] \neq p[j]$! В действительности нас интересует наибольшая

грань подстроки $p[1..j-1]$, за которой *не следует* буква $p[j]$ (если, конечно, такая грань существует). Другими словами, надо найти наибольшее целое число j' (если оно существует), такое, что

- подстрока $p[1..j'-1]$ является гранью строки $p[1..j-1]$;
- подстрока $p[1..j']$ *не является* гранью строки $p[1..j]$.

Это подводит нас к построению нового вспомогательного массива β'' , j -й элемент которого вычисляется по следующим правилам.

Если $j = m + 1$, то $\beta''[j] \leftarrow \beta'[j]$. Для значений $j \in 1..m$ $\beta''[j] \leftarrow j'$, где $j' - 1$ — длина наибольшей грани подстроки $p[1..j-1]$, такой, что $p[j'] \neq p[j]$. Если такой грани не существует, то $\beta''[j] \leftarrow 0$.

Напомним (см. раздел 1.3), что грани подстроки $p[1..j-1]$ можно определить посредством их длин, которые являются элементами убывающего множества

$$S_{j-1} = \{\beta[j-1], \beta^2[j-1], \dots, \beta^k[j-1]\}, \quad (7.1)$$

где k — целое число, такое, что $\beta^k[j-1] = 0$ (число k определяется единственным образом). Тогда для любого $j \in 2..m$ ненулевое значение j' элемента $\beta''[j]$ возможно только в том случае, если существует значение (длина) $j' - 1 \in S_{j-1}$, которое удовлетворяет условиям, перечисленным выше.

Отметим, что приведенное выше определение элементов массива β'' применимо даже тогда, когда длина $j' - 1$ наибольшей подходящей грани равна нулю. В этом случае для некоторого $j < m + 1$ и *любой* грани длиной $j' - 1$ выполняется равенство $p[j'] = p[j]$. Особое значение $\beta''[j] = 0$ показывает, что следующее сравнение будет выполняться между буквами $x[i+1]$ и $p[1]$, т.е. это тот случай, когда процедура *match* возвращает $j = 1$. Детали, связанные с минимальными изменениями в алгоритме 7.1.1, необходимые для преобразования массива β (или массива β') в массив β'' , показаны в упражнениях 7.1.9 и 7.1.10. Назовем таким образом модифицированный алгоритм 7.1.1 алгоритмом КМП*. В упражнении 7.1.11 предложено показать, что алгоритм КМП* в самом худшем случае выполняется так же, как и алгоритм КМП. Поэтому справедлива такая теорема.

Теорема 7.1.3. Алгоритм КМП* корректно вычисляет все вхождения заданного непустого паттерна $p[1..m]$ в заданную строку $x[1..n]$ за время порядка $O(n)$ с использованием дополнительной памяти объемом $\Theta(m)$, при этом выполняется не более $2n - m$ буквенных сравнений. ■

Интересно посмотреть, какой вид имеют массивы β , β' и β'' для строк Фибоначчи. Мы увидим, что здесь различия между массивами β' и β'' весьма значительны. В частности, $\beta''[j] = 0$ для всех позиций j вхождения буквы a , следующей за буквой b , поскольку каждая строка Фибоначчи начинается с буквы a

и не содержит подстрок bb . Для строки Фибоначчи f_7 массивы β , β' и β'' имеют следующий вид.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
$f_7 =$	a	b	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b	a	b	a	
$\beta =$	0	0	1	1	2	3	2	3	4	5	6	4	5	6	7	8	9	10	11	7	8	
$\beta' =$	0	1	1	2	2	3	4	3	4	5	6	7	5	6	7	8	9	10	11	12	8	9
$\beta'' =$	0	1	0	2	1	0	4	0	2	1	0	7	1	0	4	0	2	1	0	12	0	9

Алгоритмы КМП и КМП* просты, элегантны и интересны, поэтому они стимулировали многочисленные исследования методов вычисления внутренних паттернов в строковых последовательностях. Однако эти алгоритмы, к сожалению, не эффективны: показано, как теоретически [22], так и практически [212, 76], что алгоритм КМП* в среднем лишь чуть быстрее простого алгоритма 2.2.1, хотя в самом худшем случае он, конечно, выполняется значительно быстрее. В следующем разделе мы покажем другие подходы построения алгоритмов сравнения с паттернами, которые эффективны в среднем.

Упражнения 7.1

1. Напишите процедуру *match* для алгоритма КМП и докажите ее корректность.
2. Предположим, что строки x и p заканчиваются специальными буквами (которые называются *сигнальными метками* (sentinels)) S_1 и S_2 соответственно. Эти буквы различны между собой и не совпадают ни с одной буквой алфавита. Перепишите процедуру *match* так, чтобы она могла распознавать и использовать сигнальные метки. Оцените эффективность новой процедуры *match* по сравнению с аналогичной процедурой из предыдущего упражнения.
3. Предположим, что алгоритм КМП применяется к строке $x = a^n$ и паттерну $p = a^m$, $1 \leq m \leq n$, с использованием
 - а) процедуры *match* из предыдущего упражнения,
 - б) процедуры *match* из первого упражнения.

В тексте показано, что в случае б) для решения задачи необходимо n буквенных сравнений. Оцените количество буквенных сравнений, необходимых в случае а). На этом основании сделайте заключение об использовании минимального количества буквенных сравнений как критерия эффективности алгоритмов сравнения с паттернами.

4. Пусть алгоритм КМП (с исходной процедурой *match*) применяется для нахождения всех вхождений паттерна $p = a^{m-1}b$ в строку $x = a^n$, $2 \leq m \leq n$.

Покажите, что в этом случае алгоритм выполнит ровно $2n - m$ буквенных сравнений.

5. Докажите утверждение, сделанное в тексте, что в алгоритме КМП не могут одновременно достигаться верхняя граница увеличения переменной i и верхняя граница величины сдвига паттерна.
6. В тексте ничего не сказано о *нижней границе* количества буквенных сравнений, выполняемых в алгоритме КМП. Покажите, что эта нижняя граница равна $n - m + 1$. Что вы можете сказать о *среднем* количестве буквенных сравнений, выполняемых в алгоритме КМП, с учетом предварительной обработки паттерна p ?
7. Сколько буквенных сравнений потребуется алгоритму КМП для нахождения всех вхождений паттерна $p = (\lambda_1 \lambda_2 \dots \lambda_k)^{m/k}$ в строку $x = (\lambda_1 \lambda_2 \dots \lambda_k)^{n/k}$, если все буквы $\lambda_1, \lambda_2, \dots, \lambda_k$ различны?
8. Докажите, что смещение $\beta''[j]$, соответствующее несовпадению сравниваемых букв в позиции j , равно смещению $\beta''[\beta'[j]]$, соответствующему несовпадению букв в позиции $\beta'[j]$, только в том случае, если $p[j] = p[\beta'[j]]$. Отсюда будет следовать, что $\beta''[j] = \beta'[j]$ или $\beta''[j] = \beta''[\beta'[j]]$.
9. Предположим, что к паттерну p справа добавлена уникальная буква $\$$. Модифицируйте алгоритм 1.3.1 таким образом, чтобы он за время порядка $\Theta(n)$ вычислял массив β' . Другая модификация этого алгоритма за такое же время должна вычислять массив β'' . Докажите корректность новых алгоритмов, основываясь на предположении о корректности алгоритма 1.3.1.
Совет. Алгоритм для вычисления массива β' должен быть проще, чем алгоритм 1.3.1. На основе этого алгоритма, используя соотношения из предыдущего упражнения, несложно построить алгоритм для вычисления массива β'' . Можно ли вычислить массив β'' без предварительного вычисления (и сохранения) массива β' ?
10. Измените алгоритм КМП так, чтобы он использовал массив β'' , а не массив β' .
11. Вычислите массив β'' для паттерна $p = a^{m-1}b$. Затем покажите, что в самом худшем случае алгоритм КМП* требует не более $2n - m$ буквенных сравнений.
12. Вычислите массивы β' и β'' для
 - а) $p = abcdabce$;
 - б) $p = abcabcacab$;
 - в) $p = abacabadacabae$.
 Существуют ли строки, для которых $\beta'' = \beta'$?

7.2 Алгоритм Бойера–Мура

Подобно алгоритму КМП, алгоритм БМ (так для краткости будем называть алгоритм Бойера–Мура, Boyer–Moore, [38]) сдвигает паттерн p вдоль текстовой строки x слева направо. Однако, в отличие от алгоритма КМП, в алгоритме БМ буквы в паттерне p сравниваются *справа налево*, т.е. в порядке $p[m]$, $p[m - 1]$, \dots , $p[1]$. Как мы увидим далее, это, на первый взгляд, незначительное изменение коренным образом изменит стратегию определения, в случае несовпадения сравниваемых букв, следующего положения паттерна p относительно строки x . Как и алгоритм КМП, данный алгоритм зависит от предварительно вычисленных массивов, которые позволяют определить каждый сдвиг паттерна p за константное время. Далее мы увидим, что в алгоритме БМ такие массивы тесно связаны с массивом β'' , описанным в предыдущем разделе.

Прежде чем детально описывать алгоритм БМ, рассмотрим присущие ему ограничения. На практике встречается много задач сравнения с паттерном, которые чувствительны к этим ограничениям.

- Один из предварительно вычисляемых массивов алгоритма БМ (обычно этот массив обозначается как δ_1) зависит от знания вхождений в паттерн p каждой буквы алфавита строки x . Поэтому алфавит строки x должен быть конечным и известным заранее.
- Доступ к элементам массива δ_1 основан не на индексах элементов, а на буквах, т.е. к элементам этого массива обращаются как $\delta_1[\lambda]$. Для того чтобы время доступа к этим элементам было константное, необходимы предварительные вычисления для определения буквы λ . Это означает, что алфавит должен быть *индексированным* (терминология раздела 4.1). Напомним, что если алфавит упорядочен, но не индексированный, то может потребоваться время порядка $\Theta(\log m)$ для определения буквы λ в паттерне $p[1..m]$ и такое же время для доступа к элементам массива δ_1 . Иногда это может быть существенным ограничением для применения алгоритма БМ.
- Поскольку паттерн p просматривается справа налево, необходимо организовать доступ к элементам паттерна в обратном порядке. Как указывалось в обсуждении 1.2.1, может потребоваться линейное время для копирования элементов паттерна p в новый массив в обратном порядке.

Теперь приведем листинг алгоритма БМ, а затем рассмотрим этот алгоритм подробно.

Алгоритм 7.2.1 (Алгоритм Бойера–Мура)

```

▷ Нахождение всех вхождений паттерна  $p$  в строку  $x$ 
 $i \leftarrow m$ 
while  $i \leq n$  do
     $(i, j) \leftarrow hctam(i, m)$ 
    
```

```

if  $j = 0$  then output  $i + 1$ 
 $i \leftarrow i + \max\{\delta_1[x[i]], \delta_2[j]\}$ 
    
```

Надеюсь, что читатель простит мне некоторую вольность за присвоение имени *hctam*¹ процедуре, которая сравнивает паттерн p с буквами строки x , выполняя это сравнение справа налево. Отметим, что здесь сравнение с p заканчивается, когда переменная j , первоначально равная m , уменьшается до нуля, тогда как в алгоритме КМП переменная j возрастала до значения $m + 1$ начиная с некоторого начального значения — с этой точки зрения процедура *hctam* немного проще процедуры *match*, поскольку не зависит от предыдущего значения j . Как видно, алгоритм БМ прост и корректен, если, конечно, корректно работает процедура *hctam* и корректно вычисляется величина сдвига паттерна (что требует корректного вычисления массивов δ_1 и δ_2). Обобщим это утверждение в виде теоремы.

Теорема 7.2.1. При условии корректной реализации процедуры *hctam* и массивов δ_1 и δ_2 алгоритм 7.2.1 корректно вычисляет все вхождения данного непустого паттерна $p[1..m]$ в заданную строку $x[1..n]$. ■

Основная идея массива δ_1 проста. Мы предполагаем, что процедура *hctam* находит совпадение (возможно, на пустых подстроках)

$$x[i + 1..i + h] = p[m - h + 1..m]$$

для некоторого $i \in 0..n - h$ и некоторого $h \in 0..m$ и затем возвращает пару чисел $(i, j) = (i, m - h)$. По соглашению положим, что строка x и паттерн p начинаются с сигнальных меток, т.е. $x[0] = \$_1$ и $p[0] = \$_2$. Тогда можно утверждать, что для любого случая, включая $i = 0$ и $j = 0$, несовпадение

$$\lambda = x[i] \neq p[j] = \mu \tag{7.2}$$

будет найдено. Сделаем очевидное наблюдение: если какая-либо подстрока строки x , содержащая позицию i , совпадает с паттерном p , тогда в p входит буква $\lambda = x[i]$. Далее, если мы знаем позицию (обозначим ее j') самого правого вхождения буквы λ в паттерн p , тогда мы можем определить *минимальный* сдвиг p , удовлетворяющий условию, что $x[i] = \lambda$. Как показано на рис. 7.2, в этом случае $j' < j$, что ведет к сдвигу паттерна на $j - j'$ позиций. Тогда при следующем вызове процедуры *hctam* она начнет работу со сравнения букв $p[m]$ и $x[i + (m - j')]$. Таким образом, в этом случае надо изменить значение i :

$$i \leftarrow i + (m - j'). \tag{7.3}$$

¹Если читатель еще не догадался, то *hctam* — это перевернутое слово *match*. — *Примеч. пер.*

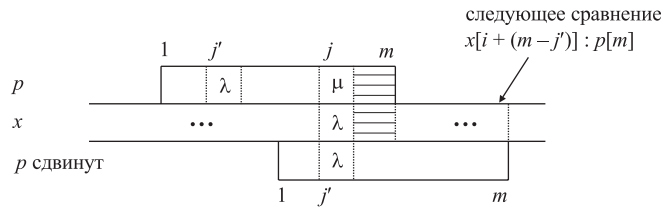


Рис. 7.2. Паттерн сдвинут так, чтобы буква λ в строке x совпала с самым правым вхождением этой буквы в паттерн

Итак, очевидно, что для каждой буквы λ (включая $\$1$) алфавита A необходимо вычислить позицию j' ее самого правого вхождения в паттерн p , положив $j' = 0$, если данная буква λ не входит в p . Тогда на предварительном этапе массив δ_1 вычисляется как

$$\delta_1[\lambda] \leftarrow m - j'. \tag{7.4}$$

Теперь, если обнаружится несовпадение (7.2) и при этом выполняется неравенство $j' < j$ (т.е. $\delta_1[\lambda] > m - j$), следующую позицию для сравнения можно вычислить так, как это сделано в алгоритме 7.2.1: $i \leftarrow i + \delta_1[\lambda]$.

Использование массивов δ_1 и δ_2

“Все это хорошо, — подумает с некоторым сомнением внимательный читатель, — но что делать в случае $\delta_1[\lambda] < m - j$?” Ответ прост: такая ситуация возможна только для $j < m$, но тогда *обязательно* $j \leq 1$. Все, что мы можем сделать в этой ситуации, — это сдвинуть паттерн p вправо на одну позицию так, чтобы следующее сравнение начиналось со сравнения буквы $p[m]$ с буквой $x[i + (m - j) + 1]$. Это эквивалентно присваиванию

$$i \leftarrow i + (m - j) + 1. \tag{7.5}$$

Сравнивая присваивания (7.3) и (7.5), находим, что в случае $j' > j$ (т.е. в случае $\delta_1[\lambda] < m - j$) выполняется неравенство $(m - j) + 1 > m - j'$. Но, как мы покажем ниже, для всех j справедливо неравенство

$$\delta_2[j] \geq (m - j) + 1. \tag{7.6}$$

Поэтому присвоение (7.5) теряет необходимость, если использовать присвоение

$$i \leftarrow i + \max\{\delta_1[x[i]], \delta_2[j]\}$$

из алгоритма 7.2.1.

Подобно массиву δ_1 , массив $\delta_2 = \delta_2[0..m]$ “вступает в игру” после того, как определено совпадение подстрок

$$x[i + 1..i + h] = p[m - h + 1..m]$$

для некоторого $i \in 0..n - h$ и некоторого $h \in 0..m$ и затем определено несовпадение (7.2). (Здесь, как и ранее, $j = m - h$.) Но с другой стороны, вычисление массива δ_2 имеет много общего с вычислением массива β'' из алгоритма КМП*, поскольку этот массив так или иначе связан с вычислением граней, но с привлечением дополнительных ограничений на грани. Если быть более точным, то после определения совпадения подстрок необходимо найти наибольшее целое j' ($j' < j$), такое, чтобы выполнялось совпадение подстрок

$$p[j' + 1..j' + h] = p[j + 1..m]$$

и при этом $\mu' = p[j'] \neq p[j] = \mu$. Если такое j' существует, то происходит присвоение

$$i \leftarrow i + (m - j'), \tag{7.7}$$

как показано на рис. 7.3. Конечно, величину сдвига $\delta_2[j] = m - j'$ можно вычислить заранее для любого $j \in 1..m$. Обращаем ваше внимание на определенную схожесть рис. 7.2 и 7.3.

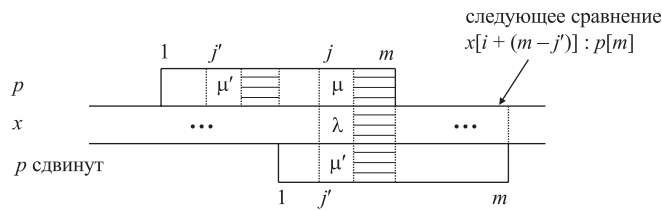


Рис. 7.3. Сдвиг первого типа

Если такого j' , удовлетворяющего перечисленным условиям, не существует (включая случай $j = 0$), тогда ищется такое наибольшее целое j' , что подстрока $p[1..j']$ является суффиксом строки $p[j + 1..m]$, т.е. $p[1..j']$ должна быть наибольшей гранью паттерна p , при этом $j' \leq h$. В этом случае, как показано на рис. 7.4, необходимо выполнить присвоение

$$i \leftarrow i + (m - j) + (m - j'). \tag{7.8}$$

Снова величины $\delta_2[j] = (m - j) + (m - j')$ можно вычислить заранее для любого j .

Необходимо потребовать, чтобы вычисление величин $\delta_2[j]$, выполненное согласно рис. 7.3 и 7.4, удовлетворяло следующим двум условиям.

- Вычисление $\delta_2[j]$ должно быть *осторожным*, т.е. сдвиг паттерна должен быть минимальным при выполнении условия, что собственная подстрока паттерна совпадает с подстрокой $x[i..i + h - 1]$ (рис. 7.3) или что собственный префикс паттерна совпадает с суффиксом подстроки $x[i..i + h - 1]$ (рис. 7.4). Минимум сдвига гарантирован тем, что, в соответствии с присвоениями (7.7) и (7.8), выбирается наибольшее значение j' .

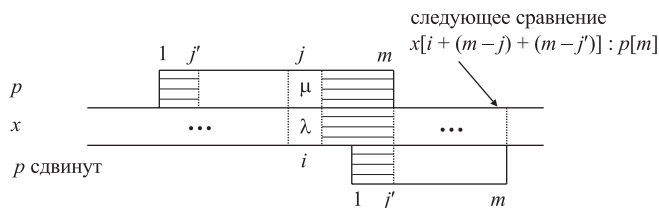


Рис. 7.4. Сдвиг второго типа

- Для любого j должно выполняться неравенство (7.6) $\delta_2[j] \geq (m - j) + 1$. Отсюда вытекает присвоение (7.7) при $j' < j$ и присвоение (7.8) при $m - j' \geq 1$. Поэтому достаточно сделать сдвиг на величину $\max\{\delta_1[\lambda], \delta_2[j]\}$. Отметим, что такая формула сдвига гарантирует, что на каждом шаге алгоритма БМ паттерн будет сдвигаться не менее, чем на одну позицию, т.е. алгоритм обязательно завершится.

Мы показали, как используются массивы δ_1 и δ_2 , и убедились, что при таком их использовании алгоритм 7.2.1 будет выполняться корректно.

Вычисление массива δ_2

Теперь вернемся к вычислению массива δ_2 . Здесь существенно, что позиция j' вычисляется как функция от заданной позиции $j \in 0..m$. Эти вычисления разбиваются на две части. Первая часть вычислений, которая соответствует сдвигу **первого типа**, показанного на рис. 7.3, выполняется только для тех значений j , для которых существует соответствующее значение j' (другими словами, эта часть вычислений выполняется тогда, когда существует соответствующая грань, которой предшествует буква μ' , не равная μ). Вторая часть вычислений, которая соответствует сдвигу **второго типа**, показанного на рис. 7.4, выполняется для тех значений j , соответствующие значения j' которых не могут быть вычислены в первой части вычислений. Чтобы понять различие между этими частями вычислений, предположим, что в результате вычислений первого типа получаем массив $g_I = g_I[1..m]$, где $j' = g_I[j]$ — значение, соответствующее значению j (если, конечно, такое значение j' существует), а в результате вычислений второго типа получаем массив $g_{II} = g_{II}[1..m]$. Вычисление массива g_{II} проще, поэтому начнем рассмотрение с этого массива.

Предположим, что с помощью алгоритма 1.3.1 за линейное время вычислен массив граней $\beta = \beta[1..m]$ паттерна $p[1..m]$. Для данного значения $j \in 1..m$, как ясно из рис. 7.4, значение $j' = g_{II}[j]$ равно длине наибольшей грани паттерна, причем эта длина не должна превышать $h = m - j$. Чтобы определить это значение, надо просто иметь доступ к элементам убывающей последовательности

$$S_m = \langle \beta[m], \beta^2[m], \dots, \beta^k[m] = 0 \rangle,$$

определенной в разделе 1.3. Тогда для каждого j , такого, что $0 < j \leq m - \beta[m]$, положим $g_{II}[j] = \beta[m]$. Для значений j из интервала $m - \beta[m] < j \leq m - \beta^2[m]$ положим $g_{II}[j] = \beta^2[m]$ и т.д. Для $j = 0$ (если имеется совпадение с паттерном) положим $g_{II}[j] = \beta[m]$. Очевидно, что массив g_{II} можно вычислить за линейное время.

Для вычисления значений $g_I[j]$, $1 \leq j \leq m$, необходимо найти наибольшее целое j' (если оно существует и меньше j), такое, что

- подстрока $p[j' + 1..j' + (m - j)]$ является гранью строки $p[j' + 1..m]$,
- подстрока $p[j'..j' + (m - j)]$ не является гранью строки $p[j'..m]$.

Если такое значение j' существует, тогда $g_I[j] \leftarrow j'$, иначе $g_I[j] \leftarrow 0$.

Чтобы лучше понять эти вычисления, рассмотрим грани суффиксов паттерна p . Читатель с хорошей памятью, конечно же, помнит упражнение 1.3.10, где введено понятие “правого массива граней”. Обозначим этот массив как $\rho = \rho[1..m]$, где $\rho[j]$ — длина наибольшей грани подстроки $p[j..m]$. Напомним (также из упражнения 1.3.10), что если β_T — массив граней транспонированной (обратной) строки

$$p_T = p[m]p[m - 1] \dots p[1],$$

тогда $\beta_T[m - i + 1] = \rho[i]$ для всех $i \in 1..m$. Следовательно, длины граней суффиксов $p[i..m]$ являются элементами убывающей последовательности

$$S_{m-i+1}^T = \langle \beta_T[m - i + 1], \beta_T^2[m - i + 1], \dots, \beta_T^k[m - i + 1] = 0 \rangle.$$

Тогда для вычисления элемента $g_I[j]$ необходимо наибольшее целое число j' такое, что $j' < j$ и число $m - j$ принадлежит последовательности $S_{m-j'}^T$, но число $m - j + 1$ не является элементом последовательности $S_{m-j'+1}^T$. Отметим, что если b^* — наибольший элемент массива β_T , тогда для любого j , такого, что $m - j > b^*$, $g_I[j] \leftarrow 0$.

Вычисление массива g_I выполняется аналогично, но не идентично, вычислению массива β'' из раздела 7.1. Вместе с тем, используя те же аргументы, что и в разделе 7.1, можно показать, что вычисление массива g_I выполняется за время порядка $\Theta(m)$. Детали этих вычислений оставим для рассмотрения в упражнении 7.2.8, а здесь сформулируем основной результат этого подраздела.

Теорема 7.2.2. Для данного паттерна $p = p[1..m]$ массив $\delta_2 = \delta_2[0..m]$ можно вычислить за время порядка $\Theta(m)$. ■

Мы отсылаем читателя, заинтересовавшегося этими сложными вычислениями, к работе [119], где впервые доказано, что массив δ_2 можно корректно вычислить за линейное время, а также к недавней работе [60], где предложено более простое доказательство этих фактов.

Примеры применения алгоритма БМ

Чтобы лучше понять, как вычисляются и используются массивы δ_1 и δ_2 в алгоритме БМ, рассмотрим пример строки Фибоначчи f_7 , представленный в предыдущем разделе.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
$f_7 = S_2$	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	
$\beta =$	0	0	1	1	2	3	2	3	4	5	6	4	5	6	7	8	9	10	11	7	8	
$\rho =$	8	7	6	5	4	3	2	1	8	7	6	5	4	3	2	1	3	2	1	0	0	
$\beta_T =$	0	0	1	2	3	1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8	
$g_I =$	0	0	0	0	0	0	0	0	0	0	0	0	8	0	0	0	0	16	0	16	20	
$g_{II} = 8$	8	8	8	8	8	8	8	8	8	8	8	8	8	8	3	3	3	3	3	1	1	0
$\delta_2 = 34$	33	32	31	30	29	28	27	26	25	24	23	22	13	25	24	23	22	5	22	5	1	

Вычисление массива δ_1 просто: $\delta_1[a] = 0$, $\delta_1[b] = 1$ и для любой буквы λ , не равной a или b , $\delta_1[\lambda] = 21$. Однако для поиска вхождения в строки Фибоначчи текстовых строк, состоящих только из букв a и b , массив δ_1 не играет никакой роли, поскольку в этом случае для всех $j \in 0..m$

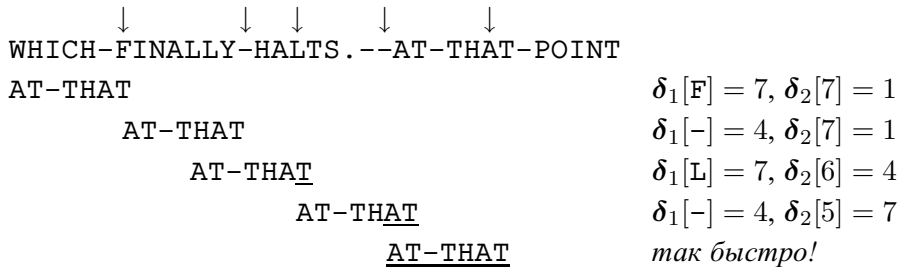
$$\max\{\delta_1[\lambda], \delta_2[j]\} = \delta_2[j]. \tag{7.9}$$

Но если паттерн будет содержать и другие буквы, отличные от букв a и b , тогда массив δ_1 необходим.

Чтобы проиллюстрировать важность обоих массивов δ_1 и δ_2 , рассмотрим пример, приведенный в исходной работе [38], где впервые описан алгоритм БМ. В этом примере паттерн $p = \text{AT-THAT}$ сравнивается с текстовой строкой

$x = \dots \text{WHICH-FINALLY-HALTS. --AT-THAT-POINT} \dots$

В этом случае $\delta_2[0..7] = 12\ 11\ 10\ 9\ 8\ 7\ 4\ 1$, тогда как $\delta_1[A] = 1$, $\delta_1[T] = 0$, $\delta_1[-] = 4$, $\delta_1[H] = 2$ и $\delta_1[\lambda] = 7$ для любых других букв λ . Использование значений массивов δ_1 и δ_2 показано на приведенной ниже схеме выполнения алгоритма БМ, где вертикальными стрелками обозначены начальные позиции буквенных сравнений для процедуры *hctam*. Отметим, что здесь для локализации паттерна потребовалось всего 15 буквенных сравнений.



Эффективность алгоритма БМ

В отличие от алгоритма КМП, алгоритм БМ в значительно большей мере отличается в теоретическом плане от простого алгоритма 2.2.1 и, как показывает предыдущий пример, на практике выполняется также значительно быстрее. Практические [212] и теоретические [25, 28] исследования показывают, что среднее количество выполняемых алгоритмом БМ буквенных сравнений является сублинейной функцией от длины текстовой строки, даже если сюда приплюсовать буквенные сравнения, выполняемые на этапе вычисления массивов δ_1 и δ_2 . Обычно для текста на естественном языке среднее количество буквенных сравнений составляет примерно $0,3n$, включая предварительную обработку паттерна, если длина паттерна больше 10, но намного меньше n . Для сравнения — алгоритм КМП в этой ситуации требует в среднем $1,0n$ буквенных сравнений (см. упражнение 7.1.6).

Чтобы среднее количество буквенных сравнений было сублинейным, для этого важную роль в алгоритме БМ играет использование массива δ_1 , особенно для текстов на английском или на других языках, имеющих большие алфавиты. Однако для анализа самого худшего случая эффектом от использования массива δ_1 можно пренебречь, поскольку этот самый “худший случай” реализуется на бинарном алфавите. Тогда, как показывает пример со строкой Фибоначчи в предыдущем подразделе и пример со строкой $x = (a^k b)^r$, рассмотренный ниже, всегда выполняется равенство (7.9). Поэтому в анализе самого худшего случая мы будем игнорировать использование массива δ_1 .

К сожалению, в самом худшем случае, независимо от того, используется или нет массив δ_1 , поведение алгоритма БМ не лучше, чем поведение простого алгоритма 2.2.1 — алгоритм БМ также может потребовать $m(n - m + 1)$ буквенных сравнений. Этот факт может подтвердить пример $x = a^n$, $p = a^m$, который успешно решается алгоритмом КМП, но в алгоритме БМ для каждого из $n - m + 1$ вхождений паттерна p требуется провести m буквенных сравнений, чтобы определить одно вхождение. Этот феномен отражает основное различие между алгоритмами КМП и БМ: если паттерн является кратной строкой, то в алгоритме КМП сравнения выполняются только на одном периоде паттерна, тогда как в алгоритме БМ сравнения производятся для каждой буквы паттерна. В разделе 8.4 будет показана модификация алгоритма БМ, которая гарантированно находит все вхождения даже периодических паттернов за линейное время. Но сейчас мы исследуем временную сложность исходного алгоритма БМ.

Анализ поведения алгоритма БМ в самом худшем случае, даже если исключить из рассмотрения периодические паттерны и проигнорировать массив δ_1 , является очень сложной задачей. В работе [137] в результате тяжелого доказательства показано, что этот алгоритм выполняет не более $7n$ буквенных сравнений для нахождения всех вхождений аperiodического паттерна. Несколько лет спустя результатом другого трудного доказательства [104] стала верхняя граница в $4n$ буквенных сравнений. И только в 1994 году эта верхняя граница смогла

опуститься до величины примерно $3n$ [53]. Доказательство последнего результата также технически сложно и трудно. Поэтому мы представим собственное доказательство облегченного результата (которое, вместе с тем, будет основываться на доказательстве из работы [53]), что для апериодических паттернов количество буквенных сравнений не превышает $4n$.

Но прежде чем приступить непосредственно к доказательству этого утверждения, рассмотрим класс примеров, для которых требуется примерно $3n$ буквенных сравнений. Рассмотрим текстовую строку $x = (a^k b)^r$, где k ($k \geq 2$) и r ($r \geq 1$) — целые числа. Применим алгоритм БМ для сравнения этой строки с паттерном $p = a^{k-1} b a^{k-1}$. Этот паттерн входит ровно $r - 1$ раз в строку x , при этом для определения каждого вхождения требуется $2k - 1$ буквенных сравнений. В упражнении 7.2.9 предложено показать, что при $r \geq 3$ алгоритм БМ выполнит ровно $(3k - 1)(r - 1)$ буквенных сравнений. Если k и r устремить к бесконечности, то отношение

$$\frac{(3k - 1)(r - 1)}{n} = \frac{(3k - 1)(r - 1)}{(k + 1)r}$$

будет стремиться к 3. Следовательно, для любого действительного числа $\varepsilon > 0$ существуют задачи сравнения с паттернами текстовых строк длиной $n = n(\varepsilon)$, для решения которых алгоритм БМ выполнит более $3n - \varepsilon$ буквенных сравнений. Эти задачи, с учетом результата статьи [53], составляют существенную часть задач для “самого плохого случая” поведения алгоритма БМ.

Теперь рассмотрим “самый плохой случай” для апериодических паттернов. Но сначала введем некоторые новые обозначения. Обозначим через t подстроку строки x , которую закончила обрабатывать процедура *hctam* на некотором шаге алгоритма БМ. Фактически $t = p[j + 1..m]$ — суффикс паттерна p длиной $t = |t| = m - j$. Заметим, что количество буквенных сравнений, выполненных при данном вызове процедуры *hctam*, равно $t + 1$ (t сравнений, при которых произошло совпадение букв, и одно сравнение, конечное, когда сравниваемые буквы не совпали). Далее обозначим через s величину сдвига паттерна, вычисляемую только с использованием массива δ_2 . Тогда или $s = (m - j') - (m - j) = j - j'$ (рис. 7.3), или $s = m - j'$ (рис. 7.4). Суффикс паттерна p длиной s обозначим как s .

Рассмотрим строку ts . При сдвиге первого типа (рис. 7.3) строка $ts = p[j' + 1..m]$ имеет грань $p[m - j]$ и, следовательно, период $s = (m - j') - (m - j)$. При сдвиге второго типа (рис. 7.4) строка ts также имеет грань длиной $m - j$, и поскольку $|ts| = (m - j') - (m - j)$, то ts имеет период $s = m - j'$.

Обозначим через f количество позиций строки x , обрабатываемых во время текущего вызова процедуры *hctam*, но такие, что буквы, стоящие в этих позициях, не встречались процедуре *hctam* ранее при ее предыдущих вызовах. Поскольку на каждом шаге алгоритма БМ паттерн сдвигается не менее, чем на одну позицию, то $f \geq 1$. Теперь покажем, что при каждом вызове процедуры *hctam* количество

$t + 1$ буквенных сравнений удовлетворяет неравенству

$$t + 1 \leq f + 3s \tag{7.10}$$

для случая, когда в нормальной форме паттерна $p = u^k u'$ показатель степени $k < 3$. (В этом случае будем говорить, что паттерн p не является 3-периодическим.) Сколько в алгоритме БМ не производилось бы сравнений, ни сумма Σf , ни сумма Σs не могут превосходить n . Тогда для случая не 3-периодического паттерна на основании неравенства (7.10) мы сможем сделать вывод, что общее количество буквенных сравнений, выполненных во время всех вызовов процедуры *hctam*, не превысит величины $n + 3n = 4n$.

Далее будем различать **длинный сдвиг** паттерна, когда $3s \geq t$, и **короткий сдвиг**, когда $3s < t$. Нетрудно заметить, что при сравнениях, приводящих к длинному сдвигу, неравенство (7.10) выполняется. Поэтому необходимо рассмотреть только те сравнения, выполняемые во время вызова процедуры *hctam*, которые приводят к короткому сдвигу паттерна.

Поскольку при сдвиге второго типа (рис. 7.4) предполагается, что паттерн p имеет период $m - j'$, короткий сдвиг требует выполнения условия

$$m - j' = s < t/3 = (m - j)/3 \leq m/3.$$

Поэтому паттерн p с необходимостью должен быть 3-периодическим, но для этого случая неравенство (7.10) мы не рассматриваем. Таким образом, для доказательства неравенства (7.10) надо исследовать только сдвиг $s = j - j'$ первого типа (рис. 7.3) при дополнительном условии, что $3(j - j') < m - j$. Такая ситуация короткого сдвига первого типа показана на рис. 7.5. Здесь мы можем убедиться, что суффикс $ts = p[j' + 1..m]$ имеет грань длиной $m - j$ и, следовательно, период $s = j - j'$.

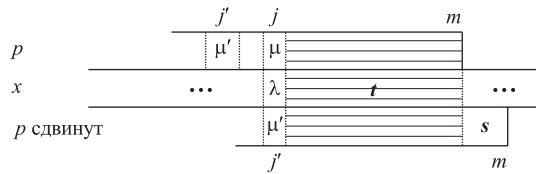


Рис. 7.5. Короткий сдвиг первого типа

Без потери общности мы можем предположить, что $s = v^k$, $k \geq 1$, где v — некратная строка. Тогда подстрока t строки x и суффикс ts паттерна p имеют период $v = |v|$. Теперь можно записать правую нормальную форму (см. упражнение 1.2.24) для строки ts : $ts = v'v^r$, где $r = \lfloor (s + t)/v \rfloor$ и v' — возможно, пустой собственный суффикс строки v , поскольку $t = v'v^{r-k}$.

Имея полную характеристику текущего сравнения подстроки t , рассмотрим сравнения, выполняемые процедурой *hctam* на предыдущих шагах алгоритма БМ.

Лемма 7.2.3. Если текущий сдвиг является сдвигом первого типа, тогда во время предыдущих выполнений процедуры *hctam*

- а) буква $p[m]$ не сравнивалась ни с одним вхождением $v[v]$ в строке t , т.е. не сравнивалась ни с одной буквой в позициях $t[t - iv]$, $i = 0, 1, \dots, r - k - 1$;
- б) между позициями паттерна p и строки t имело место не более v сравнений;
- в) буква $p[m]$ имела сравнения только с некоторыми буквами или в префиксе $t[1..v - 1]$ или в подстроке $t[t - (v - 1)..t - 1]$.

Доказательство.

- а) Предположим, что такое сравнение имело место. Поскольку при текущем вызове процедуры *hctam* $p[m]$ сравнивается с $t[t]$, самым правым вхождением $v[v]$ в t , то можно предположить, что буква $p[m]$ ранее сравнивалась с $t[t - iv]$ для некоторого целого $i \in 1..r - k - 1$ и при этом было получено совпадение

$$p[m - (t - iv) + 1..m] = t[1..t - iv],$$

а также несовпадение $p[m - (t - iv)] = \mu \neq \lambda$, где, как и ранее, λ — буква, которая в строке x предшествует подстроке t .

Зададим вопрос: какой сдвиг может быть в результате такого несовпадения? Определенно, это короткий сдвиг первого типа (рис. 7.5), поскольку достижимо значение j' , используемое для текущего сравнения. При этом можно найти такое большое значение j' , что $j' < j$. Однако отметим (см. рис. 7.5), что такой выбор значения j' ведет к такому сдвигу паттерна, что $p[m]$ должно лежать правее t . Но тогда текущие совпадения букв не будут иметь место. Получили противоречие. Делаем заключение, что буква $p[m]$ не сравнивалась ранее с любым вхождением $v[v]$ в строку t .

- б) Предположим, что $p[m]$ выравнивалась для сравнения с некоторой буквой $t[q]$, $q \in 1..t$. Как следует из утверждения а), позиция q не совпадает ни с одним вхождением $v[v]$ в строку t . Тогда может быть не более v сравнений паттерна p со строкой t , причем последнее сравнение дало несовпадение сравниваемых букв, иначе строка v равна циклическому сдвигу самой себя и, следовательно, по теореме 1.4.2 является кратной строкой, что противоречит предположению о том, что строка v не кратна.
- в) Предположим далее, что $q \notin 1..t$. Тогда q должно быть позицией в некоторой подстроке v'' строки v . Заметим, что после обнаружения несовпадения, $p[m]$ сдвигается так, чтобы провести сравнение с $v''[v]$. Тогда через конечное число выполнений процедуры *hctam* $p[m]$ будет иметь сравнение с $v[v]$, что противоречит утверждению а).

Лемма доказана. ■

Отметим, что этот результат не зависит от предположений о периодичности паттерна p — лемма справедлива для всех сдвигов первого типа. Мы используем эту лемму и предыдущее обсуждение для доказательства основного результата данного подраздела.

Теорема 7.2.4. Если паттерн p не является 3-периодическим, тогда алгоритм 7.2.1 выполняет не более $4n$ буквенных сравнений.

Доказательство. Как указывалось выше, мы можем ограничиться только рассмотрением сдвигов первого типа, где при сравнении произошло t ($t < m$) совпадений. В этом случае применима лемма 7.2.3. В частности, из утверждений δ) и ϵ) следует, что большинство позиций $[1..v - 1]$ и $[t - (2v - 2)..t - 1]$ подстроки t вычислено ранее на предыдущих этапах алгоритма — всего таких позиций $3v - 3 \leq 3s - 3$. Тогда количество f позиций, обрабатываемых впервые, удовлетворяет неравенству

$$f + 3s \geq t - (3s - 3) + 3s = t + 3.$$

Следовательно, выполняется неравенство (7.10) и поэтому количество буквенных сравнений, выполняемых в алгоритме БМ, как показано ранее, не превышает $4n$. ■

Как отмечалось выше, безусловная верхняя граница $3n$ числа буквенных сравнений устанавливается аналогично, но технически более сложным способом.

На этом мы на время прекращаем рассмотрение алгоритма Бойера–Мура. Но в следующей главе мы рассмотрим несколько вариантов этого алгоритма, включая один (раздел 8.4), для которого можно доказать линейное время выполнения даже для периодических паттернов.

Упражнения 7.2

1. Опишите структуру данных, способствующую вычислению значений $\delta_1(\lambda)$ для букв λ упорядоченного (но не индексированного) алфавита A размером α за время $O(\log \alpha)$.
2. Основываясь на предположении, что текстовая строка x (и паттерн p) начинается в нулевой позиции с символа $\$1$ (паттерн, соответственно, с символа $\$2$), напишите для алгоритма БМ процедуру *hctam* и докажите ее корректность.
3. При обсуждении массива δ_1 рассматривались случаи $j' < j$ и $j' > j$, но не случай $j' = j$. Восполните этот пробел.
4. Для заданного индексированного алфавита A и заданного паттерна p предложите алгоритм вычисления массива δ_1 . Оцените временную сложность этого алгоритма.

5. Выполните предыдущее упражнение в предположении, что алфавит не индексирован, но конечен и упорядочен, для чего используйте “пред-предварительный” этап для преобразования данного алфавита A в индексированный алфавит A' и для преобразования паттерна p в строку, определенную на алфавите A' .
6. Вычислите массив δ_1 для индексированного алфавита $A = \{a, b, c, d\}$ и паттерна $p = abcabacad$. Каково значение $\delta_1[p[m]]$ для произвольного паттерна p ? Приведите пример строки x , такой, что использование для определения сдвига только массива δ_1 всегда ведет к сдвигу паттерна p влево вдоль строки x .
7. Покажите, что “идею” массива δ_1 можно адаптировать для применения в алгоритме КМП. Опишите подробно механизм вычисления нового массива δ_1 и требования, накладываемые на этот массив для того, чтобы его можно было использовать в алгоритме КМП.
8. Приведите алгоритм вычисления за время $\Theta(m)$ массива δ_2 для заданного паттерна p . Докажите корректность этого алгоритма и оцените его временную сложность. Кроме памяти объемом $\Theta(m)$, необходимой для хранения паттерна p и массива δ_2 , сколько еще памяти может потребоваться для данного алгоритма? От каких параметров (объем алфавита, длина паттерна и т.п.) может зависеть или не зависеть объем дополнительной памяти?
9. Покажите, что для нахождения всех вхождений паттерна $p = a^{k-1}ba^{k-1}$ в строку $x = (a^k b)^r$, где k ($k \geq 2$) и r ($r \geq 3$) — целые числа, алгоритм БМ выполнит ровно $(3k - 1)(r - 1)$ буквенных сравнений. Сколько потребуется буквенных сравнений в случае $r = 1$ и $r = 2$? Предположите, что строка x и паттерн p используют начальные символы $\$1$ и $\$2$ соответственно.
10. Сколько буквенных сравнений потребуется алгоритму БМ для нахождения всех вхождений паттерна $p = (\lambda_1 \lambda_2 \dots \lambda_k)^{m/k}$ в строку $x = (\lambda_1 \lambda_2 \dots \lambda_k)^{n/k}$, если все буквы $\lambda_1, \lambda_2, \dots, \lambda_k$ различны?

7.3 Алгоритм Карпа–Рабина

Почти все из нескольких десятков разработанных после 1970 года алгоритмов сравнения с паттернами являются вариантами или алгоритма КМП или алгоритма БМ либо их обоих. Есть одно примечательное исключение из этого правила — это алгоритм Карпа–Рабина, который использует подход, совершенно отличный от подходов, использующих сдвиг паттерна вдоль просматриваемой текстовой строки.

Основная идея алгоритма КР (так для краткости мы будем называть алгоритм Карпа–Рабина, Karp–Rabin) — использование *сигнатур*. Сигнатура — это

определенное “замещающее” представление данной строки, которое должно удовлетворять следующим требованиям.

1. Должна быть малая вероятность того, что сигнатуры разных строк имеют одинаковую сигнатуру.
2. Сигнатуру можно эффективно вычислить.
3. Сигнатуру можно вычислить на основе другой сигнатуры.

Сигнатура играет такую же роль при сравнении с паттерном, что и *хеш-функция* в структурах данных [5]. Фактически, как и предложено в основополагающей работе [131], сигнатура σ строки x — это целочисленная хеш-функция, вычисленная на основе букв строки x . В алгоритме КР сигнатура $\sigma(p)$ паттерна p обычно вычисляется заранее, как и сигнатуры $\sigma(x[i..i + m - 1])$ для всех $i \in 1..n - m + 1$. Если для некоторого значения i выполняется равенство $\sigma(x[i..i + m - 1]) = \sigma(p)$, тогда вступает в действие процедура $match(i, m)$ из простого алгоритма 2.2.1 (раздел 2.2), которая проверяет, будет ли совпадение $x[i..i + m - 1] = p$.

При “хорошем” выборе функции σ , она, конечно, должна удовлетворять первому и третьему требованиям, перечисленным выше. Тогда только изредка возможно *ложное совпадение*, т.е. выполняется совпадение сигнатур, но не совпадение строк:

$$\sigma(x[i..i + m - 1]) = \sigma(p), \quad \text{но} \quad x[i..i + m - 1] \neq p.$$

Таким образом, выполнив простое сравнение двух целых чисел, для чего требуется только константное время, можно определить совпадение паттерна с подстрокой строки x . В работе [131] показано, как сигнатуры текстовых строк можно эффективно вычислять, т.е. второе требование для сигнатур также выполнимо.

Для дальнейшего обсуждения данного подхода немного упростим запись сигнатур: сигнатуру $\sigma(x[i..i + m - 1])$ будем обозначать как $\sigma(i)$. Конечно, если на каждом шаге алгоритма необходимо полностью вычислять $\sigma(i)$, то такой подход не практичен, поскольку тогда надо обеспечить доступ к каждой позиции подстроки $x[i..i + m - 1]$ и, следовательно, для вычисления каждой сигнатуры $\sigma(i)$ может потребоваться время порядка $\Theta(m)$, а всего для нахождения всех вхождений паттерна p в строку x может потребоваться время порядка $\Theta(m(n - m + 1))$.

Но в алгоритме КР предложен способ *последовательного* вычисления $\sigma(i)$, т.е. сигнатура $\sigma(i + 1)$ вычисляется на основе сигнатуры $\sigma(i)$. Сделаем упрощающее предположение, что и паттерн p и строку x можно трактовать как двоичные последовательности, состоящие из 0 и 1: при необходимости буквы любого алфавита можно закодировать в виде двоичных строк. Тогда имеет смысл рассмотреть следующие суммы:

$$\Sigma(i) = \sum_{j=1}^m 2^{m-j} x[i + j - 1]; \quad \Sigma(p) = \sum_{j=1}^m 2^{m-j} p[j]. \quad (7.11)$$

Эти суммы — просто другая запись двоичных строк длиной m в виде целых двоичных чисел такой же длины (учитывая ведущие нули). Поскольку равенство $\Sigma(i) = \Sigma(\mathbf{p})$ возможно только в том случае, когда $\mathbf{x}[i..i + m - 1] = \mathbf{p}$, алгоритм локализации паттерна \mathbf{p} в строке \mathbf{x} будет состоять только из вычисления и сравнения сумм Σ . Однако для больших паттернов (например, при $m > 32$) значения Σ также становятся слишком большими для того, чтобы их можно было бы быстро и просто сравнить. Поэтому в работе [131] предложено определять сигнатуры как вычеты по модулю q , где q — простое (достаточно большое) число:

$$\sigma(i) = \Sigma(i) \bmod q; \quad \sigma(\mathbf{p}) = \Sigma(\mathbf{p}) \bmod q. \quad (7.12)$$

Такая хеш-функция отображает 2^m всех возможных двоичных строк длиной m на множество целых чисел $0..q - 1$. Тогда каждому такому числу будет соответствовать в среднем $2^m/q$ двоичных строк или, другими словами, с вероятностью $1/q$ сигнатуры двух различных строк будут совпадать.

Чтобы показать, как вычисляются сигнатуры $\sigma(i)$, введем операции $+_q$, $-_q$ и \times_q , определяющие сложение, вычитание и умножение по модулю q соответственно. Оставим для упражнений 7.3.1 и 7.3.2 исследование свойств этих операций. Здесь отметим, что эти операции позволяют вычислить сигнатуру $\sigma(i + 1)$ на основании значения $\sigma(i)$, $i \geq 1$, за константное время по следующей формуле.

Лемма 7.3.1. $\sigma(i + 1) = 2 \times_q (\sigma(i) -_q (2^{m-1} \bmod q) \times_q \mathbf{x}[i]) +_q \mathbf{x}[i + m]$.

Доказательство этой формулы предложено дать в упражнении 7.3.3. ■

Поскольку значения $2^{m-1} \bmod q$ могут быть вычислены заранее, вычисление сигнатуры $(i + 1)$ на основании значения $\sigma(i)$ требует всего четырех операций, которые можно эффективно реализовать на компьютере. С учетом этого замечания запишем алгоритм КР.

Алгоритм 7.3.1 (Алгоритм Карпа–Рабина)

```

▷ Нахождение всех вхождений паттерна  $\mathbf{p}$  в строку  $\mathbf{x}$ 
выбрать  $q$ 
вычислить  $\sigma(\mathbf{p})$            ▷ используются формулы (7.12)
for  $i \leftarrow 1$  to  $n - m + 1$  do
    вычислить  $\sigma(i)$        ▷ используется формула леммы 7.3.1
    if  $\sigma(i) = \sigma(\mathbf{p})$  then  $j \leftarrow \text{match}(i, m)$ 
    if  $j = m + 1$  then output  $i$ 
    
```

Нет сомнения, что приведенный алгоритм КР корректен, поскольку он сводится к алгоритму 2.2.1 в случае $\sigma(i) = \sigma(\mathbf{p})$, а в случае $\sigma(i) \neq \sigma(\mathbf{p})$ по определению $\mathbf{x}[i..i + m - 1] \neq \mathbf{p}$. Поэтому имеет место следующая теорема.

Теорема 7.3.2. Алгоритм 7.3.1 корректно вычисляет все вхождения непустого паттерна p в заданную строку x . ■

Очевидно, что эффективность алгоритма КР очень сильно зависит от выбора числа q . При плохом выборе числа q процедура $match(i, m)$ может вызываться $n - m + 1$ раз (и при этом совпадения сигнатур будут ложными). Тогда для выполнения алгоритма КР потребуется время порядка $\Theta(nm)$, что не лучше, чем у простого алгоритма 2.2.1. С другой стороны, при удачном выборе числа q не будет ложных совпадений сигнатур, и поэтому алгоритм КР будет выполнен за время $O(n + km)$, где k — количество вхождений паттерна p в строку x . В работе [131] теоретически показано, что если простое число q выбрано случайным образом из интервала целых чисел $1..mn^2$ и ложные совпадения сигнатур также случайны, тогда ожидаемое время выполнения алгоритма КР практически сводится до минимума $O(n + km)$. Это ожидаемое время выполнения алгоритма КР почти сравнимо с временем выполнения алгоритма БМ. Однако заметим, что в алгоритме КР время выполнения зависит от количества k вхождений паттерна p в строку x . Поэтому если паттерн часто встречается в строке x , например $k \in \Theta(n)$, то время выполнения алгоритма имеет порядок $\Theta(nm)$.

На практике, как утверждается в статье [131], алгоритм КР будет наиболее подходящим алгоритмом сравнения с паттерном в случае больших паттернов, например, при $m \geq 200$. Если сравнивать алгоритм КР с алгоритмами КМП и БМ, то алгоритм КР требует дополнительной памяти только фиксированного объема. Для такой реализации данного алгоритма можно использовать накапливающие регистры для быстрого доступа к данным независимо от величины m , тогда как в алгоритмах КМП и БМ доступ к данным для больших значений m замедляется, поскольку данные выбираются из массива размером $\Theta(m)$. Но с другой стороны, на машинах определенной архитектуры (например, при использовании кеш-памяти) такая задержка доступа к данным может быть незначительной, а алгоритм КР может выполняться медленнее из-за интенсивного использования целочисленной арифметики [24, 103].

Как указывается в работе [131], алгоритм КР можно адаптировать для работы с двумерными паттернами, даже в случае их неравномерной формы, что ведет к эффективной реализации параллельных вычислений.

Упражнения 7.3

1. Обозначим через \bullet одну из операций $+$, $-$, \times , а через \bullet_q — соответствующую операцию $+_q$, $-_q$ или \times_q . Покажите, что для любых целых положительных чисел q , r и s выполняется равенство

$$(r \bullet s) \bmod q = (r \bmod q) \bullet_q (s \bmod q).$$

2. Для больших значений m вычисления по формулам (7.12) могут вызвать переполнение машинного слова, если суммы (7.11) вычисляются ранее, чем вычеты по модулю q . Используя операции $+_q$ и \times_q , разработайте эффективную итерационную процедуру вычисления σ , которая будет сохранять частные суммы, ограниченные величиной порядка $O(q)$.
3. Докажите лемму 7.3.1.
4. Пусть $m \leq 32$. Разработайте эффективный алгоритм, подобный алгоритму КР, который будет использовать суммы Σ вместо сигнатур σ . Какова верхняя граница выполнения такого алгоритма в самом худшем случае? Какова нижняя граница его выполнения в самом лучшем случае? Какова средняя ожидаемая его временная сложность?
5. Предположим, что p и x определены на целочисленном алфавите $\{0, 1, \dots, t\}$. Предложите наиболее приемлемые в этом случае суммы типа (7.11). Сформулируйте и докажите обобщение леммы 7.3.1.

7.4 Алгоритм Демелки–Бейза–Ятса–Гоннета

Прежде чем описывать сам алгоритм, немного остановимся на его истории. Его недавняя история началась в 1992 году, когда Бейза-Ятс и Гоннет (Baeza-Yates and Gonnet) в своей знаменитой работе [24] предложили метод, основанный на бинарных отображениях; этот метод будет описан далее в этом разделе. Эта работа инспирировала появление десятков, если не сотен других работ, которые разрабатывают самые разнообразные варианты этого метода. Однако у этой истории существует неожиданная предыстория: в работе Дёмёлки 1968 года (Dömölki, [77]) описывается практически тот же самый алгоритм, и, более того, автор ссылается на статью [78] 1965 года, в которой “этот метод описан более подробно”. Так получилось, что алгоритм Дёмёлки, опубликованный в Англии в популярном научном компьютерном журнале, остался забытым и не востребованным на 27 лет, пока его заново не открыл кто-то другой!

На этот необычный факт обратили мое внимание [190] во время моих лекций, прочитанных в университете Дебрецена в октябре-ноябре 2001 года. В честь первооткрывателя данного метода будем называть описываемый алгоритм алгоритмом Демелки–Бейза–Ятса–Гоннета (сокращенно — алгоритм ДБГ), а не алгоритмом Бейза–Ятса–Гоннета, как принято в настоящее время.

Алгоритм ДБГ подобен алгоритму КР, рассмотренному в предыдущем разделе, и также, как и алгоритм КР, отличен от алгоритмов КМП и БМ. Здесь основное внимание уделяется не наиболее быстрому способу сдвига паттерна вдоль строки x , а оптимизации вычислений на каждом шаге алгоритма, выполняющих сравнение с паттерном для каждой позиции строки x . Другими словами, алгоритм ДБГ (как и алгоритм КР) пытается минимизировать вычисления на каждом шаге,

но не пытается минимизировать количество шагов (как это делается в алгоритмах КМП и БМ). Фактически, в алгоритме ДБГ найден способ сравнения текущей буквы $x[i]$ со всеми m буквами паттерна p с помощью одной операции, которая выполняется за константное время!

Ключевую роль в алгоритме ДБГ играет бинарный массив (строка) $s = s[1..m]$, который описывает текущее *состояние* вычислений для текущей позиции i ($i \in 1..n$) в строке x [24]. В этом массиве для любого значения $j \in 1..m$ выполняется следующее соглашение:

$$s[j] = 0 \Leftrightarrow +x[i - j + 1..i] = p[1..j], \quad (7.13)$$

где i — текущая позиция в строке x . Очевидно, что если подстрока $x[i - m + 1..i]$ совпадает с паттерном p , то $s[m] = 0$. Кажется разумным до начала вычисления массива s положить все его элементы равными 1 (т.е. положить $s = 1^m$), это будет отражать тот факт, что на данном шаге алгоритма пока не производились сравнения с паттерном. Далее надо показать, как изменить значения массива s так, чтобы выполнялось соглашение (7.13).

Но сначала введем еще один массив (подобный массиву δ_1 из алгоритма БМ), предварительно сделав предположения о том, что алфавит строки x конечен, известен заранее и индексирован (см. раздел 4.1). Тогда без потери общности (см. упражнение 4.1.3) можно считать, что алфавит состоит из натуральных чисел $N_\alpha = \{1, 2, \dots, \alpha\}$. Теперь определим двухмерный массив $t = t[1..\alpha, 1..m]$, такой, что для всех $h \in 1..\alpha$ и $j \in 1..m$

$$t[h, j] = 0, \text{ если } p[j] = h, \text{ и } t[h, j] = 1 \text{ в противном случае.} \quad (7.14)$$

При таком определении массива t для любой буквы $h \in N_\alpha$ строка $t[h, 1..m]$ элементов этого массива будет состоять из нулей и единиц, причем нули будут стоять в тех позициях, в которых встречается буква h в паттерне p . Следующая лемма показывает использование массива t для вычисления массива (состояния) s в позиции i на основании известного состояния в позиции $i - 1$.

Лемма 7.4.1. Обозначим через $s^{(i)}$ строку s , определяющую состояние, соответствующее позиции i текстовой строки x . Положим, что начальное состояние описывается как $s^{(0)}[0..m] = 0(1^m)$. Тогда для всех $i \in 1..n$ и $j \in 1..m$

$$s^{(i)}[j] = s^{(i-1)}[j - 1] \vee t[x[i], j], \quad (7.15)$$

где \vee — знак логической операции OR (логическое сложение ИЛИ). (Напомним, что результат логической операции OR равен 1, если хотя бы один операнд равен 1, и 0 — в противном случае.)

Доказательство. Пусть i — текущая позиция в строке x . (Отметим, что, согласно формуле (7.15), в вычислениях не участвует значение $s^{(i-1)}[m]$.) Сначала предположим, что для некоторого $j \in 1..ms^{(i-1)}[j-1] = 0$, тогда, в соответствии с соглашением (7.13), выполняется равенство

$$x[i-j+1..i-1] = p[1..j-1].$$

Если $t[x[i], j] = 0$, тогда, согласно (7.14), $x[i] = p[j]$, и поэтому полагаем $s^{(i)}[j] \leftarrow 0$. Если же $t[x[i], j] = 1$, тогда $x[i] \neq p[j]$, и в этом случае полагаем $s^{(i)}[j] \leftarrow 1$. Каждое из этих действий совпадает с результатом выполнения операции OR.

Аналогично рассматривается случай $s^{(i-1)}[j-1] = 1$. ■

Чтобы увидеть, как работает формула (7.15) в вычислении состояний $s^{(i)}$, $i \in 1..n$, рассмотрим простой пример с паттерном $p = 212$ и строкой $x = 12112121$, определенными на алфавите $A = \{1, 2\}$. Здесь $m = 3$, $n = 8$ и $t = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$. Используя начальное состояние $s^{(0)} = 0111$, вычисляем

$$\begin{aligned} s^{(1)}[1] &= 0 \vee t[1, 1] = 1, \\ s^{(1)}[2] &= 1 \vee t[1, 2] = 1, \\ s^{(1)}[3] &= 1 \vee t[1, 3] = 1. \end{aligned}$$

Аналогично выполняются вычисления для $i = 2, 3, \dots, 8$: $s^{(2)} = 011$, $s^{(3)} = 101$, $s^{(4)} = 111$, $s^{(5)} = 011$, $s^{(6)} = 101$, $s^{(7)} = 010$, $s^{(8)} = 101$. Значение $s^{(7)}[3] = 0$ определяет совпадение p и $x[5..7]$.

Приведем листинг алгоритма ДБГ.

Алгоритм 7.4.1 (Алгоритм Демелки–Бейза–Ятса–Гоннета)

```

▷ Нахождение всех вхождений паттерна  $p$  в строку  $x$ 
 $s[0..m] \leftarrow 0(1^m)$    ▷ задание начального состояния
for  $i \leftarrow 1$  to  $n$  do
  ▷ Сдвиг на одну позицию вправо в массиве  $s$  с помощью
  ▷ процедуры rightshift, затем вычисляется новое состояние,
  ▷ оставляя без изменения  $s[0]$  и выполняя
  ▷ для остальных элементов массива  $s$  логическую
  ▷ операцию OR с операндами  $s[1..m]$  и  $t[x[i], 1..m]$ 
   $s \leftarrow \text{rightshift}(s, 1) \vee t[x[i], 1..m]$ 
  if  $s[m] = 0$  then output  $i - m + 1$ 
    
```

Алгоритм 7.4.1 основан на доказанной формуле леммы 7.4.1, поэтому мы имеем формальное право сформулировать следующее утверждение.

Теорема 7.4.2. Алгоритм 7.4.1 корректно вычисляет все вхождения непустого паттерна p в заданную строку x . ■

Чтобы оценить временную и пространственную сложность алгоритма ДБГ, необходимо ввести в рассмотрение еще один параметр, а именно, длину w машинного слова. Как видно из алгоритма 7.4.1, если $m \leq w$, то каждая из строк $s[1..m]$ и $t[h, 1..m]$ будет занимать не более одного машинного слова. Поэтому вычисление массива s может потребовать всего две операции, выполняемые за константное время: операцию сдвига на один разряд и логическую операцию OR в регистре, содержащем одно машинное слово. В случае произвольного m строки $s[1..m]$ и $t[h, 1..m]$ будут занимать по $\lceil m/w \rceil$ машинных слов, и поэтому для вычисления массива s потребуется время порядка $\Theta(\lceil m/w \rceil)$. Следовательно, имеет место такая теорема.

Теорема 7.4.3. Алгоритм 7.4.1 вычисляет все вхождения непустого паттерна p в заданную строку x за время порядка $\Theta(\lceil m/w \rceil)$, при этом используя память объемом $\Theta((\alpha + 1)\lceil m/w \rceil)$. ■

Из этой теоремы следует, что в случае больших паттернов алгоритм ДБГ выполняется за время $\Theta(nm)$, а в случае больших алфавитов требует памяти объемом $\Theta(nm)$. Поэтому в таких ситуациях данный алгоритм не применяется. Однако в случае коротких паттернов и малых алфавитов алгоритм выполняется очень быстро и эффективно. Как мы увидим в главе 10, алгоритм ДБГ можно очень легко адаптировать для вычисления приближенных паттернов — это очень полезный и гибкий алгоритм для решения задач такого класса [233].

Относительно массива $t = t[1..\alpha, 1..m]$ можно сказать, что его целесообразно вычислить на предварительном этапе, и такая операция не займет много времени.

Теорема 7.4.4. Массив t можно вычислить за время $\Theta(\alpha \lceil m/w \rceil + m)$.

Доказательство этой теоремы предложено дать в упражнении 7.4.2. ■

Напомним обсуждение 1.3.1, в котором рассматривался вопрос об отношении между размером задачи n и длиной машинного слова w и где показано, что на практике при $w \geq 32$ вполне обоснованно можно положить $\log_2 n/w \leq 2$ (или, другими словами, $w \in \Omega(\log_2 n)$). Тогда, как следует из теоремы 7.4.3, время выполнения алгоритма ДБГ составит $O(mn/\log n)$ для любых текстовых строк $x[1..n]$ — это лучшая верхняя граница времени выполнения среди всех известных алгоритмов сравнений с паттерном!

Упражнения 7.4

1. Завершите доказательство леммы 7.4.1.
2. Предложите алгоритм вычисления массива t , используемого в алгоритме ДБГ. Затем докажите теорему 7.4.4.

3. Вычислите массив t для алфавита $A = \{a, b, c, d, e\}$ и паттерна

$$p = abacabadabacabae.$$

4. Можно ли переформулировать задачу данного раздела таким образом, чтобы алгоритм ее решения использовал логическую операцию AND (логическое умножение И) вместо операции OR?

7.5 Заключение

Из четырех описанных в этой главе алгоритмов наибольшим спектром достоинств обладает, без сомнения, алгоритм Бойера–Мура. На практике, конечно, учитываются все параметры задачи (длина текстовой строки, длина паттерна, количество найденных вхождений паттерна в строку и т.д.). Алгоритм БМ при прочих равных условиях для периодических паттернов гарантирует, в теории, не более $3n$ буквенных сравнений и время выполнения порядка $\Theta(n)$. Более того, как мы видели, среднее количество буквенных сравнений в алгоритме БМ составляет только примерно одну треть от среднего количества буквенных сравнений, выполняемых алгоритмом КМП. Вместе с тем не иссякает интерес к алгоритму КМП, поскольку он гарантирует не более $2n$ буквенных сравнений в самом худшем случае и может обрабатывать множественные вхождения периодических паттернов. Если не брать во внимание периодические паттерны, то алгоритм БМ теоретически имеет единственный недостаток, который сформулирован в виде ограничений применения этого алгоритма в начале раздела 7.2. Фактически, на сегодняшний день алгоритм БМ является лучшим для точного сравнения с паттерном, даже несмотря на то что его многочисленные “потомки” подчас демонстрируют эффективность, превосходящую эффективность своего знаменитого “предка”. В следующей главе мы рассмотрим несколько таких перспективных “потомков”. В частности, мы покажем, что небольшая модификация алгоритма БМ позволяет находить множественные вхождения периодических паттернов с такой же эффективностью, как и при использовании алгоритма КМП.

ГЛАВА 8

Наследники Бойера–Мура

Тот, кто только говорит, но ничего не делает,
похож на сад, заросший сорняками.

— Аноним

В этой главе мы рассмотрим несколько вариантов алгоритма Бойера–Мура (алгоритм БМ, раздел 7.2), которые разработаны в последнюю четверть минувшего столетия и которые имеют или некоторые теоретические усовершенствования (уменьшают количество буквенных сравнений), или практические достоинства (уменьшают время выполнения). Фактически обширную литературу по этой теме можно четко разделить на теоретические и практические работы: те, которые устанавливают все более точные оценки границ числа буквенных сравнений в самом худшем случае, и те, которые предлагают все более быстрые и универсальные алгоритмы. В этой главе основное внимание уделяется именно последним работам и алгоритмам. И не потому, что я считаю теоретические работы не интересными и не важными, а просто потому, что хочу проследить эволюцию идей, ведущих к современному теоретическому пониманию проблемы сравнения текстовых строк, близкому к тому подходу, который пропагандируется в этой книге. Цель этой главы — дать “качественное” представление основных идей (многие из них являются простыми эвристиками), которые повышают практическую эффективность алгоритма БМ. И, как мы увидим, некоторые из этих идей также ведут к теоретическим улучшениям, т.е. уменьшают количество буквенных сравнений. В разделе 8.6 мы рассмотрим интересный вариант алгоритма КМП (об этом алгоритме см. раздел 7.1), который уменьшает количество буквенных сравнений

в самом худшем случае с величины $2n - m$ до величины $1,5n - 0,5m$. В последнем разделе кратко рассмотрим теоретические аспекты описанных алгоритмов.

8.1 Алгоритм БМ — цикл с перескоками

Мы начнем с алгоритма, который был представлен в исходной работе Бойера и Мура [38], но который был забыт на несколько лет. Основная идея этого алгоритма (его часто называют *циклом с перескоками* (skip loop)) заключается в том, что паттерн p сдвигается вдоль строки x путем больших “перескоков”, при этом длина перескоков определяется с помощью легко проверяемого условия, которое *позволяет надеяться*, что паттерн p не встречается в “перепрыгнутой” части строки x .

Для практической реализации данного алгоритма используется знакомый нам слегка измененный массив δ_1 (здесь его будем обозначать как δ'_1), с помощью которого получим условия перескока. Изменения в массиве δ_1 просты: для любой буквы λ положим $\delta'_1[\lambda] = \delta_1[\lambda]$, за исключением случая $\lambda = p[m]$, когда вместо присвоения $\delta_1[p[m]] \leftarrow 0$ положим $\delta'_1[p[m]] \leftarrow N$, где N — некоторое фиксированное число, такое, что $N > m + n$. Далее в алгоритме БМ (см. алгоритм 7.2.1) первый оператор в цикле **while** заменяется на следующий оператор цикла

$$\text{repeat } i \leftarrow i + \delta'_1[x[i]] \text{ until } i > n$$

Пусть i_0 — значение i , которое предшествует выполнению перескока. Если после перескока мы имеем $i > N$, то отсюда можно сделать следующие заключения:

1. $x[i - N] = p[m]$;
2. не существует позиции $i' \in i_0..(i - N) - 1$, такой, что $x[i'] = p[m]$.

Таким образом, $i - N$ — это первая позиция в строке x , которая находится правее позиции $i_0 - 1$ и в которой присутствует буква $p[m]$. В этом случае можно вернуться к обычному процессу алгоритма БМ (т.е. можем выполнить процедуру $hctam(i, m - 1)$), предварительно выполнив оператор

$$\text{if } i > N \text{ then } i \leftarrow (i - N) - 1.$$

Если же на выходе перескока $i \leq N$, тогда $i > n$, и мы снова можем утверждать, что не существует позиции i' , такой, что $i_0 \leq i' \leq n$ и $x[i'] = p[m]$. Другими словами, паттерн p не может совпасть с подстрокой $x[i - m + 1..n]$, и поэтому поиск можно закончить.

Даже это краткое описание позволит без особого труда изменить алгоритм БМ таким образом, чтобы он реализовал идею цикла с перескоками. В результате получим алгоритм, который обычно называют ВМFAST (быстрый алгоритм БМ), поскольку на практике в среднем он выполняется быстрее, чем исходный алгоритм

БМ [121]. Свойства алгоритма BMFAST и ситуации, когда он выполняется быстрее (или медленнее), чем алгоритм БМ, показаны в упражнениях 8.1.1–8.1.5.

Чтобы получить качественное представление о преимуществах последнего алгоритма, заметим, что в алгоритме BMFAST буквенное сравнение $x[i_0] : p[m]$ заменяется одним или несколькими перескоками.

- Если в действительности $x[i_0] = p[m]$, тогда в цикле с перескоками выполняется больше операторов, чем выполнялось бы буквенных сравнений. Поэтому в данном случае алгоритм BMFAST не эффективен.
- Если $x[i] = p[m]$ для некоторого $i > i_0$, тогда нахождение такого i в общем случае в алгоритме BMFAST происходит быстрее, чем в алгоритме БМ.
- Если не существует такого $i > i_0$, что $x[i] = p[m]$, то в этом случае алгоритм BMFAST также более эффективен, чем алгоритм БМ, особенно на больших алфавитах.

Предложены другие варианты циклов с перескоками, в частности в [119] предложен алгоритм, построенный на этой основе, для поиска одной буквы паттерна p в строке x , который особенно эффективен, если известно, что эта буква в строке x встречается нечасто. Этот алгоритм в общем случае выполняется быстрее, чем алгоритм БМ, но в среднем медленнее, чем алгоритм BMFAST [121].

Упражнения 8.1

1. Объясните подробно, почему имеют место утверждения 1 и 2.
2. Запишите алгоритм BMFAST и докажите его корректность в предположении, что корректен алгоритм БМ.
3. Охарактеризуйте ситуации, в которых алгоритм BMFAST будет выполняться быстрее, чем алгоритм БМ. Аналогично опишите ситуации, где этот алгоритм будет более медленным, чем алгоритм БМ. Не забудьте рассмотреть случай $x = a^n$, $p = a^m$!
4. В разделе 7.2 был показан пример $x = (a^k b)^r$, $p = a^{k-1} b a^{k-1}$, который считается “плохим” для выполнения алгоритма БМ. Оцените количество буквенных сравнений, выполняемых алгоритмом BMFAST для этого примера. Будет ли в действительности в этом случае алгоритм BMFAST быстрее, чем алгоритм БМ?
5. Всегда ли алгоритм BMFAST требует меньшего количества буквенных сравнений, чем алгоритм БМ? Обоснуйте свой ответ.
6. Запишите алгоритм, упоминавшийся в конце раздела, который реализует идею цикла с перескоками для поиска вхождения в строку одной заданной буквы.

8.2 Алгоритм Бойера–Мура–Хоспула

В предыдущем разделе показано, как небольшое изменение массива δ_1 (точнее, только одного значения $\delta_1[p[m]]$) можно использовать для реализации длинного “перескока” паттерна вдоль исследуемой строки. В этом разделе рассмотрим еще один способ вычисления значения $\delta_1[p[m]]$, который позволит реализовать алгоритм БМ без использования массива δ_2 . Это поможет избежать трудностей, связанных с вычислением массива δ_2 и описанных в разделе 7.2 (и которые снова проявят себя в разделе 8.4). Преимущества такого подхода очевидны — новый алгоритм на практике в среднем эффективнее исходного алгоритма БМ, а количество выполняемых буквенных сравнений, как показано в разделе 7.2, за исключением некоторых редких случаев, пропорционально mn .

Напомним, что путем использования массива δ_2 достигаются две цели: во-первых, алгоритм БМ никогда не возвращается назад, т.е. любое смещение паттерна выполняется только вправо; во-вторых, за исключением случая 3-периодических паттернов, алгоритм БМ выполняется за время порядка $O(n)$, т.е. достигает теоретической верхней границы времени выполнения, что невозможно сделать без использования массива δ_2 . Хоспул (Horspool, [119]) показал, что по крайней мере первую цель можно достичь, используя только массив δ_1 .

В разделе 7.2 определено, что для любой буквы λ из алфавита A $\delta_1[\lambda] = m - j'$, где j' — самая правая позиция буквы λ в паттерне p , и $\delta_1[\lambda] = 0$, если λ не встречается в паттерне. Основная идея алгоритма Бойера–Мура–Хоспула (алгоритма БМХ) заключается в том, что когда обнаружено несовпадение паттерна с очередной подстрокой строки x , паттерн p сдвигается далее не в соответствии с обнаруженной позицией несовпадения, а в соответствии с самой правой позицией в паттерне буквы, вызвавшей несовпадение. Другими словами, пусть имеется совпадение

$$x[i + 1..i + h] = p[m - h + 1..m]$$

для некоторого $h \in 0..m$, которое заканчивается несовпадением

$$\lambda = x[i] \neq p[m - h] = \mu.$$

Вместо сдвига $i \leftarrow i + \delta_1[x[i]]$, предусмотренного в алгоритме БМ, в алгоритме БМХ выполняется сдвиг

$$i \leftarrow (i + h) + \delta_1[x[i + h]].$$

Таким образом, $p[m]$ далее сравнивается с буквой в позиции i строки x , которая в текущем сравнении отстояла на $\delta_1[x[i + h]]$ позиций справа от позиции буквы $p[m]$. Такая стратегия решает одну из проблем в алгоритме БМ, порождаемую сдвигом, определяемым только на основе массива δ_1 , — теперь невозможен сдвиг паттерна *влево*, так как все значения в массиве δ_1 неотрицательны. С другой

стороны, новый подход обостряет проблему нулевого сдвига (когда $\delta_1[\lambda] = 0$), который возможен при совпадении $\lambda = p[m]$, поскольку тогда (для $h > 0$) должно выполняться совпадение $p[m] = x[i + h]$, что приводит к значению

$$\delta_1[x[i + h]] = m - m = 0.$$

Как указывалось в начале этого раздела, алгоритм БМХ решает эту проблему путем переопределения значения $\delta_1[\lambda]$ для случая $\lambda = p[m]$ (новый массив будем обозначать как δ_1''). Сначала рассмотрим ситуацию, когда совпадение $p[m] = x[i + h]$ порождено единственным вхождением λ в паттерн p . В этом случае можно не проводить сравнение p и x до тех пор, пока паттерн p не сдвинут полностью на позицию $i + h$ в строке x . Поэтому при следующем вызове процедуры *hctam* можно выполнять сравнения, начиная со сравнения $p[m]$ с буквой $x[(i + h) + m]$. В таком случае удобно положить $\delta_1''[\lambda] \leftarrow m$. В общем случае, обозначив через j' ($j' < m$) позицию самого правого вхождения буквы λ в паттерн p , имеет смысл положить $\delta_1''[\lambda] \leftarrow m - j'$. Таким способом в алгоритме БМХ создается массив δ_1'' , который совпадает с массивом δ_1 для всех букв, за исключением буквы $\lambda = p[m]$, когда значение $\delta_1''[\lambda]$ вычисляется так, как только что показано. Алгоритм БМХ легко получается из алгоритма БМ путем внесения минимальных изменений в процедуру *hctam* и в вычисление массива δ_1 . Мы оставим описание этих изменений для упражнения 8.2.1. В упражнении 8.2.3 предложено определить такие строки p и x , для которых алгоритм БМХ требует времени выполнения порядка $O(mn)$.

Относительная простота алгоритма БМХ (по сравнению с алгоритмом БМ) не только делает его привлекательным для применения на практике, но и послужила причиной многочисленных исследований его поведения. В работах [23, 25, 28] показано, что на алфавите размером α для сравнения текстовых строк длиной n и m ($n > m$) количество буквенных сравнений примерно равно (или больше) оценке $(n - m + 1)/\alpha$, которая становится все более точной при возрастании значений m и α . Так, для больших алфавитов алгоритм БМХ (несмотря на то, что в самом худшем случае он работает плохо) в среднем выполняется значительно быстрее, чем стандартный алгоритм БМ. Напомним (см. раздел 7.2), что алгоритм БМ на больших алфавитах в среднем выполняет около $0,3n$ буквенных сравнений.

Упражнения 8.2

1. В соответствии с описанием, данным в этом разделе, запишите алгоритм БМХ, указав изменения в процедуре *hctam* и в вычислении массива δ_1 .
2. Докажите корректность алгоритма БМХ, записанного в предыдущем упражнении.
3. Опишите бесконечные классы строк $p = p[1..m]$ и $x[1..n]$, таких, что алгоритм БМХ для их сравнения потребует времени порядка $\Theta(mn)$.

4. В первых двух разделах этой главы описаны различные модификации алгоритма БМ, основанные на вычисленном особым образом значении $\delta_1[p[m]]$. Можно ли скомбинировать эти модификации в единый алгоритм, который не использовал бы массив δ_2 , но выполнял бы и перескок, и сдвиг на основе измененного массива δ_1 ? Запишите такой алгоритм, докажите его корректность, оцените его асимптотическую сложность и эффективность в практическом применении.

Замечание. Автору этой книги известно, что гибридный алгоритм, предложенный в этом упражнении, никогда ранее не изучался. Ваши исследования этого алгоритма могут привести к новым открытиям! (Так уже не раз бывало, например, см. [129].)

8.3 Частота встречаемости букв и алгоритмы Бойера–Мура–Санди

Хотя все алгоритмы БМ сканируют паттерн p справа налево, нет обязательного правила именно такого порядка просмотра этих элементов. Независимо от порядка их просмотра, целью должно быть наискорейшее нахождение несовпадения паттерна с просматриваемой подстрокой строки x и кратчайшее (по количеству раз) перемещение паттерна вдоль строки x (с помощью перескоков или сдвигов). В этом разделе рассмотрим несколько вариантов алгоритма БМ, которые применимы в ситуациях, когда известны оценки частот встречаемости букв алфавита в строке x и когда значения этих частот значимо различны для разных букв. В таких ситуациях особый интерес представляют те буквы паттерна p , которые наиболее редко встречаются в строке x . Такие буквы, присутствующие в паттерне, с наименьшей вероятностью будут совпадать с буквами строки x , и на этой основе можно организовать большие перескоки или сдвиги паттерна p вдоль строки x . В частности, используя такой подход, можно значительно ускорить сравнение текстов на естественных языках, у которых частота встречаемости букв, как правило, известна.

8.3.1 Сравнение со стражем

Предположим, что на предварительном этапе алгоритма найдена позиция $j_{\text{стр}}$, такая, что буква $p[j_{\text{стр}}]$ из паттерна p — это буква, с наименьшей частотой встречаемая в текстовых строках, определенных на некотором алфавите. В этом случае процедуру *hctam* следует применять только тогда, когда найдено совпадение *стража* (guard) $p[j_{\text{стр}}]$ с соответствующей буквой в строке x . Другими словами, в этом случае можно в алгоритме БМ использовать следующий фрагмент кода.

```
if  $x[i - (m - j_{\text{стр}})] = p[j_{\text{стр}}]$  then
```

```

    (i, j) ← hctam(i, m)
    if j = 0 then output i + 1
else
    j ← jстр
    i ← i + max{δ1[x[i]], δ2[j]}

```

Доказательство корректности такого варианта алгоритма БМ мы оставим для упражнения 8.3.1 (доказательство не очевидно).

Идея “стража”, по-видимому, впервые предложена в работе [142], затем она была независимо сформулирована в работе [121]. Согласно тестовым проверкам, результаты которых представлены в [121], использование стража приводит к явному увеличению скорости сравнения текстов на естественных языках.

8.3.2 Наиболее эффективные перескоки

Из сказанного выше следует вывод, что наиболее быстрый поиск в строке x буквы $p[j]$ будет тогда, когда $j = j_{\text{стр}}$. Также представляет интерес предположение, сделанное в упражнении 8.3.5, что алгоритм со стражем можно эффективно реализовать как цикл с перескоками, где значение $j_{\text{стр}}$ будет использоваться для того, чтобы исключить проверку всех позиций в строке x во время перескока.

Здесь мы рассмотрим еще один подход к организации перескоков, описанный в [121], — это модификация алгоритма BMFAST, где “основная” буква $p[m]$ заменена на наиболее редко встречаемую букву $p[j]$.

Предположим, нам надо найти вхождения паттерна $p = \text{marzipan}$ в тексте x на английском языке, набранном в нижнем регистре и из которого, для простоты, удалены пробелы и знаки препинания. Алгоритм BMFAST выполняет перескоки длиной $\delta'_1[x[i]]$, если текущая буква $x[i]$ совпадает с $p[m] = p[8] = n$. Длина перескоков может быть от 1 (если $x[i] = a$) до 8 (если в паттерне p нет буквы $x[i]$) или до N (если $x[i] = n$). Поскольку в английском языке буква n шестая по частоте встречаемости, то перескоки длиной N будут совершаться относительно часто. Так как такие перескоки расточительны (значительно менее эффективны, чем просто распознать равенство $x[i] = p[m] = n$), то желательно уменьшить частоты таких перескоков.

Один из возможных способов уменьшить количество перескоков длиной N — это искать в строке x букву, совпадающую с буквой $p[1] = m$, но еще лучше выбрать для сравнения букву, которая расположена ближе к правому краю паттерна и которая имеет низкую частоту появления в строке x . “Кандидатами” для такого сорта букв могут быть буква $p[6] = p$ и, особенно, буква $p[4] = z$. (Мы не рассматриваем буквы $p[7] = a$ и $p[5] = i$, поскольку эти буквы имеют еще более высокую частоту встречаемости в английском языке, чем буква n .) Для организации сравнения с этими буквами можно вычислить массив δ'_1 так,

как будто паттерн p содержит только 6 (если использовать букву $p[6] = p$) или 4 буквы (при выборе буквы $p[4] = z$). Конечно, в этом случае средняя длина перескоков уменьшится (поскольку исключены перескоки длиной N), а длина наибольшего перескока составит только 6 (соответственно, 4) позиций. Таким образом, возникает необходимость оценить и сравнить “стоимость” уменьшения длины перескоков и “стоимость” уменьшения количества перескоков длиной N . В работе [121] проведен соответствующий анализ и предложена методика выбора для сравнения такой буквы паттерна p , которая минимизирует (по вероятности) количество перескоков. Для выполнения предварительных действий по выбору такой буквы требуется время порядка $\Theta(m)$, но это все равно приводит к относительно небольшому уменьшению общего времени сравнения паттерна и текстовой строки — уменьшение составляет всего 3–5% [121].

8.3.3 Первый алгоритм Бойера–Мура–Санди

Как показано в упражнении 8.3.5, использование стража в действительности скрывает перескоки. Все, что мы рассмотрели в этом разделе ранее, — это просто набор способов перескока, основанных на частоте встречаемости букв. Теперь рассмотрим более революционную идею: переопределение в алгоритме БМ массива δ_1 таким образом, чтобы полностью исключить необходимость выполнения буквенных сравнений с паттерном p в порядке справа налево. (Этот новый алгоритм назовем БМС1, т.е. первый алгоритм Бойера–Мура–Санди.)

Подобно всем другим модификациям алгоритма БМ, обобщающая идея алгоритма БМС1 также проста [217]. Предположим, что при текущем сравнении паттерна p и строки x буква $p[m]$ выровнена по букве $x[i_0]$. Тогда после обнаружения несовпадения букв следующая попытка сравнения не сможет достигнуть успеха, пока буква $x[i_0 + 1] = \lambda$ будет совпадать с такой же буквой λ , входящей в паттерн p . Следовательно, если j' — самая правая позиция буквы λ в паттерне p , то возможен “безопасный” сдвиг паттерна вдоль строки x на $m - j' + 1$ позиций. Это наблюдение подводит к мысли вычислить на предварительном этапе алгоритма БМС1 массив Δ_1 следующим образом: для каждой буквы $\lambda \in A$ положим

$$\Delta_1[\lambda] \leftarrow m - j' + 1,$$

где j' — самая правая позиция буквы λ в паттерне p , и $\Delta_1[\lambda] \leftarrow 0$, если буквы λ нет в паттерне p . Отметим, что для любой буквы λ

$$\Delta_1[\lambda] = \delta_1[\lambda] + 1.$$

Нетрудно также заметить, что если при текущем сравнении, когда буква $p[m]$ выравнивалась по букве $x[i_0]$, обнаружено несовпадение, то величина $\Delta_1[x[i_0 + 1]]$ определяет сдвиг паттерна p вправо для следующего сравнения, когда буква $p[m]$ выравнивается по букве $x[i_0 + \Delta_1[x[i_0 + 1]]]$.

Как показано в работе [S90], массив Δ_1 имеет ряд преимуществ по сравнению с массивом δ_1 .

1. Подобно элементам массива δ_1'' из алгоритма БМХ, элементы $\Delta_1[\lambda]$ всегда больше нуля (если λ входит в p), и поэтому сдвиг паттерна вдоль строки x всегда будет ненулевым.
2. Можно ожидать, что в среднем величины сдвигов, рассчитанные на основе массива Δ_1 , будут значительно большими, чем рассчитанные на основе массива δ_1 . Напомним, что в алгоритме БМ обнаружение несовпадения в позиции i ведет к сдвигу паттерна p вдоль строки x на величину

$$\delta_1[x[i]] - (m - j) = \Delta_1[x[i]] - 1 - (m - j), \quad (8.1)$$

где $m - j$ — длина совпадающих подстрок перед обнаружением несовпадения в позиции i . Поскольку в среднем

$$\Delta_1[x[i_0 + 1]] = \Delta_1[x[i]],$$

то из (8.1) следует, что сдвиг на основе массива Δ_1 будет превышать сдвиг на основе массива δ_1 на величину $m - j + 1$, где $m - j$ — средняя длина совпадающих подстрок.

3. Основное преимущество применения массива Δ_1 : поскольку теперь величина сдвига паттерна не зависит от позиции обнаруженного несовпадения (при текущем сравнении), следовательно, величина сдвига не зависит и от порядка, в котором будут сравниваться буквы паттерна. Таким образом мы освободились от зависимости алгоритма БМ от порядка выполнения сравнений справа налево!

Как видно из этого списка “достоинств”, алгоритм БМС1 можно легко реализовать так, чтобы сдвиг определялся только значениями массива Δ_1 . Кроме того, можно ожидать, что этот алгоритм будет эффективнее алгоритма БМХ, поскольку средняя величина сдвига в алгоритме БМС1 больше, чем в БМХ. Первые два пункта списка “достоинств” показывают, что в простейшей форме алгоритм БМС1 может поддерживать просмотр паттерна справа налево. Однако третий пункт дает свободу в выборе порядка просмотра паттерна. Поэтому в принципе можно найти такой порядок, который приведет к более быстрому выполнению алгоритма, хотя бы в среднем. Это заставляет опять вспомнить о частоте встречаемости букв. Рассмотрим более эффективную версию алгоритма БМС1, где используется знание о частоте букв.

В наиболее общей версии алгоритма БМС1 применяется перестановка позиций $1..m$ паттерна p , которая хранится в массиве $\pi = \pi[1..m]$. Хотя возможна любая перестановка позиций $1..m$, рационально применить перестановку, основанную на частоте встречаемости букв. Такая перестановка определяется следующим образом: для каждого $j \in 1..m$ определяем $\pi[j] = j'$, где j' — позиция буквы

$p[j']$, имеющей j -ю по величине частоту встречаемости. Здесь предполагается, что позиции в паттерне считаются справа налево. Например, для паттерна

$$\begin{array}{cccccccc}
 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\
 p & = & t & a & r & z & i & p & a & n
 \end{array}$$

получаем $\pi = 4\ 6\ 1\ 3\ 8\ 5\ 7\ 2$.

Для использования массива π в алгоритме БМС1 и в соответствующей процедуре сравнения достаточно просто $p[j]$ заменить на $p[\pi[j]]$. Эти изменения в алгоритме предлагается сделать в упражнении 8.3.7. Отметим, что теперь сравнение паттерна p и строки x можно всегда проводить в естественном порядке $j = 1, 2, \dots, m$, поскольку для выполнения сканирования паттерна слева направо или справа налево можно применить перестановки $\pi = 1\ 2 \dots n$ или $\pi = nn - 1 \dots 1$ соответственно.

Алгоритм БМС1 особенно привлекателен для сравнения паттернов со строками, определенными на алфавитах, для которых известно распределение частот встречаемости букв, и особенно если это распределение явно асимметрично. Его также можно скомбинировать со способами вычисления эффективных перескоков, подобными тому, который используется в алгоритме BMFAST. Такой гибридный алгоритм не только быстрый, но и относительно простой. Однако алгоритм БМС1 имеет и обратную сторону: возрастает время выполнения (по сравнению с алгоритмом БМ), поскольку используется не прямой доступ к элементам паттерна p посредством отображения $p[\pi[j]]$. Это дополнительное время вычислений увеличивается, когда распределение частот встречаемости букв алфавита не имеет сильно выраженной асимметрии. Кроме того, как показано в упражнении 8.3.6, требуется не менее $\Theta(n)$ дополнительного времени для предварительного вычисления массива π , что может составить значительную величину при больших значениях m .

8.3.4 Второй алгоритм Бойера–Мура–Санди

В этом подразделе мы покажем, как освободить алгоритм БМ от требования сканировать паттерн справа налево, используя для этого модификацию массива δ_2 . В результате получим второй алгоритм Бойера–Мура–Санди (назовем его для краткости БМС2), который в некоторых случаях выполняется даже быстрее, чем алгоритм БМС1.

Сначала заметим, что массив Δ_1 из алгоритма БМС1 можно скомбинировать с массивом δ_2 из стандартного алгоритма БМ. Оставляя пока в стороне вопрос о вычислении массива Δ_1 , заменим в алгоритме БМ присвоения $i_0 \leftarrow i \leftarrow m$, выполняемые в начале каждого шага сравнения, на присвоения

$$i_0 \leftarrow i \leftarrow \max \{i_0 + \Delta_1[x[i_0 + 1]], i + \delta_2[j]\}.$$

Как указывается в работе [219], это должно улучшить алгоритм в случае, когда известно распределение частот встречаемости букв алфавита. Можно ввести массив Δ_2 , который будет обобщением массива δ_2 , если использовать $\pi[j]$ вместо j . Если, как и выше, i_0 будет обозначать позицию в строке x , с которой выравнивается позиция m в паттерне p , тогда для всех $h \in 1..m\pi'$ $[h] = i_0 - (m - \pi[h])$ будет обозначать позицию в строке x , с которой выравнивается позиция $p[\pi[h]]$. Допуская некоторую небрежность в записи, кратко можно записать ($\forall h \in 0..m$)

$$p[\pi[1..h]] = p[\pi[1]]p[\pi[2]] \dots p[\pi[h]],$$

$$x[\pi'[1..h]] = x[\pi[1..h] + (i_0 - m)] = x[\pi'[1]]x[\pi'[2]] \dots x[\pi'[h]].$$

Тогда каждое частичное совпадение паттерна p и строки x можно выразить как равенство $x[\pi'[1..j - 1]] = p[\pi[1..j - 1]]$ для некоторого $j \in 0..m + 1$, при котором обнаруживается несовпадение

$$x[\pi'[1..j]] \neq p[\pi[1..j]]. \quad (8.2)$$

Так же, как и в алгоритме БМ, мы можем считать, что x и p заканчиваются специальными несовпадающими символами $\$1$ и $\$2$ соответственно. Тогда при $j = m + 1$ (в этом случае паттерн p совпадает с соответствующей подстрокой строки x) выполняются равенства

$$x[\pi'[m + 1]] = \$1, \quad p[\pi[m + 1]] = \$2.$$

Следовательно, и в таком случае справедливо условие (8.2).

Используя эти обозначения и соглашения, опишем далее вычисление массива Δ_2 как обобщение вычислений, соответствующих сдвигам первого и второго типов в алгоритме БМ.

- **Сдвиг первого типа.** Если существуют наибольшее положительное целое число $j' < \pi[j]$ и смещение $d = \pi[j] - j'$, такие, что

$$p[\pi[1..j - 1] - d] = p[\pi[1..j - 1]],$$

и при этом $p[j'] \neq p[\pi[j]]$, тогда сдвиг, порождаемый несовпадением в позиции $\pi[j]$, определяет значение d j -го элемента массива Δ_2 , т.е. $\Delta_2[j] \leftarrow d$.

- **Сдвиг второго типа.** Если не существует целого j' , такого, как описано в предыдущем пункте, тогда находится наибольшая грань $p[\pi[1..j']]$ строки $p[\pi[1..m]]$ ($0 \leq j' \leq m - 1$), удовлетворяющая следующим условиям: $j' < j$, $d = \pi[m] - \pi[j'] > 0$ и

$$d = \pi[m - 1] - \pi[j' - 1] = \pi[m - 2] - \pi[j' - 2] = \dots = \pi[m - j' + 1] - \pi[1].$$

Используя полученное таким образом значение d , снова полагаем $\Delta_2[j] \leftarrow d$.

Поскольку возможна ситуация, когда $j' = 0$ (т.е. строка $p[\pi[1..m]]$ имеет только пустую грань), необходимо положить $\pi[0] = 0$ для того, чтобы приведенные выше условия всегда имели смысл. В действительности, как показано в упражнении 8.3.9, если перестановка π соответствует порядку возрастания частот встречаемости букв алфавита, то почти всегда будет выполняться равенство $j' = 0$, т.е. в этом случае будем иметь максимальный возможный сдвиг $d = \pi[m]$.

Нетрудно заметить, что вычисление массива Δ_2 , выполненное для “переставленного” паттерна $p[\pi[1..m]]$, в точности соответствует вычислению массива δ_2 для исходного паттерна $p[1..m]$. Основная трудность в вычислении массива Δ_2 возникает из-за того, что здесь позиции $\pi[h]$ и $\pi[h+1]$ не обязаны быть смежными, и поэтому необходимо искать величину смещения d . Вследствие этих дополнительных вычислений и непрямого доступа к буквам паттерна, предварительное вычисление массива Δ_2 требует больших затрат, чем вычисление массива δ_2 , хотя порядок сложности вычисления массива Δ_2 остается на уровне $\Theta(m)$. Таким образом, несмотря на то что вычисление массива Δ_2 требует больших затрат, ничего принципиально нового в этих вычислениях нет. Поэтому подробно мы их рассматривать не будем.

Когда массив Δ_2 вычислен, сдвиг в алгоритме БМС2 можно определить по формуле

$$i_0 \leftarrow i_0 + \max\{\Delta_1[x[i_0 + 1]], \Delta_2[\pi[j]]\}$$

или по упрощенной формуле (предложенной в [S90])

$$i_0 \leftarrow i_0 + \Delta_2[\pi[j]].$$

Как показано в статье [121], даже значительно усовершенствованные версии алгоритма БМС2 конкурентоспособны только в случае длинных паттернов, когда возрастающая стоимость предварительных вычислений становится ограничительным фактором для любых алгоритмов. Эти исследования вызывают сомнения в высокой эффективности алгоритма БМС2 и приводят к разработке значительно более простого алгоритма [121], который мы назовем БМС2*.

В этом алгоритме для определения величины перескока выбирается буква $p[j]$ или как самая правая буква паттерна (тогда $j = m$, так делается в алгоритме BMFAST, см. раздел 8.1), или на основе частоты встречаемости букв (см. выше подраздел “Наиболее эффективные перескоки”). После выполнения перескока алгоритм либо завершается, если перескок вышел за пределы строки x , либо буква $p[j]$ выравнивается по букве $x[i_0 - (m - j)]$ и начинается побуквенное сравнение. (Здесь, как и ранее, считается, что на предыдущем шаге позиция m паттерна выравнивалась по позиции i_0 в строке x .) Затем, после того как процедура *hctam* найдет несовпадение, в алгоритме БМС2* выполняется сдвиг $i_0 \leftarrow i_0 + \Delta[j]$, где $\Delta[j]$ определяется следующим образом:

- если существует наибольшее целое число $j' < j$, такое, что $p[j'] = p[j]$, тогда $\Delta[j] = j - j'$;
- в противном случае $\Delta[j] = j$.

Как видно из этого описания, алгоритм БМС2* очень прост с несложным предварительным этапом. Но несмотря на свою простоту этот алгоритм оказался очень эффективным на практике [121], особенно в случае использования стража в процедуре сравнения и эффективной реализации перескоков. Детали этого алгоритма оставляем для упражнения 8.3.11.

Упражнения 8.3

1. Докажите корректность алгоритма со стражем, приведенного в подразделе “Сравнение со стражем”, основываясь на следующих фактах.
 - В алгоритме используется значение $\delta_1[x[i]]$ независимо от того, установлено равенство $x[i] = p[m]$ или нет.
 - Значение $\delta_2[j]$ используется независимо от того, установлено равенство

$$x[i - (m - j) + 1..i] = p[j + 1..m]$$

или нет.

2. Профессор Шрдлу (Etaoin Shrdlu) вычислил массив вероятностей $prob[\lambda]$, с которыми буквы λ некоторого алфавита A встречаются в произвольном месте произвольного текста, основанного на этом алфавите. Запишите алгоритм, вычисляющий $j_{стр}$ на основе массива $prob$, и определите его сложность.
3. Рассмотрите модификацию алгоритма со стражем, где, кроме замены $j \leftarrow \leftarrow j_{стр}$, используется замена $i \leftarrow i - (m - j_{стр})$. Сравните эффективность модифицированного алгоритма с исходным алгоритмом.
4. Измените алгоритм КМП так, чтобы он использовал страж.
5. Измените алгоритм БМ со стражем таким образом, чтобы сравнение с $p[j_{стр}]$ выполнялось внутри цикла, определяющего перескок. Будет ли эта версия алгоритма работать быстрее, чем алгоритм со стражем?
6. Разработайте алгоритм, который использовал бы массив $prob$ профессора Шрдлу для вычисления массива перестановок $\pi[1..m]$ (см. алгоритм БМС1). Какова временная сложность полученного алгоритма? Как следует расположить элементы в массиве $prob$, чтобы уменьшить временную сложность алгоритма?
7. Запишите алгоритм БМС1 с применением массива перестановок π и соответствующую процедуру сравнения. Какова временная сложность алгоритма БМС1 в самом худшем случае?

8. Модифицируйте алгоритм БМС1 так, чтобы он использовал перескоки, как описано в разделе 8.1. Какие трудности необходимо преодолеть в этой модификации, которые не встречались в алгоритме BMFAST?
9. Покажите, что если перестановка π букв паттерна p упорядочивает их в порядке возрастания частоты встречаемости этих букв, тогда строка $p[\pi[1..m]]$ имеет пустую грань, за исключением случая, когда $p = \lambda^m$ для некоторой буквы λ . Докажите более общий результат: данное утверждение справедливо, если какая-либо буква входит в паттерн только подряд (без перебивки другими буквами).
10. Опишите изменения в вычислениях алгоритма КМП, где вместо массива β' будет использоваться его аналог массив B' , основанный на перестановке π позиций паттерна p . Запишите новый алгоритм КМП и процедуру *match*. Как вы думаете, будет ли новый алгоритм более эффективным на практике и в теоретическом плане, чем исходный алгоритм КМП?
11. Запишите алгоритм БМС2* для случаев, когда для определения величины перескока выбирается буква $p[j]$ как самая правая буква паттерна (тогда $j = m$) и на основе частоты встречаемости букв (тогда $j = j_{\text{стр}}$). Для этих случаев запишите процедуру *hctam* и алгоритм вычисления массива Δ . Какова временная сложность алгоритма БМС2* в самом худшем случае?

8.4 Алгоритм Бойера–Мура–Гелила

В предыдущих разделах этой главы рассматривались, в основном, простые (хотя и интересные) модификации массива δ_1 из алгоритма БМ. В этом разделе рассмотрим модификации массива δ_2 .

Напомним читателю, что в разделе 7.2 возможная периодичность паттерна p мешала установить безусловную верхнюю границу в $4n$ количества буквенных сравнений, выполняемых в алгоритме БМ. Здесь мы опишем другую, очень простую модификацию алгоритма БМ, которую назовем алгоритмом Бойера–Мура–Гелила (сокращенно, алгоритм БМГ). Затем сделаем простую модификацию алгоритма БМГ и получим алгоритм, который назовем, возможно нескромно, алгоритмом Бойера–Мура–Гелила–Смита (сокращенно, алгоритм БМГС). Последний алгоритм позволяет установить верхнюю границу буквенных сравнений без всяких условий, накладываемых на паттерн p . В результате, с помощью алгоритма БМГС получим доказательство более строгого варианта теоремы 7.2.4.

Основная идея алгоритма БМГ [100] заключается в том, что при сравнении паттерна p со строкой x можно использовать наше знание о периоде паттерна для исключения лишних сравнений. Так, если при очередном сравнении паттерна со строкой обнаружены совпадения $x[i + 1..i + m] = p[1..m]$ и несовпадение $x[i] \neq p[0]$, то осуществляем сдвиг вправо на величину $\delta_2[0] = m - \beta[m]$, т.е. на

период паттерна (см. раздел 1.2). После выполнения сдвига мы знаем, что префикс $p[1..\beta[m]]$ совпадает с соответствующей подстрокой строки x , выровненной по этому префиксу, и поэтому нет необходимости проводить здесь сравнения — достаточно проверить совпадение со строкой x суффикса $p[\beta[m] + 1..m]$ длиной $m - \beta[m]$.

Вариант Гелила (Galil) алгоритма БМ чрезвычайно прост: после завершения очередного шага сравнения значение $\beta[m]$ наибольшей грани паттерна p используется для определения последних $m - \beta[m]$ позиций паттерна, которые будут проверяться при следующем вызове процедуры *hctam*. Напомним факт из раздела 7.2, что значение $\beta[m]$ можно вычислить на предварительном этапе алгоритма как побочный продукт вычисления массива δ_2 . В представленном ниже алгоритме 8.4.1 видны все изменения, которые необходимо сделать в алгоритме БМ, основное из которых будет введение новой переменной *jump* (скачок), значение которой равно длине префикса паттерна, исключаемого из сравнения при следующем вызове процедуры *hctam*. В этом алгоритме предполагается, что процедура *hctam* при завершении сравнений возвращает значение $j = 0$ независимо от того, равняется или нет переменная *jump* нулю. Соответствующую модификацию процедуры *hctam* и доказательство корректности этого алгоритма предложено дать в упражнениях 8.4.1 и 8.4.2.

Алгоритм 8.4.1 (Алгоритм Бойера–Мура–Гелила)

▷ Нахождение всех вхождений паттерна p в строку x

$i \leftarrow m; \text{jump} \leftarrow 0$

while $i \leq n$ **do**

$(i, j) \leftarrow \text{hctam}(i, \text{jump}, m)$

if $j = 0$ **then**

output $i + 1; \text{jump} \leftarrow \beta[m]$

else

$\text{jump} \leftarrow 0$

$i \leftarrow i + \max\{\delta_1[x[i]], \delta_2[j]\}$

На большинстве тестовых задачах алгоритм БМГ выполнялся медленнее, чем алгоритм БМ, поскольку в основном цикле выполняется еще обработка переменной *jump*. Достоинства этого алгоритма проявляются, только когда большой период паттерна p часто встречается в строке x . Однако алгоритм БМГ полезен в теоретическом плане: в работе [100] на его основе доказана верхняя граница в $4n$ буквенных сравнений в самом худшем случае для произвольного паттерна, независимо от того, периодический он или нет. Далее мы покажем, как небольшое изменение в этом алгоритме позволяет получить данную границу.

Переход от алгоритма БМГ к алгоритму БМГС естественен и понятен: вместо постоянного использования переменной *jump* следует использовать ее толь-

ко в случае сдвига второго типа (см. рис. 7.4). Сдвиг этого типа легко распознается благодаря тому факту, что в этом случае выполняется неравенство $\delta_2[j] \geq m$. Поэтому если после очередного шага сравнения выполняется неравенство $\delta_2[j] \geq \delta_1[x[i]]$, то значение переменной $jump$ для следующего сравнения должно устанавливаться следующим образом: если $\delta_2[j] \geq m$, тогда $jump \leftarrow j'$, в противном случае $jump \leftarrow 0$, где, как показано на рис. 7.4, j' — это длина наибольшей грани подстроки $p[j + 1..m]$. Поскольку для сдвига второго типа справедливо равенство $\delta_2[j] = (m - j) + (m - j')$, отсюда получаем формулу для вычисления значения j' :

$$j' = (m - j) - (\delta_2[j] - m).$$

Подробности реализации модифицированного алгоритма БМГ (или, другими словами, алгоритма БМГС) оставим для упражнения 8.4.3. Отметим только, что алгоритм БМГС не намного сложнее алгоритма БМГ, но имеет некоторое преимущество (уменьшает количество буквенных сравнений) в случае реализации сдвига второго типа. Фактически, в алгоритме БМГС гарантировано, что в случае сдвига второго типа на следующем шаге число буквенных сравнений будет не более, чем на предыдущем. Во всех других случаях количество буквенных сравнений будет не больше, чем в алгоритме БМГ (и, следовательно, не больше, чем в алгоритме БМ). Таким образом, можно ожидать, что на практике алгоритм БМГС в среднем будет выполняться быстрее, чем алгоритм БМГ. В теоретическом плане алгоритм БМГС приводит к доказательству теоремы 7.2.4 для случая произвольного паттерна.

Напомним, что вследствие леммы 7.2.3 можно утверждать, что для *любого* паттерна p (т.е. независимо от того, периодический он или нет) в случае сдвига первого типа (см. рис. 7.3) всегда имеем

$$t + 1 \leq f + 3s, \tag{8.3}$$

где $t = m - j$ — длина текущего совпадающего суффикса p , f — количество совпадающих позиций и s — длина сдвига. Мы также видели, что неравенство (8.3) справедливо во всех случаях, когда результаты сравнения порождают “длинные” сдвиги (тогда $3s \geq t$). Таким образом, для доказательства теоремы 7.2.4 необходимо только рассмотреть сдвиги второго типа (см. рис. 7.4) для 3-периодического паттерна p , который в правой нормальной форме можно представить как $p = v'v^k$, $k \geq 3$, где v — некратная строка. Как и в разделе 7.2 длину строки v будем обозначать как v , т.е. $v = |v|$.

После того как мы освежили свою память фактами из главы 7, можно приступить к основному результату данного раздела.

Теорема 8.4.1. Выполнение алгоритма БМГС требует не более $4n$ буквенных сравнений.

Доказательство. Как мы видели, утверждение теоремы выполняется для всех паттернов, которые не являются 3-периодическими, и в случае длинных сдвигов и сдвигов первого типа. Поэтому предположим, что имеется период $v \leq m/3$ и короткий сдвиг s второго типа, порождаемый текущим сравнением. Поскольку сдвиг короткий, то $3s < t$, и так как это сдвиг второго типа, то, при найденном несовпадении символа $p[m-t]$ с соответствующей буквой строки x , этот сдвиг должен быть не менее $m-t$, т.е. $s \geq m-t$. На основе этих двух неравенств мы заключаем, что для выполненного сравнения 3-периодического паттерна, которое порождает короткий сдвиг второго типа, справедливы неравенства

$$t \geq 3m/4 \geq 9v/4. \tag{8.4}$$

Теперь рассмотрим сравнение, которое непосредственно предшествовало текущему сравнению. Если такого предыдущего сравнения не было или оно порождает сдвиг второго типа, то выполнение алгоритма БМГС гарантирует, что в *текущем* сравнении будут проверяться только те позиции строки x , которые уже проверялись на предыдущем сравнении. Другими словами, в этом случае для текущего сравнения выполняется равенство $t = f$, так что неравенство (8.3) будет выполняться. Поэтому надо рассмотреть только предыдущее сравнение, порождающее сдвиг первого типа, как показано на рис. 8.1.

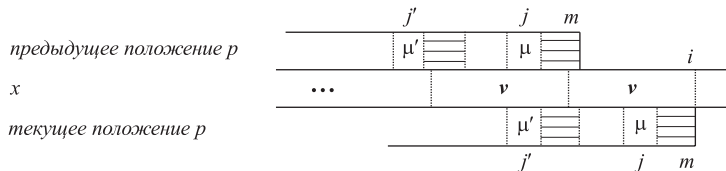


Рис. 8.1. Предыдущий сдвиг первого типа, при текущем сравнении $t \geq 9v/4$

Предположим, что на предыдущем шаге выполнено $t_{\text{пред}} = m - j$ буквенных сравнений и после этого рассчитан сдвиг $s_{\text{пред}} = j - j'$. Если, кроме того, предположить, что $t_{\text{пред}} \geq v$, тогда должна существовать подстрока

$$u = p[j' + 1..j' + (m - j)]$$

длиной не менее v , такая, что $u = p[j + 1..m] = u'v$ для некоторой строки u' . Поэтому подстрока u должна иметь суффикс v . Но поскольку буква $p[j'] = \mu'$ не совпадает с буквой μ , суффикс v не может выравняться ни с каким из k вхождений подстроки v в паттерн $p = v'v^k$. Поэтому v совпадает со своим циклическим сдвигом $R_j(v)$, $0 < j \leq n - 1$ (см. раздел 1.4). Но по теореме 1.4.2 это возможно только тогда, когда v является кратной строкой, что противоречит

определению данной строки v . На основе этого противоречия делаем заключение, что неверно предположение о том, что $t_{\text{пред}} \geq v$. Поэтому далее считаем, что $t_{\text{пред}} = |u| < v$.

Теперь сосредоточим внимание на значении $s_{\text{пред}}$. Сначала заметим, что подстрока u не может быть суффиксом строки v , поскольку она не может находиться левее несовпадения $\mu' \neq \mu$. Следовательно, строка $\mu'u$ будет или подстрокой строки v (тогда она при каждом вхождении в строку v располагается левее подстроки $p[j + 1..m]$), или подстрокой некоторого циклического сдвига строки v (тогда она также будет подстрокой строки v^2). В обоих случаях самая правая позиция $j' + (m - j)$ в строке $\mu'u$ будет находиться от позиции m в p на меньшем расстоянии, чем v . Таким образом, для предыдущего сдвига имеем

$$s_{\text{пред}} = m - (j' + (m - j)) = j - j' < v.$$

Из неравенств (8.4) следует, что подстрока v^2 строки x совпадает с суффиксом v^2 паттерна p при текущем выравнивании (позиция i на рис. 8.1). Поскольку $s_{\text{пред}} < v$, то подстрока v строки x должна быть левее крайней правой позиции паттерна при предыдущем его положении, как показано на рис. 8.1. Вследствие корректной работы алгоритма БМ (и, следовательно, алгоритма БМГС) такая ситуация невозможна (так как это спровоцировало бы сдвиг паттерна влево). Отсюда делаем вывод, что на предыдущем шаге не мог быть сдвиг первого типа. ■

Читатель может сравнить это доказательство с доказательством леммы 7.2.3 — и здесь и там используются одни и те же идеи.

Как будет показано в следующем разделе, алгоритм БМГС можно рассматривать как упрощенную версию более общего алгоритма, так называемого алгоритма Турбо-БМ, который в действительности требует всего $2n$ буквенных сравнений.

Кроме алгоритма Турбо-БМ, алгоритм БМГ инспирировал разработку по крайней мере еще одного алгоритма (назовем его БМГАГ, см. [16]). В этом алгоритме сохраняется информация о нескольких предыдущих этапах сравнения, что дает возможность сравнить паттерн p со строкой x путем использования не более $2n - m + 1$ буквенных сравнений. К сожалению, этот результат достигается алгоритмом БМГАГ только при условии выполнения дополнительных вычислений порядка $9n$ “небуквенных” сравнений и других вычислений на предварительном этапе алгоритма. С учетом этих обстоятельств данный алгоритм не может быть конкурентом алгоритму Турбо-БМ ни в практическом, ни в теоретическом планах. Однако дальнейшее совершенствование алгоритма БМГАГ, предложенное в работах [65, 60], позволило значительно сократить и упростить вычисления предварительного этапа, а также уменьшить верхнюю границу буквенных сравнений до $1,5n$.

Упражнения 8.4

1. Перепишите процедуру *hctam* для ее использования в алгоритме БМГ.
2. Докажите корректность алгоритма БМГ на основе доказанной корректности алгоритма БМ.
3. Запишите алгоритм БМГС и докажите его корректность.
4. Применяв небольшую изобретательность, можно попытаться уменьшить количество буквенных сравнений в алгоритме БМГС, выполняемых в случае, когда на предыдущем шаге имел место сдвиг первого типа. Основная трудность, которую надо преодолеть, заключается в том, что если в случае реализации на предыдущем шаге сдвига второго типа известно, что подстрока совпадающих символов является префиксом паттерна p , то в случае сдвига первого типа можно только утверждать, что эта подстрока входит в p . Здесь не обойтись только одной дополнительной переменной *jump* — надо еще хранить значение самой правой позиции подстроки совпадающих символов. Запишите этот новый алгоритм (назовем его БМГС*) и процедуру *hctam* для этого алгоритма. Очевидно, что алгоритм БМГС* никогда не будет использовать больше буквенных сравнений, чем алгоритм БМГС. Но будет ли на практике этот алгоритм более эффективным, чем алгоритм БМГС?

8.5 Алгоритм Турбо-БМ

Как сказано в предыдущем разделе, в алгоритме БМГС можно уменьшить количество буквенных сравнений (см. упражнение 8.4.4 и краткое описание алгоритма БМГС*). В этом разделе мы покажем, как можно модифицировать алгоритм БМГС путем организации хранения информации о предыдущем этапе сравнения. В этом случае не только уменьшается количество используемых буквенных сравнений, но иногда и увеличивается длина сдвига паттерна p вдоль строки x . Этот новый алгоритм назовем алгоритмом Турбо-БМ, а увеличенные сдвиги — *турбо-сдвигами* [73].

Как показано на рис. 8.2 и 8.3, после частичного совпадения подстроки $t = p[j + 1..m]$ с буквами строки x происходит соответственно сдвиг первого или второго типов, и далее наблюдается следующее совпадение подстроки t' с буквами строки x . Здесь, как и ранее, используем обозначения $t = |t|$, $s = |s|$ и $t' = |t'|$.

Случай $t' = s$ в алгоритме БМГС* рассматривается особо, поскольку тогда или суффикс ts паттерна p совпадает с буквами строки x и остается проверить только подстроку $p[1..j']$ (если предыдущий сдвиг был сдвигом первого типа), или весь паттерн p совпадает с буквами строки x (если предыдущий сдвиг был сдвигом второго типа).

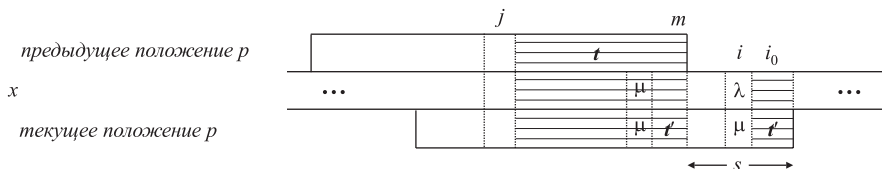


Рис. 8.2. Предыдущий сдвиг первого типа

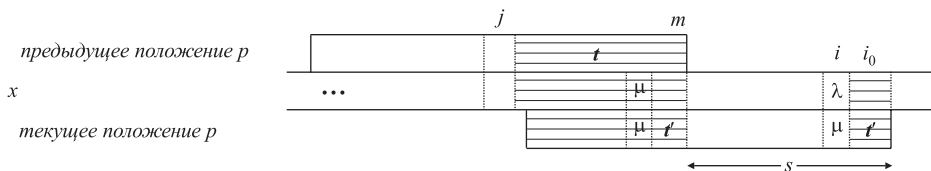


Рис. 8.3. Предыдущий сдвиг второго типа

Чтобы познакомиться с турбо-сдвигами, рассмотрим случай $t' < s$, когда зафиксировано несовпадение

$$p[m - t'] = \mu \neq \lambda = x[i_0 - t']. \tag{8.5}$$

Здесь строка $\mu t'$ является суффиксом строки s . Привлекая дополнительное условие $t' < t$, мы видим, что строка $\mu t'$ является также суффиксом строки t , совпадающей с подстрокой строки x , и, в частности, для “общей” буквы μ имеем

$$p[m - t'] = p[(m - s) - t'] = \mu = x[(i_0 - s) - t']. \tag{8.6}$$

Сравнивая (8.5) и (8.6), видим, что буквы μ и λ в строке ts отделены друг от друга на s позиций. Напомним (см. раздел 7.2), что строка ts с необходимостью имеет период s . Поэтому буквы μ и λ , стоящие в этих позициях, не могут совпасть в ts . Следовательно, можно без опаски сдвинуть паттерн вправо на $i - s$ позиций, где находится буква μ .

На рис. 8.2 видно, что в случае, когда предыдущее сравнение порождает сдвиг первого типа, букву $p[m - s - t]$ можно сдвинуть вправо на $t - t'$ позиций, тогда для следующего сравнения буква $p[m]$ выравняется с буквой $x[i + t]$. Это соответствует турбо-сдвигу

$$i \leftarrow i + t. \tag{8.7}$$

С другой стороны, на рис. 8.3 видно, что в случае, когда предыдущее сравнение порождает сдвиг второго типа, букву $p[1]$ можно сдвинуть вправо на $m - s - t'$ позиций, т.е. для следующего сравнения $p[m]$ выравняется с $x[i + m - s]$. Это соответствует турбо-сдвигу

$$i \leftarrow i + (m - s). \tag{8.8}$$

Отметим, что для сдвига первого типа $t \leq m - s$, тогда как для сдвига второго типа $m - s \leq t$. Поскольку $t = m - j$, присваивания

$$\begin{aligned} s &\leftarrow \delta_2[j] - (m - j) \\ jump &\leftarrow m - \max\{j, s\} \end{aligned}$$

для каждого типа сдвига определяют число позиций на расстоянии s от правого конца паттерна p (выровненного для следующего сравнения), не требующих проверки. Это число хранится в переменной $jump$.

Теперь пришло время записать весь алгоритм Турбо-БМ.

Алгоритм 8.5.1 (Алгоритм Турбо-Бойера–Мура)

▷ *Нахождение всех вхождений паттерна p в строку x*

$i \leftarrow m; jump \leftarrow 0$

while $i \leq n$ **do**

$(i, j) \leftarrow hctam(i, s, jump, m)$

if $j = 0$ **then output** $i + 1$

$shift \leftarrow \max\{jump, \delta_1[x[i]]\}$

$i \leftarrow i + \max\{\delta_2[j], shift\}$

if $\delta_2[j] < shift$ **then**

$jump \leftarrow 0$

else

$s \leftarrow \delta_2[j] - (m - j); jump \leftarrow m - \max\{j, s\}$

Отметим, что в данном алгоритме не требуется этап предварительных вычислений, который присутствовал в алгоритме БМ, — для хранения дополнительной нужной информации используются переменные s и $jump$. Верхняя граница в $4n$ буквенных сравнения, показанная на алгоритме БМГС*, конечно, справедлива и для алгоритма Турбо-БМ. Но фактически, как отмечалось ранее, с небольшими усилиями можно показать [73], что этот алгоритм требует не более $2n$ буквенных сравнений. Теоретически в самом худшем случае алгоритм Турбо-БМ выполняет не больше буквенных сравнений, чем алгоритм КМП (см. раздел 7.1), но в среднем, как отмечено ранее (раздел 7.5), его реализация должна быть более быстрой.

В работе [73] также предлагается модификация алгоритма Турбо-БМ, где используется ориентированный ациклический граф слов (см. раздел 5.3.1) паттерна p для определения более эффективных значений переменных s и $jump$. Такой алгоритм на практике выполняется более быстро, чем алгоритм Турбо-БМ, но он требует предварительного этапа для построения графа слов.

Упражнения 8.5

1. Перепишите процедуру $hctam$ для ее использования в алгоритме Турбо-БМ (и убедитесь, что это очень простая процедура!).

8.6 Наследники алгоритма КМП

После того как мы рассмотрели несколько вариантов алгоритма БМ, в этом разделе опишем одну модификацию алгоритма КМП (см. раздел 7.1), которая значительно (от $2n$ до $1,5n$) уменьшает количество буквенных сравнений, выполняемых в самом худшем случае в процессе решения задачи нахождения всех вхождений паттерна p в текстовую строку x . Хотя при описании этого алгоритма мы будем основываться на идеях, высказанных в работе [13], впервые он был представлен в тезисах к докторской диссертации Ханката (Hancart, [110]). Поэтому данный алгоритм будем называть алгоритмом Кнута–Морриса–Пратта–Ханката (сокращенно, КМПХ).

Основная идея алгоритма КМПХ вытекает из наблюдения, сделанного в начале раздела 7.1, где в первый раз был затронут вопрос о сложности алгоритма КМП. Напомним, что сравнение паттерна $p = a^{m-1}b$ со строкой $x = a^n$ требует $2n - m$ буквенных сравнений, т.е. столько же, сколько требуется в алгоритме КМП в самом худшем случае, как показано в разделе 7.1 позднее. Чтобы исключить такие случаи, в алгоритме КМПХ паттерн p разбивается на две части, которые затем обрабатываются более-менее самостоятельно:

- префикс максимальной длины вида $\lambda^{m'-1}$, где $1 < m' < m$, λ — некоторая буква алфавита;
- непустой суффикс $p[m'..m]$, где $p[m'] \neq \lambda$.

В случае, когда $p = \lambda^m$, полагаем $m' \leftarrow 1$, т.е. здесь префикс пустой и суффикс $p[m'..m] = \lambda^m$.

Используем обозначение в виде тройки (i, j_L, j_R) для отслеживания частичных совпадений обеих частей паттерна (левой и правой частей) с подстрокой $x[i + 1..i + m]$, как показано на рис. 8.4:

- j_L — позиция в префиксе паттерна p , такая, что $j_L \in 1..m'$ и

$$p[1..j_L - 1] = x[i + 1..i + (j_L - 1)];$$

- j_R — позиция в суффиксе паттерна p , такая, что $j_R \in m'..m + 1$ и

$$p[m'..j_R - 1] = x[i + m'..i + (j_R - 1)].$$

Начальное значение тройки (i, j_L, j_R) на момент старта алгоритма КМПХ полагается равным $(0, 1, m')$, отражая тот факт, что в этот момент не зафиксировано каких-либо совпадений с символами строки x ни префикса $p[1..m' - 1]$, ни суффикса $p[m'..m]$.

Выполнение алгоритма КМПХ подобно выполнению алгоритма ПМК: паттерн также сдвигается вдоль строки x слева направо без возможности вернуться назад, на каждом этапе сравнение паттерна p со строкой x выполняется в окне

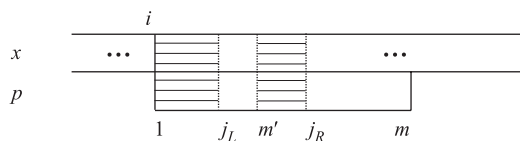


Рис. 8.4. Частичные совпадения $p[1..j_L - 1]$ и $p[m'..j_R - 1]$ левой и правой частей паттерна соответственно

длиной m , которое начинается от позиции $i + 1$ строки x . Как и в алгоритме КМП, в данном алгоритме сравнение букв происходит слева направо, но по отдельности в двух “подокнах”, определяемых значением m' . Подробнее эти вычисления можно описать как процесс изменения значений троек $t = (i, j_L, j_R)$, при этом естественно выделить три различных случая в зависимости от значения j_R .

1. $j_R = m'$ (это означает, что на текущем этапе сравнения ни одна из букв суффикса $p[m'..m]$ не совпадает с буквами строки x).

В этом случае сравниваем $p[j_R]$ с $x[i + j_R]$. Если обнаружено совпадение, то получаем новую тройку $t = (i, j_L, j_R + 1)$. Если же совпадения нет, то тройка имеет вид $t = (i + 1, j'_L, m')$, где $j'_L = \max\{1, j_L - 1\}$. Последняя тройка соответствует сдвигу паттерна на одну позицию вправо, при этом после сдвига выполняется сравнение в левом подокне с $p[1..j'_L - 1] = \lambda^{j'_L - 1}$.

2. $m' < j_R < m + 1$ (т.е. на текущем этапе сравнения подстрока $p[m'..j_R - 1]$ совпадает с буквами строки x).

Далее сравниваем $p[j_R]$ с $x[i + j_R]$. Если обнаружено совпадение, то получаем новую тройку $t = (i, j_L, j_R + 1)$ (такую же, как в п. 1).

Теперь предположим, что $p[j_R] \neq x[i + j_R]$. В таком случае следует выполнить сдвиг паттерна вправо на максимально возможное число позиций. Поскольку

$$x[i + m'] = p[m'] \neq p[j] = \lambda, \quad \forall j \in 1..m' - 1,$$

этот сдвиг не может иметь длину менее m' .

Напомним, что в алгоритме КМП далее вычисляется $\beta''[j]$, при этом или $\beta''[j] - 1$ равно длине наибольшей грани подстроки $p[1..j - 1]$, такой, что $p[\beta''[j]] \neq p[j]$, или $\beta''[j] = 0$, если такой грани не существует. Обозначим через q период подстроки $p[1..j_R - 1]$. Тогда (см. раздел 1.2) $(j_R - 1) - q$ определяет длину наибольшей грани подстроки $p[1..j_R - 1]$. Так как $j_R - 1 - q \geq \beta''[j_R] - 1$, то

$$j_R - \beta''[j_R] \geq q.$$

Но из определения m' следует, что $q \geq m'$. Поэтому

$$j_R - \beta''[j_R] \geq m'. \quad (8.9)$$

Так же как и в стандартном алгоритме КМП, сдвиг будет равен $j_R - \beta''[j_R]$ позициям. После сдвига сначала со строкой x будет сравниваться префикс паттерна длиной $\beta''[j_R] - 1$. Если $\beta''[j_R] - 1 < m'$, то $p[1..\beta''[j_R] - 1] = \lambda^{\beta''[j_R]-1}$, и новая тройка будет иметь вид

$$t = (i + (j_R - \beta''[j_R]), \max\{\beta[j_R], 1\}, m').$$

Если же $\beta''[j_R] - 1 < m'$, то новая тройка вычисляется как

$$t = (i + (j_R - \beta''[j_R]), m', \beta''[j_R]).$$

3. $j_R = m + 1$ (суффикс $p[m'..m]$ полностью совпадает с буквами строки x). Если, кроме того, $j_L = m'$, то это означает, что паттерн $p = p[1..m' - 1]p[m'..m]$ полностью совпадает с подстрокой строки x , которая начинается с позиции $i + 1$ (т.е. найдено одно из вхождений паттерна p в строку x). Если конец строки x не достигнут, то для следующего этапа сравнений тройка t вычисляется так же, как в п. 2 для случая $p[j_R] \neq x[i + j_R]$. Если $j_L < m'$, тогда выполняется сравнение $p[j_L]$ с $x[i + j_L]$. Если обнаружено совпадение этих букв, то новая тройка равна $t = (i, j_L + 1, j_R)$. В противном случае t вычисляется так же, как в п.2.

Вычисления, описанные в этих пунктах, используются в следующем алгоритме (процедура *match* здесь похожа на аналогичную процедуру в алгоритме КМП, ее оставим для упражнения 8.6.2).

Алгоритм 8.6.1 (Алгоритм Кнута–Морриса–Пратта–Ханката)

- ▷ Нахождение всех вхождений паттерна p в строку x
- $(i, j_L, j_R) \leftarrow (0, 1, m')$
- while** $i \leq n - m$ **do**
 - ▷ Поиск совпадений подстроки x с суффиксом $p[m'..m]$
 - $j_R \leftarrow \text{match}(i, j_R, m)$
 - ▷ Если $p[m'..m]$ полностью совпадает с подстрокой x ,
 - ▷ то ищутся совпадения с префиксом $p[1..m' - 1]$
 - if** $i = m + 1$ **then**
 - $j_L \leftarrow \text{match}(i, j_L, m' - 1)$
 - if** $j_L = m'$ **then output** $i + 1$
 - ▷ Переопределение троек t в соответствии с пп. 1–3
 - $i \leftarrow i + (j_R - \beta''[j_R])$
 - if** $j_R = m'$ **then**

```

     $j_L \leftarrow \max\{1, j_L - 1\}$ 
else if  $\beta''[j_R] \leq m'$  then
     $(j_L, j_R) \leftarrow (\max\{\beta''[j_R], 1\}, m')$ 
else
     $(j_L, j_R) \leftarrow (m', \beta''[j_R])$ 

```

Отметим, что в этом алгоритме переменная i на каждом шаге возрастает не менее, чем на 1, поэтому алгоритм обязательно должен завершиться. Предполагая, что тройки t пересчитываются корректно, можно утверждать, что справедлива следующая теорема.

Теорема 8.6.1. Алгоритм 8.6.1 корректно вычисляет все вхождения заданного непустого паттерна p в заданную строку x . ■

Чтобы вычислить количество буквенных сравнений, выполняемых в алгоритме КМПХ, введем величину C_i , равную количеству буквенных сравнений, которые выполняются на текущем i -м шаге цикла **while**. Если обозначить через (i, j_L, j_R) значения тройки t в начале цикла, а через (i, j_L, j_R^*) — значения после выполнения процедуры $match(i, j_R, m)$, тогда вычислить C_i можно следующим образом.

1. Если $j_R^* = m'$, тогда на текущем шаге цикла **while** имеет место только одно буквенное сравнение в процедуре $match$. Поэтому здесь $C_i = 1$.
2. Если $j_R^* \in m' + 1..m$, тогда $C_i = j_R^* - j_R + 1$.
3. Если $j_R^* = m + 1$, то это означает, что найдено полное совпадение суффикса $p[m'..m]$ с соответствующей подстрокой строки x . Поэтому процедура $match$ вызывается также для проверки префикса $p[1..m' - 1]$ как $match(i, j_L, m' - 1)$. В этом случае, если не используется специальный символ $p[m + 1] = \$$,

$$C_i \leq (m - j_R + 1) + (m' - j_L).$$

Для удобства последующих вычислений, введем еще одну величину

$$Q_i = 3i/2 + j_L/2 + j_R, \tag{8.10}$$

изменение ΔQ_i которой за время выполнения одного шага цикла **while** можно вычислить по следующим формулам.

1. Если $j_R^* = m'$, тогда, как нетрудно заметить, $i^* = i + 1$ и $j_R^* = j_R$, при этом значение j_L может уменьшиться не более, чем на 1. Поэтому

$$\Delta Q_i \geq (3 \times 1)/2 - 1/2 = 1.$$

2. В случае, когда $j_R^* > m'$ и $\beta''[j_R^*] \leq m'$, паттерн p сдвигается вправо на $j_R^* - \beta''[j_R^*]$ позиций, тогда как j_L изменится не более, чем на $\max\{\beta''[j_R^*], 1\}$, а j_R уменьшится до величины m' . Следовательно,

$$\begin{aligned} \Delta Q_i &\geq 3(j_R^* - \beta''[j_R^*])/2 + (\beta''[j_R^*] - j_L)/2 + (m' - j_R) = \\ &= (j_R^* - \beta''[j_R^*]) + (j_R^* - j_L)/2 + (m' - j_R). \end{aligned}$$

3. Наконец, если $j_R^* > m'$ и $\beta''[j_R^*] > m'$, тогда находим

$$\begin{aligned} \Delta Q_i &= 3(j_R^* - \beta''[j_R^*])/2 + (m' - j_L)/2 + (\beta''[j_R^*] - j_R) = \\ &= (j_R^* - j_R) + (j_R^* - \beta''[j_R^*])/2 + (m' - j_L)/2. \end{aligned}$$

Теперь покажем, что количество буквенных сравнений, выполняемых на любом шаге цикла **while**, не превосходит величины ΔQ_i .

Лемма 8.6.2. В цикле **while** алгоритма 8.6.1 для любого i выполняется неравенство $C_i \leq \Delta Q_i$.

Доказательство. Рассмотрим пять случаев, составленных из трех возможных вариантов вычисления величины C_i и трех вариантов вычисления величины ΔQ_i . Обозначим эти случаи как (11), (22), (23), (32) и (33), где первая цифра обозначает вариант вычисления величины C_i , а вторая — величины ΔQ_i .

(11) В этом случае мы немедленно получаем $\Delta Q_i \geq 1 = C_i$.

(22) В этом случае

$$\Delta Q_i - C_i \geq (m' - \beta''[j_R^*] - 1) + (j_R^* - j_L)/2 \geq -1/2,$$

поскольку $j_R^* > m' \geq j_L$ и $\beta''[j_R^*] \leq m'$. Но так как разность $\Delta Q_i - C_i$ должна быть целым числом, отсюда получаем $\Delta Q_i - C_i \geq 0$.

(23) Здесь, поскольку $j_R^* - \beta''[j_R^*] \geq 1$ и $m' - j_L \geq 0$, имеем

$$\Delta Q_i - C_i \geq (j_R^* - \beta''[j_R^*])/2 + (m' - j_L)/2 - 1 \geq -1/2.$$

Поэтому, как и в случае (22), $\Delta Q_i - C_i \geq 0$.

(32) В этом случае $j_R^* = m + 1$, поэтому

$$\Delta Q_i - C_i \geq (m+1)/2 - \beta''[m+1] = ((m+1) - \beta''[m+1])/2 - \beta''[m+1]/2.$$

Далее используем неравенство (8.9), которое справедливо и в случае $j_R^* = m + 1$, и напомним, что в варианте 2 вычисления ΔQ_i выполняется неравенство $\beta''[m+1] \leq m'$. Поэтому $\Delta Q_i - C_i \geq (m' - m')/2 = 0$, что и требовалось доказать.

(33) В данном случае имеем

$$\Delta Q_i - C_i \geq ((m + 1) - \beta''[m + 1])/2 - (m' - j_L)/2.$$

Снова применяя неравенство (8.9), получаем $\Delta Q_i - C_i \geq (m' - m' + j_L)/2 \geq 0$.

Лемма доказана. ■

Теперь можем сформулировать и доказать основной результат данного раздела.

Теорема 8.6.3. Алгоритм 8.6.1 вычисляет все вхождения непустого паттерна p в заданную строку x за время порядка $O(n - m)$, при этом требуется $\Theta(m)$ дополнительной памяти и выполняется не более $\lfloor (3n - m)/2 \rfloor$ сравнений с элементами строки x .

Доказательство. Очевидно, что, просуммировав по всем шагам цикла **while**, получим значение $C = \sum_i C_i$ общего количества буквенных сравнений, выполняемых в алгоритме 8.6.1. Из леммы 8.6.2 следует, что $C \leq \sum_i \Delta Q_i = Q_{\text{кон}} - Q_{\text{нач}}$, где $Q_{\text{кон}}$ — значение Q в конце последнего шага цикла **while**, $Q_{\text{нач}}$ — значение Q в начале первого шага цикла. Поскольку начальная тройка имеет вид $t = (0, 1, m')$, то из определения (8.10) следует, что $Q_{\text{нач}} = 1/2 + m'$. Значение $Q_{\text{кон}}$ удовлетворяет неравенству

$$Q_{\text{кон}} \leq 3(n - m)/2 + m'/2 + m.$$

Отсюда получаем требуемую оценку

$$C \leq (3n - m)/2 - (m' + 1)/2. \quad \blacksquare$$

Упражнения 8.6

1. Запишите алгоритм с временем выполнения порядка $O(m)$, вычисляющий m' для паттерна $p[1..m]$.
2. Перепишите процедуру *match* для ее использования в алгоритме КМПХ так, чтобы минимизировать количество используемых буквенных сравнений (т.е. вопросы эффективности кода должны стоять на втором плане после минимизации числа буквенных сравнений).

8.7 Создайте собственный алгоритм!

Даже не очень внимательный читатель мог заметить, что наше обсуждение возможностей “наследников” алгоритмов БМ и КМП является далеко не исчерпывающим, даже если принять во внимание только вопросы, относящиеся к их

использованию на практике. Вне поля нашего внимания остались, например, такие вопросы. Как изменятся алгоритмы Турбо-БМ и КМПХ, если в них ввести перескоки? Как сравнить алгоритм Турбо-БМ с алгоритмом БМХ? Как сравнить алгоритмы, основанные на частоте встречаемости букв, с алгоритмами, где используются стратегии поиска, основанные на значениях массива δ_2 ? Быстрее ли в среднем выполняется алгоритм Турбо-БМ по сравнению с алгоритмом БМГС*? Если ответ на последний вопрос утвердительный (как предполагается в работе [121]), то означает ли это, что алгоритмы с одной (но хорошей) стратегией определения сдвига (например, только на основе массива Δ_1) предпочтительнее алгоритмов, которые используют несколько стратегий определения сдвига? Как влияют на эффективность алгоритмов их программная реализация и выбор аппаратной и программной платформ?

На самом деле при изучении алгоритмов точного сравнения паттернов со строками можно сформулировать десятки подобных вопросов, связанных с эффективностью их выполнения. К счастью, часть сизифова труда по получению ответов на такие вопросы взяли на себя авторы работы [121], предложив классификацию алгоритмов точного сравнения паттернов со строками, в рамках которой можно представить и сравнить многие, в том числе наиболее важные алгоритмы этой категории. Позднее большой экспериментальный материал был опубликован в работах [153, 154, 177], а также представлен на Web-узле [47]. В этом разделе мы покажем классификацию алгоритмов, предложенную в [121]. Согласно этой классификации, алгоритмы локализации паттернов p в строках x можно достаточно точно описать на основании трех основных параметров:

- тип перескока;
- порядок, в котором выполняется сравнение букв в паттерне;
- способ определения сдвига.

В табл. 8.1–8.3 представлены обозначения и соответствующие описания значений этих параметров, на основе которых можно определить многие (но не все) алгоритмы, описанные в этой и предыдущей главах.

Таблица 8.1. Типы перескоков

Тип	Описание
<i>нет</i>	
δ'_1	как в алгоритме BMFAST
<i>редкий</i>	перескок к следующей редко встречающейся букве [119]
<i>мин.ст</i>	перескок к букве с минимальной стоимостью [121]

Таблица 8.2. Порядок сравнения

Порядок	Описание
<i>впр</i>	слева направо, как в алгоритме КМП
<i>обр</i>	справа налево, как в алгоритме БМ
π	использование массива перестановок, основанного на частоте встречаемости букв
<i>впр + стр</i>	<i>впр</i> совместно со стражем (см. раздел 8.3)
<i>обр + стр</i>	<i>обр</i> совместно со стражем
$\pi + стр$	массив перестановок совместно со стражем

Таблица 8.3. Способы определения сдвигов

Способ	Описание
+1	постоянный сдвиг паттерна на одну позицию вправо
β'	сдвиг определяет грань подстроки совпадающих букв (как в алгоритме КМП)
β''	для вычисления сдвига используется грань и определенная позиция (как в алгоритме КМП)
$\delta_1\delta_2$	сдвиг, как в исходном алгоритме БМ
δ_2	сдвиг, как в алгоритме БМ, но только на основе δ_2
δ_1''	сдвиг на основе модифицированного массива δ_1 (алгоритм БМХ)
Δ_1	сдвиг на основе модифицированного массива δ_1 (подход Санди)
Δ_2	сдвиг на основе модифицированного массива δ_2 с использованием массива π (подход Санди)
Δ	сдвиг на основе перескока (подход Санди)
<i>турбо</i>	сдвиг на основе массивов δ_1 , δ_2 и турбо-сдвига

Используя эти обозначения, можно охарактеризовать алгоритмы сравнения с паттерном в терминах троек

$\langle \text{перескок, порядок сравнения, сдвиг} \rangle$,

где каждый параметр тройки принимает одно из значений, приведенных в соответствующих таблицах. Например, простой алгоритм из раздела 2.2 можно описать в виде тройки $\langle \text{нет, впр, +1} \rangle$, а одна из характеристик алгоритма BMFAST имеет вид $\langle \delta_1', \text{обр, } \delta_1\delta_2 \rangle$, тогда как алгоритм БМС2* можно представить как $\langle \text{наим.ст, обр, } \Delta \rangle$.

Отметим, что общее количество различных троек (следовательно, и количество различных алгоритмов), которые можно составить на основе значений пара-

метров, данных в этих трех таблицах, составляет $240 = 4 \times 6 \times 10$. Даже с учетом того, что некоторые тройки невозможны (например, $\langle \text{нет}, \text{обр}, \beta' \rangle$), эта классификация предлагает к рассмотрению и изучению огромное количество возможных алгоритмов. На этой основе вы можете попытаться “сконструировать” собственный алгоритм. Успехов вам!

Упражнения 8.7

1. На основе приведенной классификации разработайте собственный алгоритм и определите его эффективность в среднем и в самом худшем случае.
2. Покажите, что алгоритм КМПХ, описанный в предыдущем разделе, не охватывается приведенной классификацией. Как можно включить его в эту классификацию?

8.8 Точные оценки сложности алгоритмов точного сравнения с паттерном

В последних двух главах мы изучали, в основном с практической точки зрения, алгоритмы вычисления всех вхождений заданного паттерна $p = p[1..m]$ в заданную текстовую строку $x = x[1..n]$. Основываясь на результатах экспериментальных исследований, было показано, что все варианты алгоритма Бойера–Мура в среднем ведут себя “почти линейно” — среднее количество буквенных сравнений, используемых этими алгоритмами, составляет cn , где c , как правило, равно 0,3 или меньше, а в алгоритме БМХ равно $1/\alpha$, где α — размер алфавита. Более того, мы видели, что по крайней мере один из вариантов алгоритма БМ (а именно, алгоритм Турбо-БМ) требует не более $2n$ буквенных сравнений в самом худшем случае, тогда как другой алгоритм (алгоритм БМГАГ, [65]) — не более $1,5n$ буквенных сравнений. Но в то же время вариант алгоритма КМП (алгоритм КМПХ) требует только $(3n - m)/2$ буквенных сравнений. Таким образом, одинаковую временную сложность в самом худшем случае, по крайней мере измеренную как количество буквенных сравнений, имеют как варианты алгоритма БМ, так и варианты алгоритма КМП. Из этих наблюдений естественным образом возникают следующие фундаментальные теоретические вопросы.

1. Какое минимальное число буквенных сравнений требуется в самом худшем случае алгоритмам сравнения с паттерном?
2. Существуют ли алгоритмы, на которых достигается эта минимальная граница?
3. Какое минимальное число буквенных сравнений требуется в среднем алгоритмам сравнения с паттерном?

Первоначально ответ на первый вопрос был дан в тот же год, когда были опубликованы алгоритмы КМП и БМ: как показано в работе [196], для любого фиксированного паттерна $p = p[1..m]$ и любого алгоритма, который корректно определяет, является ли паттерн p подстрокой произвольных текстовых строк $x = x[1..n]$, существует такая строка x^* , в которой необходимо проверить не менее $n - m + 1$ позиций. Однако не существует алгоритма, который гарантировал бы, что ему требуется просмотреть *не более* $n - m + 1$ позиций. Поэтому остается открытым вопрос, насколько точна эта нижняя граница. Далее рассмотрим вопрос о нижней границе временной сложности алгоритмов, которые только определяют, входит ли заданный паттерн p в строку x , но не обязательно находят все вхождения паттерна в строку.

В последнюю четверть минувшего столетия и позднее основные усилия исследователей были направлены на поиск ответов на вопросы 1 и 2, сформулированные выше. Значительный вклад в этот поиск был сделан в работах [54, 55, 96, 97, 52, 56]. Вариантами ответов на эти вопросы является поиск оптимального соотношения между временной и пространственной сложностью алгоритмов [99, 66]. В частности, в статье [52] представлен алгоритм (назовем его алгоритмом СН), который вычисляет все вхождения паттерна p в строку x с проверкой не более

$$\left(1 + \frac{8}{3m + 3}\right)n \tag{8.11}$$

букв строки x . Более того, этот алгоритм является *онлайновым*, хотя и в особом понимании этого термина, — буквы текстовой строки могут появляться только в окне длиной m , которое монотонно перемещается слева направо без возврата назад. Если эффективность алгоритмов измерять как количество используемых буквенных сравнений, то алгоритм СН на сегодняшний день наиболее эффективен. Он имеет только один недостаток: необходимость этапа предварительных вычислений с временем выполнения порядка $O(m^2)$.

Позднее в работе [50] в качестве более точного ответа на первый вопрос было показано, что временная сложность алгоритма СН очень близка к теоретической нижней границе числа буквенных сравнений. Для онлайн-алгоритмов, в которых применяются два вида сравнения (паттерн-текст и текст-текст), нижняя граница составляет

$$\left(1 + \frac{9}{4m + 4}\right)n, \tag{8.12}$$

причем эта оценка справедлива только для тех паттернов, длина которых удовлетворяет условию $m = 36k + 35$, $k \geq 0$. Для онлайн-алгоритмов, использующих только сравнения типа паттерн-текст, доказана оценка нижней границы, слегка превышающая приведенную. Таким образом, все онлайн-алгоритмы

должны выполнять количество буквенных сравнений, не меньшее, чем указано в (8.12).

Для общих алгоритмов, использующих сравнения как типа паттерн-текст, так и типа текст-текст, доказана оценка нижней границы в самом худшем случае

$$\left(1 + \frac{2}{m+1}\right)n, \quad (8.13)$$

которая справедлива для всех паттернов длиной $m = 2k + 1$, $k \geq 2$.

На основе приведенных оценок делаем вывод, что асимптотически (т.е. для достаточно больших m и n) так называемая “точная оценка сложности алгоритмов сравнения с паттерном”, по крайней мере в терминах буквенных сравнений в самом худшем случае, установлена практически точно. Из формул (8.11)–(8.13) следует, что наилучший алгоритм сравнения с паттерном должен выполнять

$$(1 + d/m) \times n \quad (8.14)$$

буквенных сравнений в самом худшем случае, где для онлайн-алгоритмов $9/4 \leq d \leq 8/3$ и для произвольных алгоритмов $2 \leq d \leq 8/3$.

Теперь вернемся к “забытому” третьему вопросу. Поскольку анализ поведения алгоритмов в среднем кажется более сложным, чем анализ поведения в самом худшем случае, может показаться неожиданным, что в настоящее время известен точный ответ на этот вопрос. Сначала в работе [137] был описан простой алгоритм сравнения с паттерном, который в среднем требовал

$$O\left(\frac{\log_{\alpha} m}{m}n\right) \quad (8.15)$$

буквенных сравнений, где $\alpha = |A|$ — размер алфавита и среднее берется по всем возможным строкам, построенным на алфавите A . Два года спустя в статье [234] было показано, что для достаточно больших m и $n > 2m$ верхняя граница (8.15) также является нижней границей минимального количества буквенных сравнений, которые используются в среднем любым алгоритмом сравнения с паттерном. В [152] описан вариант алгоритма БМ, на котором достигается эта граница.

Сравнивая оценку (8.15) с оценкой (8.14), нетрудно заметить, что

$$\frac{\log_{\alpha} m}{m} < \frac{m+d}{m}.$$

Отсюда вытекает, что поведение в среднем наилучшего возможного (в среднем) алгоритма асимптотически лучше (примерно на величину $\log_{\alpha} m/m$) поведения в самом худшем случае наилучшего возможного (в самом худшем случае) алгоритма. Впрочем, этот результат был интуитивно ожидаем.

Упражнения 8.8

1. Объясните, почему оценки (8.12) и (8.13) справедливы только для определенных значений m и почему эти оценки нельзя использовать как нижние границы для всех паттернов.
2. Выведите оценку (8.14) из оценок (8.11)–(8.13).

ГЛАВА 9

Алгоритмы вычисления расстояния между строками

Человек думает, что чем тверже он произносит слова, тем тверже его позиция.

— Герман Мелвилл (1819–1891)

В двух предыдущих главах мы рассмотрели алгоритмы точного сравнения с паттернами. Теперь приступим к изучению алгоритмов *приближенного* сравнения заданного паттерна p с заданной строкой $x = x[1..n]$. Однако к этой теме мы подойдем непрямым путем. Чтобы найти приближенное совпадение паттерна p с подстрокой u строки x , сначала надо определиться с понятием “приближенное совпадение”. Как показано в разделе 2.2, существует, по крайней мере, три подхода к определению этого понятия.

- Приближенность совпадения между паттерном p и подстрокой u измеряется с помощью “расстояния” $d(p, u)$. В разделе 2.2 определено несколько разновидностей такого расстояния: расстояние Хемминга $d = d_H$, расстояние преобразования $d = d_E$, расстояние Левенштейна $d = d_L$ и взвешенное расстояние $d = d_W$. Конечно, возможны и другие типы расстояния (некоторые из них упоминались в разделе 2.2).
- Паттерн может содержать метасимволы, которые представляют не некоторые конкретные буквы, а определенные классы букв. Таким образом паттерн сам становится приближенным. В разделе 2.2 описаны такие метасимволы,

как символ подстановки \bullet , представляющий произвольную одиночную букву, и символ подстановки $*$, представляющий произвольную строку букв.

- Еще более “приближенными” будут паттерны, являющиеся регулярными выражениями, которые содержат метасимволы $|$ и $*$. Эти метасимволы позволяют проводить сравнение с альтернативными подстроками и кратными подстроками соответственно (см. раздел 2.2).

Исходя из первого подхода к определению “приближенности”, для нахождения совпадений необходимо уметь вычислять расстояния $d(p, u)$ для произвольной подстроки u строки x . В этом случае задача вычисления расстояний между строками становится естественной подзадачей общей проблемы приближенного сравнения с паттерном.

Поэтому в данной главе в качестве прелюдии к алгоритмам приближенного сравнения с паттерном изучим алгоритмы вычисления расстояния $d(x_1, x_2)$ между заданными строками $x_1 = x_1[1..n_1]$ и $x_2 = x_2[1..n_2]$. Если x_1 и x_2 являются последовательностями ДНК или РНК, тогда, как мы увидим далее, вычислив расстояние $d(x_1, x_2)$, можно также определить то, что называется *оптимальным выравниванием* (optimal alignment) двух строк. Мы также обнаружим, что некоторые описываемые в этой главе алгоритмы выполняют двойные роли: они вычисляют расстояния между строками и *одновременно* выполняют приближенное сравнение с паттерном.

Как показано в работе [120], для расстояния, определяемого произвольной матрицей весов W , нижняя граница временной сложности задачи вычисления расстояний между строками равна $\Omega(n \log n)$. Эта же граница применима к алгоритмам приближенного сравнения с паттерном (раздел 13.3). В данной и следующей главах приведем несколько алгоритмов, на которых достигается эта нижняя граница, однако при этом будут использованы некоторые специальные свойства расстояний Хемминга и преобразования.

Далее в этой главе детально рассмотрим четыре алгоритма, два из которых универсальны и непосредственно вычисляют расстояния $d(x_1, x_2)$ для произвольных расстояний d . Два других алгоритма вычисляют (не обязательно единственную) наибольшую общую подпоследовательность $LCS(x_1, x_2)$, на основании которой, как показано в разделе 2.2, можно вычислить расстояние Левенштейна. Все эти алгоритмы существенно используют методы динамического программирования, что открывает широкий простор для различных теоретических и практических исследований.

В следующей главе мы покажем, как универсальные алгоритмы вычисления расстояний, описанные ниже, можно преобразовать в алгоритмы приближенного сравнения с паттерном.

9.1 Базисная рекурсия

Как предложено показать в упражнении 9.1.1, существует алгоритм прямого вычисления расстояния Хемминга $d_H(\mathbf{x}_1, \mathbf{x}_2)$ между двумя строками \mathbf{x}_1 и \mathbf{x}_2 одинаковой длины n за время порядка $\Theta(n)$. Поэтому в данном разделе исследуем только расстояния Левенштейна, преобразования и взвешенное. Поскольку алгоритмы, описываемые далее, применимы ко всем этим трем типам расстояний, то естественно обозначать расстояние просто как d , предполагая при этом, что это расстояние удовлетворяет условиям метрики (2.2)–(2.5). Из этих условий можно получить некоторые свойства операций редактирования, выполняемых на одиночных (не пустых) буквах λ и μ .

- Вставка (замена пустой строки ε на букву λ): $d(\varepsilon, \lambda) > 0$.
- Удаление (замена буквы λ на пустую строку ε): $d(\lambda, \varepsilon) > 0$.
- Подстановка (замена буквы λ на букву μ): $d(\lambda, \mu) > 0$ тогда и только тогда, когда $\lambda \neq \mu$.

Также напомним, что если $d = d_L$ или d_E , то $d(\varepsilon, \lambda) = d(\lambda, \varepsilon) = 1$ для любой буквы λ ; тогда как для $d = d_L d(\lambda, \mu) = 2$ и для $d = d_E d(\lambda, \mu) = 1$.

Отметим также, что расстояние между отдельными буквами можно рассматривать как *стоимость* выполнения соответствующих операций редактирования. В общем случае расстояние между двумя строками является минимальной стоимостью преобразования одной строки в другую.

Все алгоритмы, представленные в этой главе, являются алгоритмами *динамического программирования*. Это означает, что они сравнивают входные строки $\mathbf{x}_1 = \mathbf{x}_1[1..n_1]$ и $\mathbf{x}_2 = \mathbf{x}_2[1..n_2]$ слева направо, вычисляя для каждой пары позиций i и j минимальную стоимость преобразования подстроки $\mathbf{x}_1[1..j]$ в подстроку $\mathbf{x}_2[1..j]$, основываясь при этом на уже вычисленных стоимостях (расстояниях) $d(\mathbf{x}_1[1..i], \mathbf{x}_2[1..j - 1])$, $d(\mathbf{x}_1[1..i - 1], \mathbf{x}_2[1..j])$ и $d(\mathbf{x}_1[1..j - 1], \mathbf{x}_2[1..j - 1])$. Эта базисная рекурсия объясняется ниже.

Для дальнейших вычислений удобно ввести двухмерный *массив стоимостей* $\mathbf{c} = \mathbf{c}[0..n_1, 0..n_2]$, в котором для всех $i \in 0..n_1$ и $j \in 0..n_2$ $\mathbf{c}[i, j] = d(\mathbf{x}_1[1..i], \mathbf{x}_2[1..j])$. Начальные значения для этого массива определяются следующим образом.

- $\mathbf{c}[0, 0] = 0$ — это минимальная стоимость преобразования пустой строки в пустую строку.
- $\mathbf{c}[0, j] = \sum_{1 \leq h \leq j} d(\varepsilon, \mathbf{x}_2[h])$ — минимальная стоимость вставки первых j букв строки \mathbf{x}_2 в пустую строку.
- $\mathbf{c}[i, 0] = \sum_{1 \leq h \leq i} d(\mathbf{x}_1[h], \varepsilon)$ — минимальная стоимость удаления первых i букв строки \mathbf{x}_1 таким образом, чтобы сформировать пустую строку.

Для большего понимания массива стоимостей рассмотрим пример массива $\mathbf{c}[0..5, 0..6]$, соответствующего вычислению расстояния Левенштейна $d_L(\mathbf{x}_1, \mathbf{x}_2)$,

где $x_1 = rests$, $x_2 = stress$.¹ Как видно из табл. 9.1, искомое расстояние равно

$$c[5, 6] = d_L(x_1[1..5], x_2[1..6]) = 3.$$

Таблица 9.1. Массив стоимостей для $d_L(rests, stress)$

j	0	1	2	3	4	5	6	
i	ε	s	t	r	e	s	s	
0	ε	0	1	2	3	4	5	6
1	r	1	2	3	2	3	4	5
2	e	2	3	4	3	2	3	4
3	s	3	2	3	4	3	2	3
4	t	4	3	2	3	4	3	4
5	s	5	4	3	4	5	4	3

Отметим, что в этом примере минимальная стоимость преобразования $rest \rightarrow stre$, т.е. $c[4, 4]$, равна 4. Эта стоимость получается путем прибавления стоимости одного удаления (буквы t из $rest$) к минимальной стоимости $c[3, 4] = 3$ преобразования $res \rightarrow stre$. Иначе стоимость $c[4, 4]$ можно получить путем прибавления стоимости одной вставки (буквы e в $stre$) к минимальной стоимости $c[4, 3] = 3$ преобразования $rest \rightarrow str$. Таким образом,

$$c[4, 4] = c[3, 4] + 1 = c[4, 3] + 1.$$

Мы также видим, что стоимость $c[5, 6] = 3$ полного преобразования $rests \rightarrow stress$ можно получить из минимальной стоимости $c[4, 5] = 3$ преобразования $rest \rightarrow stres$, поскольку стоимость преобразования конечных букв $s \rightarrow s$ равна нулю. Таким образом, $c[5, 6] = c[4, 5] = 0$.

Нетрудно заметить, что фактически минимальную стоимость можно получить тремя различными путями. Например, для преобразования $rest \rightarrow stress$ имеем

$$c[4, 6] = c[3, 6] + 1 = c[4, 5] + 1 = c[3, 5] + 2.$$

Последние равенства являются частным случаем отношений базисной рекурсии, которая позволяет вычислить любой элемент массива стоимостей c на основе начальных значений $c[0, 0]$, $c[0, j]$, $j = 1, 2, \dots, n_2$ и $c[i, 0]$, $i = 1, 2, \dots, n_1$. Имеет место следующая лемма.

¹Здесь шутка автора: *rest* — покой, сон, *stress* — стресс, напряжение. — Примеч. пер.

Лемма 9.1.1. Для всех $i \in 1..n_1, j \in 1..n_2$

$$c[i, j] = \min\{c[i - 1, j] + d(x_1[i], \varepsilon), \\ c[i, j - 1] + d(\varepsilon, x_2[j]), \\ c[i - 1, j - 1] + d(x_1[i], x_2[j])\}.$$

Доказательство. Основой доказательства служит факт, установленный в упр. 2.2.13: во время преобразования строк операции вставки и удаления можно выполнять в любом порядке. Этот факт порождает три возможности выполнения операций вставки и удаления при манипулировании буквами $x_1[i]$ и $x_2[j]$.

- Первая возможность — удаляется буква $x_1[i]$, стоимость операции равна $d(x_1[i], \varepsilon)$, которая прибавляется к стоимости $c[i - 1, j]$.
- Вторая возможность — вставляется буква $x_2[j]$, стоимость операции равна $d(\varepsilon, x_2[j])$, которая прибавляется к стоимости $c[i, j - 1]$.
- Третья возможность — подстановка, т.е. замена буквы $x_1[i]$ буквой $x_2[j]$ (буква $x_1[i]$ удаляется, буква $x_2[j]$ вставляется). Стоимость этой операции $d(x_1[i], x_2[j])$ прибавляется к стоимости $c[i - 1, j - 1]$. Конечно, если $x_1[i] = x_2[j]$, то $d(x_1[i], x_2[j]) = 0$.

Поскольку не существует других возможных манипуляций буквами $x_1[i]$ и $x_2[j]$ (с учетом “несимметричности” этих букв), то минимальная стоимость $c[i, j]$ должна быть минимумом этих трех перечисленных возможностей. ■

Повторим, что основой доказательства этой леммы служит факт, что операции вставки и удаления можно выполнять в любом порядке. Прилежный читатель, который выполнил упражнение 2.2.13, мог заметить, что это свойство операций вставки и удаления *не зависит* от условия (2.4) симметричности метрики. Оно имеет место даже в том случае, когда для некоторых строк x_1 и x_2 $d(x_1, x_2) \neq d(x_2, x_1)$. Таким образом, лемма 9.1.1 справедлива даже тогда, когда не выполняется условие симметричности. Это замечание подводит нас к мысли, что стоимость преобразований по своей природе не является полным аналогом расстояния. Поскольку мы определили стоимость в терминах преобразования строки x_1 в строку x_2 , то возможны случаи, когда стоимость “обратного” преобразования строки x_2 в строку x_1 не будет совпадать со стоимостью “прямого” преобразования. Такая ситуация возникает, например, в приложениях молекулярной биологии, где матрица весов W в некоторых случаях может быть несимметричной, что часто порождает неравенства $d_W(x_1, x_2) \neq d_W(x_2, x_1)$.

Наконец, отметим, что массив стоимостей и схема рекурсии для его вычисления, данная в лемме 9.1.1, очень просты и естественны. Поэтому они имеют большое применение не только для вычисления расстояний между строками (а значит, и LCS), но и, как основа почти всех алгоритмов, для вычисления наименьшей

общей сверхстроки для множества строк и для нахождения “наилучшего” выравнивания набора строк. Эти NP-полные задачи весьма важны для вычислительной биологии.

Упражнения 9.1

1. Приведите алгоритм вычисления расстояния Хемминга между двумя строками длиной n , выполняемый за время порядка $\Theta(n)$.
2. Докажите формулы для вычисления начальных значений $c[0, 0]$, $c[0, j]$ и $c[i, 0]$ массива стоимостей.

9.2 Алгоритм Вагнера–Фишера

Алгоритм Вагнера–Фишера (Wagner–Fischer) вычисляет массив стоимостей и, следовательно, расстояние между двумя заданными строками $x_1[1..n_1]$ и $x_2[1..n_2]$. Он, конечно, основывается на лемме 9.1.1. Его временная и пространственная сложности имеют порядок $\Theta(n_1 n_2)$ (доказательство этого утверждения оставим для упр. 9.2.1). Отметим, что этот алгоритм в определенном смысле онлайнный (раздел 4.1), поскольку в процессе вычисления расстояния $d(x_1, x_2)$ также вычисляются расстояния между каждой парой префиксов строк x_1 и x_2 . В разделе 9.3 покажем простой прием, который позволит в алгоритме Вагнера–Фишера уменьшить требуемый объем памяти до порядка $\Theta(\min\{n_1, n_2\})$.

Как указывается в работе [202], алгоритм вычисления расстояния между строками, основанный на массиве стоимостей, неоднократно независимо переоткрывался при проведении исследований в таких областях, как распознавание речи, автоматическая проверка орфографии, молекулярная биология и др. [228, 192, 227, 201, 203, 195, 113, 229]. Среди всех этих вариантов одного, по сути, алгоритма мы выбрали работу Вагнера и Фишера [229] для представления здесь данного алгоритма вследствие того, что алгоритм Вагнера–Фишера, кроме основного назначения, выполняет также “трассировку”, позволяя точно отображать последовательность выполняемых операторов редактирования, а также $LCS(x_1, x_2)$.

Предположим, что для заданных строк $x_1[1..n_1]$ и $x_2[1..n_2]$ уже вычислен массив стоимостей $c[0..n_1, 0..n_2]$. Для произвольной позиции $q = [i, j]$, $1 \leq i \leq n_1$, $1 \leq j \leq n_2$, элемента массива c определим *перемещение* $v : q \rightarrow q'$ как функцию отображения позиций q в предыдущую позицию q' согласно следующим правилам.

- Если $c[i, j] = c[i - 1, j] + d(x_1[i], \varepsilon)$, то $q' = [i - 1, j]$;
- иначе, если $c[i, j] = c[i, j - 1] + d(\varepsilon, x_2[j])$, то $q' = [i, j - 1]$;
- в противном случае $q' = [i - 1, j - 1]$.

Таким образом, перемещение определяет одну из трех возможных позиций в массиве c , исходя из которой с помощью одной операции редактирования получено текущее значение $c[i, j]$, при этом операторы редактирования выбираются согласно порядку

$$\langle \text{удаление, вставка, подстановка} \rangle. \tag{9.1}$$

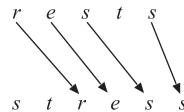
Поскольку функция v уменьшает на единицу по крайней мере одно из значений i или j , то существует такое целое число $k \in 1..i + j$, что суперпозиция $v^k([i, j])$ приведет к позиции $[i', j']$ (т.е. $v^k([i, j]) = [i', j']$), где $\min\{i', j'\} = 0$. В примере из предыдущего раздела (см. табл. 9.1) мы можем отследить последовательность

$$v([5, 6]) = [4, 5]; v([4, 5]) = [3, 5]; v([3, 5]) = [2, 4]; v([2, 4]) = [1, 3]; v([1, 3]) = [0, 2].$$

В этом случае $k = 5$ и $v^5([5, 6]) = [0, 2]$. Теперь, если для каждой операции перестановки (но не операций удаления и вставки) из этой последовательности перемещений записать ее текущее значение $[i, j]$, то в результате получим так называемый *след* (trace), который обычно обозначается буквой τ . В нашем примере получаем след

$$\tau([5, 6]) = \{[5, 6], [3, 5], [2, 4], [1, 3]\},$$

который соответствует подстановкам



В этом примере каждая “подстановка” просто замещает одну букву на саму себя. Вследствие этого стоимость таких подстановок равна нулю. Но, как показано в упр. 9.2.3, это не является обязательным свойством следа — след может содержать позиции, на которых выполняются подстановки с ненулевой стоимостью.

В общем случае след $\tau([n_1, n_2])$ определяет все подстановки (включая все подстановки нулевой стоимости), которые содержатся в последовательности операторов редактирования минимальной стоимости, необходимых для преобразования строки x_1 в строку x_2 . Отметим, что если позиция i , соответствующая букве $x_1[i]$, не входит в след τ , то это означает, что данная буква удаляется в процессе преобразования $x_1 \rightarrow x_2$. Но если в следе τ не встречается позиция j , соответствующая букве $x_2[j]$, то тогда эта буква вставляется в процессе преобразования. Таким образом след можно использовать для прямого и полного определения всей последовательности операторов редактирования минимальной стоимости, выполняемых в процессе преобразования $x_1 \rightarrow x_2$. В нашем примере отсутствующим в следе позициям соответствуют буква $x_1[4] = t$, которая удаляется, и подстрока $x_2[1..2] = st$, которая вставляется. Эту ситуацию, если использовать символ

пробела Δ , наглядно можно представить в виде следующей таблицы.

$$\begin{array}{cccccc} \Delta & \Delta & r & e & s & t & s \\ s & t & r & e & s & \Delta & s \end{array}$$

Это представление показывает оптимальное (минимальной стоимости) выравнивание строк x_1 и x_2 . (Как указывалось ранее, задача нахождения наилучшего выравнивания строковых последовательностей весьма важна для вычислительной биологии.) Мы оставим для упр. 9.2.6 детальное описание алгоритмов, находящих последовательность операторов редактирования и вычисляющих след.

Далее отметим, что в следе подстановки нулевой стоимости определяют общую подпоследовательность строк x_1 и x_2 , хотя, как показано в упр. 9.2.3, эта общая подпоследовательность не обязана быть *наибольшей* общей подпоследовательностью $LCS(x_1, x_2)$. Но если мы ограничимся только расстоянием Левенштейна $d_L = d$, для которого $d(\lambda, \mu) = 2$ на любой паре несовпадающих букв λ и μ , то нетрудно видеть, что для любых строк $x_1[1..n_1]$ и $x_2[1..n_2]$ выполняется равенство

$$c[n_1, n_2] = n_1 + n_2 - 2|LCS(x_1, x_2)|. \tag{9.2}$$

Отсюда следует, что при условии $d = d_L$ длину $LCS(x_1, x_2)$ можно вычислить по формуле

$$|LCS(x_1, x_2)| = (n_1 + n_2 - c[n_1, n_2])/2. \tag{9.3}$$

Как показано в упр. 9.2.4, саму подпоследовательность $LCS(x_1, x_2)$ несложно вычислить из соответствующего следа, если $d = d_L$.

В работе [229] приведено еще одно интересное применение массива стоимостей, а именно для проверки орфографии текстов. Здесь используется расстояние $d = d_W$ с матрицей весов W , в которой отображен тот факт, что при печати на стандартной клавиатуре ошибочные буквы (опечатки) чаще всего замещают наиболее вероятные буквы (т.е. те, которые наиболее часто встречаются в текстах на данном языке), при этом ошибочные буквы, как правило, на клавиатуре располагаются рядом с “правильными” буквами. Подобная ситуация имеет место в матрицах весов для последовательностей ДНК. Здесь наиболее часто встречается пара нуклеотидов C и G , которая может замещать более редкие пары C и T или G и A . В обоих этих приложениях становится важной операция подстановки специального типа: перестановка соседних символов. Как показано в [229], если перестановка включена в набор основных операторов редактирования, то расстояние между строками x_1 и x_2 можно вычислить за время порядка $\Theta(n_1 n_2)$, предполагая при этом, что стоимость перестановки несовпадающих букв λ и μ не меньше, чем $(d(\varepsilon, \lambda) + d(\mu, \varepsilon))/2$.

Упражнения 9.2

1. Приведите алгоритм вычисления массива стоимостей для двух заданных строк $x_1[1..n_1]$ и $x_2[1..n_2]$, выполняемый за время порядка $\Theta(n_1n_2)$.
2. Докажите, что если для расстояния d выполняется условие симметричности, то тогда массив стоимостей для преобразования строки x_2 в строку x_1 совпадает с транспонированным массивом стоимостей для преобразования строки x_1 в строку x_2 . Проверьте свое доказательство путем вычисления массива c для строк $x_1 = stress$ и $x_2 = rests$.
3. Для расстояния преобразования $d = d_E$ вычислите массив стоимостей и след, которые соответствуют строкам $x_1 = sowsear$ и $x_2 = silkpurse$. Покажите, что в данном примере не все позиции в следе соответствуют подстановкам нулевой стоимости. Также покажите, что общая подпоследовательность, определяемая на основе следа, не является подпоследовательностью $LCS(x_1, x_2)$.
4. Для расстояния Левенштейна $d = d_L$ вычислите массив стоимостей и след, которые соответствуют строкам $x_1 = sowsear$ и $x_2 = silkpurse$. Покажите, что в данном примере общая подпоследовательность, определяемая на основе следа, является подпоследовательностью $LCS(x_1, x_2)$.
5. Для чего необходимо было устанавливать порядок выполнения операторов редактирования (9.1)? Влияет ли этот порядок на правильность вычислений следа и соответствующей последовательности выполняемых операторов редактирования?
6. Предложите алгоритм, который вычислял бы след на основе ранее вычисленного массива стоимостей $c[0..n_1, 0..n_2]$ за время порядка $O(n_1 + n_2)$. Предложите также алгоритм с временем выполнения $\Theta(n_1 + n_2)$ для вычисления последовательности операторов редактирования.
7. Докажите равенство (9.2).
8. Предложите алгоритм, который вычислял бы $LCS(x_1, x_2)$ на основании ранее вычисленного следа для расстояния $d = d_L$ за время порядка $O(\min\{n_1, n_2\})$.
9. Разработайте матрицу весов W для алфавита английского языка (только для строчных букв) таким образом, чтобы учесть возрастающую вероятность замены одной буквы на другую, если этим буквам соответствуют рядом расположенные клавиши на стандартной клавиатуре пишущей машинки. Например, можно положить

$$W[d, e] = W[d, r] = W[d, s] = W[d, f] = W[d, x] = W[d, c] = 1/2,$$

тогда как $W[d, \lambda] = 1$ для всех других букв λ . Затем вычислите массив стоимостей для $x_1 = sowsear$ и $x_2 = silkpurse$ при условии, что $d = d_W$.

9.3 Алгоритмы Хешберга

Хешберг (Hirschberg) [117] предложил алгоритм непосредственного вычисления $\text{LCS}(x_1, x_2)$ за время порядка $\Theta(n_1 n_2)$ без явного вычисления массива стоимостей c . Он также показал, что этот алгоритм требует память объемом только $\Theta(\min\{n_1, n_2\})$.

Для описания этого алгоритма сначала определим массив $\gamma[0..n_1, 0..n_2]$ (аналог массива стоимостей), в котором для всех $i \in 0..n_1$ и $j \in 0..n_2$

$$\gamma[i, j] = |\text{LCS}(x_1[1..i], x_2[1..j])|.$$

Таким образом, $\gamma[i, j]$ равняется длине наибольшей общей подпоследовательности подстрок $x_1[1..i]$ и $x_2[1..j]$. Поскольку пустая строка ε имеет LCS нулевой длины с любой другой строкой, то очевидно, что $\gamma[i, 0] = \gamma[0, j] = 0$. Другие элементы массива γ вычисляются на основе следующей леммы (аналог леммы 9.1.1).

Лемма 9.3.1. Для расстояния Левенштейна $d = d_L$ и для всех $i \in 1..n_1$ и $j \in 1..n_2$

- $\gamma[i, j] = \gamma[i - 1, j - 1] + 1$, если $x_1[i] = x_2[j]$,
- $\gamma[i, j] = \max\{\gamma[i - 1, j], \gamma[i, j - 1]\}$ в противном случае.

Доказательство. Для доказательства достаточно заметить, что для расстояния $d = d_L$ выполняются равенства (9.2) и (9.3), если в них n_1 и n_2 заменить соответственно на $i \in 1..n_1$ и $j \in 1..n_2$, а затем следует применить лемму 9.1.1 для случая $d = d_L$. ■

Алгоритм 9.3.1 (Алгоритм Н1)

```

▷ Вычисление длины  $\text{LCS}(x_1, x_2)$ 
▷ Элементам первой строки и первого столбца
  массива  $\gamma$  присваиваются нулевые значения
for  $i \leftarrow 0$  to  $n_1$  do
   $\gamma[i, 0] \leftarrow 0$ 
for  $j \leftarrow 0$  to  $n_2$  do
   $\gamma[0, j] \leftarrow 0$ 
▷ Вычисление массива  $\gamma$ 
for  $i \leftarrow 1$  to  $n_1$  do
  for  $j \leftarrow 1$  to  $n_2$  do
    if  $x_1[i] = x_2[j]$  then
       $\gamma[i, j] \leftarrow \gamma[i - 1, j - 1] + 1$ 
    else
       $\gamma[i, j] \leftarrow \max\{\gamma[i - 1, j], \gamma[i, j - 1]\}$ 

```

Алгоритм Н1 является простой другой формой записи леммы 9.3.1. Мы привели этот тривиальный алгоритм не для того, чтобы бросить вызов бесспорному

интеллекту читателя, а для того, чтобы позднее можно было сравнить его с другой улучшенной версией этого алгоритма, также основанной на лемме 9.3.1.

Хешберг заметил, что в алгоритме Н1 для определения элементов строки i массива γ необходимо знать только элементы строки $i - 1$. Поэтому в любой момент выполнения алгоритма достаточно хранить в памяти только две строки этого массива. Для этого можно задать массив $\gamma'[0..1, 0..n_2]$ размерностью $2 \times n_2$, где хранились бы текущая (i -я) и предыдущая ($(i - 1)$ -я) строки массива γ . Введем *триггерную* (переключающую) *переменную* $i' \in 0..1$, с помощью которой осуществлялся бы доступ к элементам $\gamma'[i', j]$ текущей строки, тогда как значение $1 - i'$ определяло бы элементы $\gamma'[1 - i', j]$ предыдущей. В результате этих нововведений получим новый алгоритм 9.3.2, который возвращает последнюю строку $\gamma[n_1, 0..n_2]$ массива γ . Этот алгоритм запишем как функцию Н2.

Алгоритм 9.3.2 (Алгоритм Н2)

▷ *Вычисление длины* $\text{LCS}(x_1, x_2)$

function Н2(n_1, x_1, n_2, x_2)

▷ *Присвоение значения 0 элементу* $\gamma'[1, 0]$ *и элементам строки* $\gamma'[0..1, 0..n_2]$

$\gamma'[1, 0] \leftarrow 0$

for $j \leftarrow 0$ **to** n_2 **do**

$\gamma'[0, j] \leftarrow 0$

▷ *Инициализация триггерной переменной*

$i' \leftarrow 0$

▷ *Вычисление* $\gamma'[1 - i', j]$

for $i \leftarrow 1$ **to** n_1 **do**

$i' \leftarrow 1 - i'$

for $j \leftarrow 1$ **to** n_2 **do**

if $x_1[i] = x_2[j]$ **then**

$\gamma'[i, j] \leftarrow \gamma'[1 - i', j - 1] + 1$

else

$\gamma'[i, j] \leftarrow \max\{\gamma'[1 - i', j], \gamma'[i', j - 1]\}$

▷ $\gamma'[i, j] = |\text{LCS}(x_1, x_2[1..j])|$, $j = 0, 1, \dots, n_2$

return $\gamma'[i', 0..n_2]$

Отметим, что использование триггерной переменной делает алгоритм Н2 значительно более эффективным (по сравнению с алгоритмом Н1), поскольку позволяет уменьшить требуемый объем памяти до величины порядка $\Theta(n_2)$. Однако тот же результат можно получить, если подобные вычисления выполняются не относительно строк, а относительно столбцов массива γ . Поэтому можно говорить о том, что в данном алгоритме требуемый объем памяти составляет величину порядка $\Theta(\min\{n_1, n_2\})$. В упр. 9.3.3 предложено применить подобный подход

к алгоритму Вагнера–Фишера для того, чтобы в этом алгоритме уменьшить требуемый объем памяти до величины такого же порядка.

До сих пор все было хорошо. Однако напомним, мы пока можем вычислить только *длину* наибольшей общей подпоследовательности $\text{LCS}(x_1, x_2)$, а не саму эту подпоследовательность. Для решения этой задачи можно предложить алгоритм, где рекурсивно вызывалась бы функция H2. Для описания такого алгоритма необходимо ввести еще несколько понятий. Обозначим через Γ_{ij} подпоследовательность $\text{LCS}(x_1[1..i], x_2[1..j])$, а через $\Gamma_{ij}^* = \text{LCS}(x_1[i+1..n_1], x_2[j+1..n_2])$. Тогда $\gamma[i, j] = |\Gamma_{ij}|$ и $\gamma^*[i, j] = |\Gamma_{ij}^*|$. Кроме того, обозначим

$$\gamma_i^{\max} = \max_{0 \leq j \leq n_2} \{\gamma[i, j] + \gamma^*[i, j]\}.$$

Далее мы докажем, что для всех $i \in 0..n_1$ $\gamma_i^{\max} = \gamma[n_1, n_2]$. Отсюда следует, что наибольшую общую последовательность строк x_1 и x_2 можно найти путем разделения строки x_1 на префикс $x_1[1..i]$ и суффикс $x_1[i+1..n_1]$ (позиция i выбирается произвольно) и последующего вычисления функций $\text{H2}(i, x_1[1..i], n_2, x_2)$ и $\text{H2}(n_1 - i, x_1[i+1..n_1], n_2, x_2)$, что позволит найти $\gamma[i, n_2]$ и $\gamma^*[i, 0]$.

Эту схему вычислений можно применить рекурсивно. Тогда, если на каждом шаге рекурсии выбирать $i = \lfloor n_1/2 \rfloor$, после $\log_2 n_1$ шагов вычисление $\gamma[n_1, n_2]$ можно свести к суммированию нулей и единиц. При этом для каждой полученной единицы можно проверить, найдена ли буква, входящая в LCS. Более четко и детально этот алгоритм мы опишем ниже, а пока докажем следующую лемму.

Лемма 9.3.2. Для всех $i \in 0..n_1$ $\gamma_i^{\max} = \gamma[n_1, n_2]$.

Доказательство. Для произвольного номера i определим j_0 , такое, что $\gamma_i^{\max} = \gamma[i, j_0] + \gamma^*[i, j_0]$ для соответствующих LCS Γ_{ij_0} и $\Gamma_{ij_0}^*$. Тогда строка $\Gamma = \Gamma_{ij_0} \Gamma_{ij_0}^*$ будет наибольшей общей подпоследовательностью строк x_1 и x_2 длиной

$$\gamma_i^{\max} \leq \gamma[n_1, n_2]. \tag{9.4}$$

С другой стороны, если $\Gamma_{n_1 n_2}$ — LCS строк x_1 и x_2 , то $\Gamma_{n_1 n_2}$ является подпоследовательностью строки x_2 , и ее можно записать в форме $\Gamma_1 \Gamma_2$, где Γ_1 — подпоследовательность подстроки $x_1[1..i]$, Γ_2 — подпоследовательность подстроки $x_1[i+1..n_1]$, а значение i такое же, как в неравенстве (9.4). Тогда существует номер j_1 , такой, что Γ_1 будет подпоследовательностью подстроки $x_2[1..j_1]$, а Γ_2 — подпоследовательностью подстроки $x_2[j_1+1..n_2]$. В этом случае

$$\begin{aligned} |\Gamma_1| &\leq |\text{LCS}(x_1[1..i], x_2[1..j_1])| = \gamma[i, j_1]. \\ |\Gamma_2| &\leq |\text{LCS}(x_1[i+1..n_1], x_2[j_1+1..n_2])| = \gamma^*[i, j_1]. \end{aligned}$$

Следовательно, мы вправе записать

$$\gamma[n_1, n_2] = |\Gamma_{n_1 n_2}| = |\Gamma_1| + |\Gamma_2| \leq \gamma[i, j_1] + \gamma^*[i, j_1] \leq \gamma_i^{\max}. \quad (9.5)$$

Из неравенств (9.4) и (9.5) следует, что при любом выборе i выполняется равенство $\gamma_i^{\max} = \gamma[n_1, n_2]$, что и требовалось доказать. ■

Основываясь на результате доказанной леммы, можно записать алгоритм вычисления наибольшей общей последовательности (алгоритм 9.3.3) и доказать его корректность. Этот алгоритм снова реализуется в виде функции, которая возвращает искомую LCS.

Алгоритм 9.3.3 (Алгоритм НЗ)

▷ Вычисление $LCS(x_1, x_2)$

function НЗ(n_1, x_1, n_2, x_2)

▷ Вычисление LCS для строк длиной 0 и 1

if $\min\{n_1, n_2\} = 0$ **then return** $\Gamma \leftarrow \varepsilon$

else if $n_1 = 1$ **then**

if $\exists j \in 1..n_2$, такое, что $x_1[1] = x_2[j]$ **then**

return $\Gamma \leftarrow x_1[1]$

else

return $\Gamma \leftarrow \varepsilon$

▷ Шаг рекурсии: x_1 разбивается пополам

else

$i \leftarrow \lfloor n_1/2 \rfloor$

$\gamma'_1[1..n_2] \leftarrow H2(i, x_1[1..i], n_2, x_2[1..n_2])$

$\gamma'_2[1..n_2] \leftarrow H2(n_1 - i, x_1[i + 1..n_1], n_2, x_2[1..n_2])$

Нахождение j_0 , при котором достигается минимум $\gamma'_1[j] + \gamma'_2[j]$

▷ Рекурсивный вызов функции НЗ для решения двух подзадач

$\Gamma_1 \leftarrow H3(i, x_1[1..i], j_0, x_2[1..j_0])$

$\Gamma_2 \leftarrow H3(n_1 - i, x_1[i + 1..n_1], n_2 - j_0, x_2[j_0 + 1..n_2])$

▷ Формирование окончательного решения на основе решений

▷ двух подзадач

return $\Gamma \leftarrow \Gamma_1 \Gamma_2$

Теорема 9.3.3. Алгоритм 9.3.3 корректно вычисляет $LCS(x_1, x_2)$.

Доказательство. Если n_1 или n_2 равно нулю, то в этом случае LCS является пустой строкой. Если $n_1 = 1$, тогда LCS не будет пустой строкой только тогда, когда $x_1[1]$ совпадает с какой-либо буквой строки $x_2[1..n_2]$.

Если $n_1 > 1$, значение $\gamma[n_1, n_2] = |LCS(x_1, x_2)|$, получаемое после двух вызовов функции H2, вычисляется корректно вследствие леммы 9.3.2. Эти вычисления определяют позицию $j_0 \in 0..n_2$, которая разбивает строку x_2 на две части.

К каждой части рекурсивно применяется функция НЗ для вычисления LCS Γ_1 и Γ_2 , соответствующих этим частям строки x_2 . Окончательный результат получается путем объединения Γ_1 и Γ_2 . ■

Чтобы определить временную сложность алгоритма НЗ, рассмотрим $T(n_1, n_2)$ — время выполнения этого алгоритма. Поскольку алгоритм Н2 выполняется за время порядка $\Theta(n_1, n_2)$, существует константа K , такая, что

$$T(n_1, n_2) \leq K n_1 n_2 + T(n_1/2, j_0) + T(n_1/2, n_2 - j_0).$$

(Здесь сделано допущение, что j_0 вычисляется за время порядка $\Theta(n_2)$, а время порядка $O(n_2)$, необходимое для нахождения j , такого, что $x_1[1] = x_2[j]$, учтено в константе K .) Предположим, что для всех $i \in 1..n_1 - 1$ и $j \in 1..n_2 - 1$ существует другая константа K' , такая, что $T(i, j) < K'ij$. Тогда

$$T(n_1, n_2) \leq K n_1 n_2 + K'(n_1/2)(j_0 + n_2 - j_0) = (K + K'/2)n_1 n_2,$$

и, если положим $K' > 2K$, получим $T(n_1, n_2) < K'n_1 n_2$. Отсюда по индукции мы заключаем, что алгоритм НЗ выполняется за время $\Theta(n_1 n_2)$.

Для оценки объема памяти, требуемого для алгоритма НЗ, заметим, что этот алгоритм можно организовать таким образом, что для реализации всех рекурсивных вызовов будет достаточно памяти, где хранятся строки x_1 и x_2 , а также два вектора γ'_1 и γ'_2 . Таким образом, мы доказали такую теорему.

Теорема 9.3.4. Алгоритм 9.3.3 вычисляет наибольшую общую подпоследовательность для строк $x_1[1..n_1]$ и $x_2[1..n_2]$ за время порядка $\Theta(n_1 n_2)$ с использованием памяти объемом $\Theta(n_2)$. ■

Хешберг разработал три алгоритма вычисления LCS, первый из которых опубликован в 1977 году [114], два других опубликованы двумя годами позднее [115]. Первый алгоритм описан выше и является в общем случае самым медленным из трех. Два других работают быстрее, когда длина $\text{lcs} = |\text{LCS}(x_1, x_2)|$ наибольшей общей подпоследовательности мала либо, наоборот, достаточно большая. Чтобы быть более точным, предположим, что $n_1 < n_2$, и пусть $\ell = \text{lcs}(x_1, x_2)$. Тогда второй алгоритм Хешберга вычисляет LCS за время порядка $O(\ell n_2 + n_2 \log \alpha')$, где α' ($\alpha' \leq \alpha$) — количество различных букв в строке x_2 . Третий алгоритм Хешберга выполняется за время порядка $O((n_1 + 1 - \ell)\ell \log n_2)$. Второй алгоритм эффективен, когда LCS короткая и алфавит не велик, но может потребовать в самом худшем случае времени порядка $\Theta(n_1 n_2)$. Третий алгоритм эффективен, когда LCS или очень короткая или очень длинная и может потребовать в самом худшем случае времени порядка $\Theta(n_1^2 \log n_2)$.

Упражнения 9.3

1. Докажите лемму 9.3.1.
2. Опишите изменения, которые необходимо сделать в алгоритме H2 для того, чтобы вычисления производились относительно столбцов (а не относительно строк, как в алгоритме 9.3.2).
3. Предложите усовершенствования в алгоритме Вагнера–Фишера, аналогичные алгоритму H2.
4. При доказательстве неравенства $T(n_1, n_2) < K'n_1n_2$ для алгоритма H2 сказано “по индукции мы заключаем...”. Запишите эту индукцию в явном виде.
5. Покажите по индукции, что общее количество вызовов функции H3 в точности равно $2n_1 - 1$.
6. Пусть k и k' — фиксированные действительные числа, удовлетворяющие условию $0 < k < k' < 1$. Рассмотрите класс \mathcal{C} всех пар (x_1, x_2) строковых последовательностей, таких, что

$$kn \leq \text{lcs}(x_1, x_2) \leq k'n,$$

где $n = \min\{|x_1|, |x_2|\}$. Охарактеризуйте сложность второго и третьего алгоритмов Хешберга, выполняемых на строках класса \mathcal{C} .

9.4 Алгоритм Ханта–Шиманского

В этом разделе мы представим еще один алгоритм, который непосредственно вычисляет наибольшую общую подпоследовательность $\text{LCS}(x_1, x_2)$ для строк $x_1[1..n_1]$ и $x_2[1..n_2]$ без явного вычисления расстояния $d(x_1, x_2)$. Конечно, если LCS найдена, то расстояние Левенштейна d_L можно вычислить на основании равенства (9.2).

Если алфавит упорядочен, алгоритм Ханта–Шиманского (Hunt–Szymanski, сокращенно — алгоритм ХШ) [122] имеет временную сложность $O((n+r)\log n)$ и пространственную $O(n+r)$, где $n = \max\{n_1, n_2\}$, r — общее количество совпадений $x_1[i] = x_2[j]$, $i \in 1..n_1$, $j \in 1..n_2$. В самом худшем случае (например, если $x_1 = x_2 = a^n$) алгоритм ХШ выполняется за время порядка $\Theta(n^2 \log n)$ и требует памяти порядка $\Theta(n^2)$. Однако обычно $r \in O(n)$, и если алфавит упорядочен, то алгоритм выполняется за время порядка $\Theta(n \log n)$ и требует памяти порядка $\Theta(n)$, что значительно лучше того, что могут предложить алгоритмы из предыдущего раздела.

Как и многие другие алгоритмические новшества, этот алгоритм основан на очень простой идее [2]: LCS определяется с помощью “наибольшей” строго убы-

вающей и непрерывной ломаной линии, проходящей через узлы сетки, представляющие совпадения $x_1[i] = x_2[j]$. На рис. 9.1 иллюстрируется эта идея на примере $LCS(rests, stress) = ress$. Здесь прочерченная непрерывная линия, определяющая LCS, содержит наибольшее возможное число диагональных “строго убывающих” сегментов. На этом рисунке можно найти *все* возможные общие подпоследовательности строк x_1 и x_2 . Так строки *rests* и *stress* имеют еще общую подпоследовательность *sts*. Отметим также, что хотя строки *rests* и *stress* и имеют относительно длинную LCS, но значение r относительно мало, удовлетворяя неравенствам $r < n_1 + n_2 < 2n$.

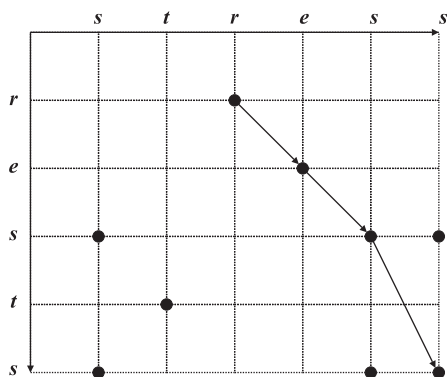


Рис. 9.1. Линия, формирующая LCS

Конечно, если работать непосредственно с сеткой размером $n_1 \times n_2$, то не получим улучшения нижней временной границы $\Omega(n_1 n_2)$ предыдущих алгоритмов. Вместо этого используется двухмерный массив j , в котором элемент $j[i, h]$ равен длине самого короткого префикса строки x_2 , который имеет LCS длиной h с подстрокой $x_1[1..i]$. Другими словами, $j[i, h]$ равен наименьшему целому числу j^* , такому, что

$$|LCS(x_1[1..i], x_2[1..j^*])| = h.$$

Например, для $x_1 = rests$ и $x_2 = stress$ имеем $j[4, 1] = 1$, $j[4, 2] = 2$, $j[4, 3] = 5$, при этом элементы $j[4, h]$ не определены для $h > 3$. Аналогично $j[5, 1] = 1$, $j[5, 2] = 2$, $j[5, 3] = 5$, $j[5, 4] = 6$ и $j[5, h]$ не определены для $h > 4$. В общем случае элементы $j[i, h]$ определены только тогда, когда $h \leq |LCS(x_1[1..i], x_2)| \leq i$. Для всех позиций $[i, h]$, где $j[i, h]$ не определены, “принудительно” положим $j[i, h] = n_2 + 1$. Для полной определенности также положим $j[i, 0] = 0$ для всех $i \in 1..n_1$, $j[0, 0] = 0$ и $j[0, h] = n_2 + 1$ для $h \in 1..n_1$. Отметим, что значения элементов массива j легко получить из рисунка, подобного рис. 9.1.

Установим некоторые свойства массива j , которые будут использованы в алгоритме Ханга–Шиманского.

Лемма 9.4.1. Для всех целых $i \in 1..n_1$ и $h \in 1..i$ выполняются неравенства

$$j[i-1, h-1] + 1 \leq j[i, h] \leq j[i-1, h].$$

Доказательство тривиально при $i = 1$, поэтому будем считать, что $i > 1$.

Заметим, что если элемент $j[i-1, h]$ определен, то также определены $j[i-1, h-1]$ и $j[i, h]$. Если подстроки $x_1[1..i-1]$ и $x_2[1..j[i-1, h]]$ имеют LCS длиной h , тогда (по определению массива j) подстроки $x_1[1..i]$ и $x_2[1..j[i-1, h]]$ имеют общую подпоследовательность длиной не менее h . Отсюда следует, что $j[i, h] \leq j[i-1, h]$.

Далее, поскольку элемент $j[i, h]$ определен, подстроки $x_1[1..i]$ и $x_2[1..j[i, h]]$ имеют общую подпоследовательность длиной h . Удалив самые правые буквы из этих подстрок, можно уменьшить длину общей подпоследовательности не более, чем на единицу. Поэтому подстроки $x_1[1..i-1]$ и $x_2[1..j[i, h]-1]$ имеют общую подпоследовательность и, следовательно, LCS длиной не менее $h-1$. Это означает, что $j[i-1, h-1] \leq j[i, h]$. Что и требовалось доказать. ■

Мы используем этот результат для того, чтобы как вычисляется значение $j[i, h]$, если уже известны значения $j[i-1, h-1]$ и $j[i-1, h]$. Этим будет построена основа рекурсии для алгоритма XIII.

Лемма 9.4.2. Для всех целых $i \in 1..n_1$ и $h \in 1..|LCS(x_1[1..i], x_2)|$ значение $j[i, h]$ вычисляется по следующему правилу:

$$j[i, h] \text{ равняется наименьшему целому } j^* \in j[i-1, h-1] + 1..j[i-1, h], \quad (9.6)$$

такому, что $x_1[i] = x_2[j^*]$ (если такое j^* существует);

$$j[i, h] = j[i-1, h], \text{ если такого } j^* \text{ не существует.}$$

Доказательство. Сначала рассмотрим, как более простой, второй вариант вычисления $j[i, h]$, когда такого j^* не существует. Поскольку $j[i, h]$ определяет *наименьший* префикс строки x_2 , то самым правым элементом в LCS для подстрок $x_1[1..i]$ и $x_2[1..j[i, h]]$ будет буква $x_2[j[i, h]]$. Вследствие леммы 9.4.1 $j[i, h] \in j[i-1, h-1] + 1..j[i-1, h]$. Но по нашему предположению не существует j^* из этого интервала, такого, что $x_2[j^*] = x_1[i]$, поэтому $x_1[i] \neq x_2[j[i, h]]$. Следовательно, та же самая LCS будет общей подпоследовательностью и для подстрок $x_1[1..i-1]$ и $x_2[1..j[i, h]]$ — в таком случае $j[i, h] \geq j[i-1, h]$. Вследствие леммы 9.4.1 справедливо также обратное неравенство. Отсюда делаем заключение, что $j[i, h] = j[i-1, h]$.

Теперь предполагаем, что существует наименьшее целое j^* из интервала $j[i-1, h-1] + 1..j[i-1, h]$, такое, что $x_1[i] = x_2[j^*]$. Из этого предположения вытекает, что подстроки $x_1[1..i]$ и $x_2[1..j^*]$ имеют общую подпоследовательность,

состоящую из $\text{LCS}(x_1[1..i-1], x_2[1..j[i-1, h-1]])$ длиной $h-1$ и буквы $x_2[j^*]$. Отсюда заключаем, что $j[i, h]$ не превышает j^* .

Теперь покажем, что $j[i, h]$ не меньше j^* . Предположим, что имеет место строгое неравенство $j[i, h] < j^*$. Тогда, вследствие леммы 9.4.1, $j[i-1, h-1] < j[i, h] < j^*$. Поскольку j^* — это *наименьшее* целое из интервала $j[i-1, h-1] + 1..j[i-1, h]$, такое, что $x_1[i] = x_2[j^*]$, то отсюда заключаем, что $x_1[i] \neq x_2[j[i, h]]$. Так как по условию $|\text{LCS}(x_1[1..i], x_2[1..j[i, h]])| = h$, поэтому подстроки $x_1[1..i-1]$ и $x_2[1..j[i, h]]$ также имеют общую подпоследовательность длиной h . Следовательно, $j[i-1, h] \leq j[i, h+1]$, а лемма 9.4.1 уточняет, что $j[i-1, h] = j[i, h+1]$. Но это невозможно, поскольку, как мы видели, $j[i, h+1] < j^* < j[i-1, h]$. Это противоречие и доказывает, что $j[i, h]$ не меньше j^* . Имея также доказанным утверждение, что $j[i, h]$ не превышает j^* , делаем вывод, что $j[i, h] = j^*$. Теорема доказана. ■

Лемма 9.4.2 фактически устанавливает для массива j такой же результат, что и лемма 9.3.1 для массива γ значения в строке i можно вычислить, зная только значения в строке $i-1$. В алгоритмах Н2 и Н3 этот результат привел к уменьшению требуемого объема памяти от величины порядка $\Theta(n_1 n_2)$ до величины порядка $\Theta(n_2)$. Здесь этот результат также позволяет уменьшить требуемый объем памяти и, кроме того, время выполнения алгоритма, по крайней мере в среднем.

Напомним, что в алгоритме Н2 был введен массив $\gamma'[1..2, 1..n_2]$ для замены массива γ . В алгоритме ХШ также заменяем массив j одномерным массивом $j'[0..n_2]$. Этот алгоритм можно описать в виде последовательности таких шагов.

Шаг 1. Вычисление массива MATCHLIST

Из леммы 9.4.2 следует, что для вычисления значений $j[i+1, 0..n_1]$ на основе значений $j[i, 0..n_1]$ надо для каждого значения $i \in 1..n_1$ знать такую позицию j в строке x_2 , для которой $x_1[i] = x_2[j]$. Эту информацию в подходящей форме будем получать из массива MATCHLIST, в котором каждая позиция i является указателем на список значений j , таких, что $x_1[i] = x_2[j]$. Для дальнейших вычислений на шаге 3 желательно, чтобы эти значения j были упорядочены в убывающем порядке.

Отметим, что если строки x_1 и x_2 определены на упорядоченном алфавите, то массив MATCHLIST можно вычислить за время порядка $O(n \log n + r)$ и с объемом памяти порядка $\Theta(n_1 + r)$, для чего можно использовать эффективные алгоритмы сортировки или поиска по дереву. (Здесь, как и выше, $n = \max\{n_1, n_2\}$.) Но если алфавит общего вида (не упорядочен), то вычисление массива MATCHLIST в самом худшем случае может потребовать времени порядка $O(n^2)$. Детали вычисления этого массива оставим для упр. 9.4.1. (См. также задачу 4.1 из раздела 4.1.)

Шаг 2. Начальное задание значений массива j'

Полагаем $j'[0] \leftarrow 0$ и для всех $h \in 1..n_2$ полагаем $j'[h] \leftarrow n_2 + 1$ (тем самым показываем, что элементы $j'[h]$ не определены).

Шаг 3. Вычисление значений $j[i, h]$ на основе значений $j[i-1, h]$, $i = 1, 2, \dots, n_1$. Эти вычисления выполняются путем замены значений в предыдущем массиве $j[0..n_1]$, соответствующем номеру $i-1$, на аналогичные значения, соответствующие номеру i . Напомним (лемма 9.4.2), что $j[i, h] = j[i-1, h]$ (тогда значение $j'[h]$ не изменяется) только в том случае, если не существует номера j^* из интервала $j[i-1, h-1] + 1..j[i-1, h]$, такого, что $x_1[i] = x_2[j^*]$. Поэтому для текущего значения i и для каждого значения $j^* \in \text{MATCHLIST}[i]$ по значениям массива j' надо найти такое значение h^* , что

$$j'[h^* - 1] + 1 \leq j^* \leq j'[h^*].$$

Тогда, согласно лемме 9.4.2,

$$j'[h^*] \leftarrow j^*. \tag{9.7}$$

Каждое такое присвоение, которое изменяет значение $j'[h^*]$, соответствует одной из r точек сетки (как показано на рис. 9.1), где $x_1[i] = x_2[j^*]$. Это присвоение можно интерпретировать как создание конечного узла (листа) $(i, j^*, \text{LINK}[h^* - 1])$ на синтаксическом дереве (см. раздел 2.1), где $\text{LINK}[h^* - 1]$ — указатель на узел, созданный ранее при присвоении (9.7) для $j'[h^* - 1]$. (Для $h^* = 1$ полагаем, что $\text{LINK}[h^* - 1] = \text{NULL}$.) Когда лист создан, указатель на него хранится в $\text{LINK}[h^*]$. После обработки всех $i \in 1..n_1$ наибольшее значение h , такое, что $j'[h] < n_2 + 1$, равняется длине $\text{LCS}(x_1, x_2)$, а $\text{LINK}[h]$ указывает путь длиной h по дереву от листа до корня. Этот путь и определяет LCS.

Шаг 4. Запись $\text{LCS}(x_1, x_2)$ в обратном порядке

На дереве находится первый узел, который соответствует наибольшему значению h , такому, что $j'[h] < n_2 + 1$. Затем записывается путь от этого узла до корня дерева — это и будет $\text{LCS}(x_1, x_2)$, записанная в обратном порядке.

Более подробно каждый описанный шаг показан в алгоритме 9.4.1.

Алгоритм 9.4.1 (Алгоритм Ханга–Шиманского)

▷ Вычисление $\text{LCS}(x_1, x_2)$

▷ Шаг 1: Вычисление массива MATCHLIST

Для каждого $i \in 1..n$ вычисляется $\text{MATCHLIST}[i]$ — указатель на список монотонно убывающих значений j , таких, что $x_1[i] = x_2[j]$

▷ Шаг 2: Начальное задание значений массива j'

$j'[0] \leftarrow 0$

for $i \leftarrow 1$ **to** n_1 **do**

```

     $j'[i] \leftarrow n_2 + 1$ 
    LINK[0]  $\leftarrow$  NULL
    ▷ Шаг 3: Вычисление значений  $j[i, h]$  на основе значений  $j[i - 1, h]$ ,  $1 \leq i \leq n_1$ 
    for  $i \leftarrow 1$  to  $n_1$  do
        for all  $j^* \in \text{MATCHLIST}[i]$  do
            поиск  $h^*$ , такого, что  $j'[h^* - 1] + 1 \leq j^* \leq j'[h^*]$ 
            if  $j^* < j'[h^*]$  then
                 $j'[h^*] \leftarrow j^*$ 
                LINK[ $h^*$ ]  $\leftarrow$  новый_узел( $i, j^*, \text{LINK}[h^* - 1]$ )
    ▷ Шаг 4: Запись  $\text{LCS}(x_1, x_2)$  в обратном порядке
    поиск наибольшего  $h$ , такого, что  $j'[h] < n_2 + 1$ 
    ПУТЬ  $\leftarrow$  LINK[ $h$ ]
    while ПУТЬ  $\neq$  NULL do
        output пара  $(i, j)$  из узла, указанного указателем ПУТЬ
        продвижение указателя ПУТЬ из текущего узла далее к корню дерева
        в соответствии с указателем LINK[ $h - 1$ ]
    
```

Чтобы пояснить действия, выполняемые алгоритмом XIII, рассмотрим наш предыдущий пример, где $x_1 = \text{rests}$ и $x_2 = \text{stress}$. Здесь, как видно из рис. 9.1,

```

MATCHLIST[1] = <3>
MATCHLIST[2] = <4>
MATCHLIST[3] = <6, 5, 1>
MATCHLIST[4] = <2>
MATCHLIST[5] = <6, 5, 1>
    
```

Значения элементов массива j' после i -й итерации следующие.

```

Начальные значения : 077777
     $i = 1$  : 037777
     $i = 2$  : 034777
     $i = 3$  : 014577
     $i = 4$  : 012577
     $i = 5$  : 012567
    
```

Подчеркнутые значения в массиве j' определяют LCS в обратном порядке:

$(i = 5, j^* = 6), (i = 3, j^* = 5), (i = 2, j^* = 4), (i = 1, j^* = 3)$.

Теорема 9.4.3. Для двух заданных строк $x_1[1..n_1]$ и $x_2[1..n_2]$, определенных на упорядоченном алфавите, алгоритм 9.4.1 корректно вычисляет $\text{LCS}(x_1, x_2)$ за

время порядка $O((n+r)\log n)$ с использованием памяти объемом порядка $O(n+r)$, где $n = \max\{n_1, n_2\}$, r — общее число пар (i, j) , таких, что $x_1[i] = x_2[j]$.

Доказательство корректности алгоритма 9.4.1 непосредственно следует из утверждений лемм 9.4.1 и 9.4.2.

Для анализа времени выполнения алгоритма и объема требуемой памяти рассмотрим каждый шаг алгоритма в отдельности. Как мы уже отмечали, если алфавит упорядочен, то на шаге 1 требуется время порядка $O(n\log n + r)$ и память порядка $\Theta(n+r)$. Шаг 2 выполняется за время порядка $O(n_1)$ и требует $\Theta(n_1)$ памяти для хранения массива j' .

На шаге 3 для выполнения цикла **for** требуется время порядка $O(n_1)$, плюс необходимо время для выполнения таких операций:

- а) вычисление r значений массива MATCHLIST,
- б) поиск значения h^* в массиве j' ,
- в) не более r модификаций массива j' и создание не более r узлов дерева.

Для выполнения операции а) требуется время порядка $\Theta(r)$, для выполнения операции б) можно применить алгоритм бинарного поиска со временем выполнения порядка $O(r\log n_1)$, для операции в) необходимо $O(r)$ времени. Таким образом, общее время выполнения шага 3 не превышает величины порядка $O(n_1 + r\log n_1)$ с использованием памяти объемом порядка $O(r)$.

Шаг 4 требует $O(n_1)$ времени для определения позиции h в массиве j' и $O(r)$ времени для вывода LCS. Суммируя время и объемы памяти, необходимые для выполнения каждого шага алгоритма, получаем требуемый результат. ■

В работе [17] предложен вариант алгоритма XIII, который выполняется за время порядка $O(n_1\log n_2 + t\log(2n_1n_2/t))$, где t — количество “главных” совпадений между буквами строк x_1 и x_2 . Поскольку $t \leq r$, максимальное время выполнения этого алгоритма составляет величину порядка $\Theta(n_1n_2)$, тогда как минимальное время — величину порядка $\Theta(n_1\log n_2)$.

Упражнения 9.4

1. Предложите алгоритм вычисления массива MATCHLIST в алгоритме XIII, предполагая, что
 - а) алфавит упорядочен,
 - б) алфавит не упорядочен.
2. В каждом варианте алгоритма укажите его асимптотическую сложность в самом худшем и самом лучшем случаях и приведите примеры строк x_1 и x_2 , на которых достигаются эти границы.

3. В алгоритме XIII, как и в других алгоритмах вычисления LCS, строки x_1 и x_2 можно менять местами. На основе этого покажите, что временную и пространственную границы сложности этого алгоритма, данные в теореме 9.4.1, можно немного усилить.

9.5 Алгоритм Укконена–Майерса

В этом разделе мы возвращаемся к основной идее алгоритма Вагнера–Фишера: итерационному процессу вычисления массива стоимостей $c = c[0..n_1, 0..n_2]$, основанному на рекуррентных соотношениях леммы 9.1.1. Однако здесь мы представим массив c в виде ориентированного ациклического графа G с одной вершиной-источником, соответствующей позиции $[0, 0]$ массива c , и вершинами-стоками, соответствующими позициям $[n_1, j]$ и $[i, n_2]$ ($1 \leq i \leq n_1, 1 \leq j \leq n_2$) этого массива, включая позицию $[n_1, n_2]$. Мы снова используем лемму 9.1.1 для нахождения пути наименьшей стоимости (“кратчайшего” пути) от $c[0, 0]$ до $c[n_1, n_2]$. Одновременно с этим используем само значение расстояния $d = d(x_1, x_2)$ как “инструмент” для распознавания тех вершин и дуг графа G , которые невозможно эффективно вычислить. В результате применения этих соображений получим алгоритм, который назовем алгоритмом Укконена–Майерса (сокращенно, УМ). Этот алгоритм, независимо предложенный Укконеном (Ukkonen) [223] и Майерсом (Myers) [189], вычисляет d за время порядка $O(nd)$, где $n = \min\{n_1, n_2\}$. Поскольку обычно значение d значительно меньше значения n , то алгоритм УМ представляет несомненный интерес.

Существенно, что этот же алгоритм можно использовать в других часто встречающихся на практике ситуациях, где значение d представляет верхнюю границу (или пороговое значение) неизвестного значения расстояния $d(x_1, x_2)$ — можно модифицировать алгоритм таким образом, чтобы он возвращал или значение $d(x_1, x_2) \leq d$, или значение ∞ (если это расстояние вычислить невозможно).

Как уже указывалось, вершины графа $G = (V, E)^2$ соответствуют позициям $[i, j]$ массива c для всех $i \in 0..n_1$ и $j \in 0..n_2$. Каждая вершина $[i, j] \in V$ имеет не более трех входящих дуг $([i', j'], [i, j])$:

$$\begin{aligned} ([i-1, j], [i, j]) \in E, & \quad \text{если} \quad c[i, j] - c[i-1, j] = d(x_1[i], \varepsilon); \\ ([i, j-1], [i, j]) \in E, & \quad \text{если} \quad c[i, j] - c[i, j-1] = d(\varepsilon, x_2[j]); \\ ([i-1, j-1], [i, j]) \in E, & \quad \text{если} \quad c[i, j] - c[i-1, j-1] = d(x_1[i], x_2[j]). \end{aligned} \tag{9.8}$$

Для любой вершины $[i, j] \in V$ в соответствии с леммой 9.1.1 существует хотя бы одна входящая дуга только в том случае, если $[i, j] \neq [0, 0]$. Припишем каждой дуге $([i', j'], [i, j])$ вес $c[i, j] - c[i', j']$.

²Здесь использовано стандартное обозначение графа $G = (V, E)$ как совокупность двух множеств: V — множество вершин, E — множество дуг графа. — *Примеч. ред.*

Назовем граф G *графом зависимостей* массива c . Нетрудно заметить, что

- $|E| \leq 3(|V| - 1)$;
- граф G является подграфом графа, структура которого в точности соответствует структуре массива c (в графе G отсутствуют дуги, веса которых не равны минимальной стоимости);
- существуют ориентированные пути от вершины $[0, 0]$ до каждой вершины $[i, j] \in V$;
- сумма весов дуг, составляющих ориентированный путь от вершины $[0, 0]$ до вершины $[i, j]$, равна стоимости последовательности операций редактирования минимальной стоимости, необходимых для преобразования подстроки $x_1[1..i]$ в подстроку $x_2[1..j]$;
- граф G имеет одну вершину-источник $[0, 0]$ и не более $n_1 + n_2 - 1$ вершин-стоков, все из которых соответствуют или элементам строки n_1 , или элементам столбца n_2 массива c .

На рис. 9.2 показан граф зависимостей для примера $x_1 = rests$ и $x_2 = stress$, где использовано расстояние Левенштейна.

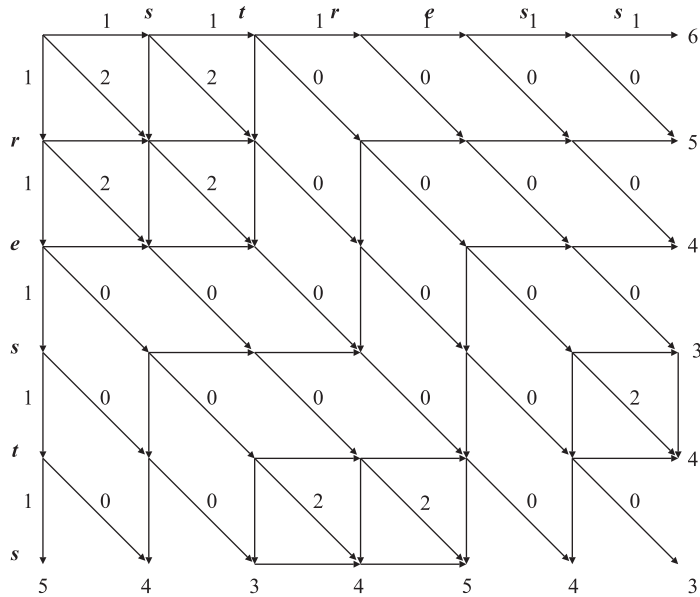


Рис. 9.2. Ориентированный граф зависимостей G для $d_L(rests, stress)$

Теперь покажем, как можно эффективно удалить из графа G вершины и дуги, которые не лежат на пути минимальной стоимости, проходящем от вершины $[0, 0]$

до вершины $[n_1, n_2]$. Покажем, что эти вычисления можно выполнить за время порядка $\Theta(nd)$.

Сначала заметим, что для каждой вершины $[i', j']$, лежащей на пути от вершины $[0, 0]$ до вершины $[i, j]$, выполняется неравенство $c[i, j] \geq c[i', j']$. Следовательно, на пути минимальной стоимости к вершине $[n_1, n_2]$ не может находиться вершины $[i', j']$, для которых $c[i', j'] > c[i, j]$, и поэтому все такие вершины и все пути, которые ведут только к таким вершинам, можно удалить из графа G . Справедливо и обратное утверждение: все вершины и все дуги, которые не удовлетворяют условию $c[i', j'] > c[i, j]$, должны лежать на пути к вершине $[n_1, n_2]$ и поэтому *должны* представлять некоторую последовательность операторов редактирования минимальной стоимости для расстояния $d(x_1, x_2)$. Как показано на рис. 9.3 для нашего примера, это простое наблюдение может значительно упростить граф G .

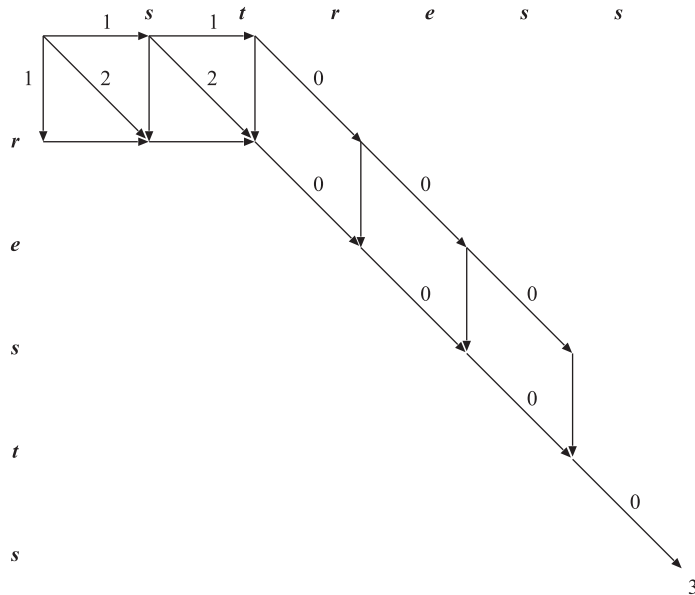


Рис. 9.3. Редуцированный граф зависимостей для $d_L(\text{rests}, \text{stress})$

Если внимательно рассмотреть рис. 9.3, то нетрудно заметить, что в редуцированном графе все операции вставки и удаления имеют стоимость, в точности равную 1. Также можно заметить, что в этом графе нет вершин, отстоящих более чем на $\lceil d/2 \rceil = 2$ “шага” от диагональной линии, которая начинается в вершине $[0, 1]$ и заканчивается в вершине $[n_1, n_2] = [5, 6]$. Более точно: в редуцированном графе не может быть более $\lfloor (n_1 - n_2 + d)/2 \rfloor$ вершин, лежащих ниже этой диагональной линии, и не может быть более такого же числа вершин, лежащих выше этой линии.

Поскольку сама диагональная линия состоит не более чем из n_1 дуг, то отсюда следует, что в редуцированном графе $O(n_1d)$ вершин (и столько же дуг). Поэтому задачу построения редуцированного графа можно решить за время порядка $O(n_1d)$ с использованием объема памяти такого же порядка.

В оставшейся части этого раздела предыдущие нестрогие рассуждения выразим в строгой математической форме, что также даст нам возможность оценить функцию расстояния посредством положительных значений стоимости операций вставки и удаления. Это в свою очередь приведет нас к желаемому алгоритму. Для упрощения выкладок далее в этом разделе будем предполагать, что $n_2 \geq n_1$. Для функции расстояния d , если она удовлетворяет условиям метрики, это предположение не является ограничением, поскольку в таком случае $d(\mathbf{x}_1, \mathbf{x}_2) = d(\mathbf{x}_2, \mathbf{x}_1)$ для любой пары строк \mathbf{x}_1 и \mathbf{x}_2 .

Пусть d^* — минимальная стоимость операций удаления и вставки, выполняемых по всем буквам алфавита A , т.е.

$$d^* = \min_{\lambda \in A} \{d(\lambda, \varepsilon), d(\varepsilon, \lambda)\} > 0.$$

Обозначим через $\Delta_{i,j}$ диагональ (идушую слева направо), на которой лежит вершина $[i, j]$. Думаю, что читатель не будет протестовать, если $\Delta_{i,j}$ припишем значение $j - i$. Таким образом, $\Delta_{i,j}$ может принимать целочисленные значения из интервала от $-n_1$ до n_2 .

Теперь рассмотрим ориентированный путь π на графе G от вершины $[i', j']$ до вершины $[i, j]$. Если $\Delta_{i,j} - \Delta_{i',j'} \geq 0$, тогда путь π содержит не менее $\Delta_{i,j} - \Delta_{i',j'}$ операций удаления. Если же $\Delta_{i,j} - \Delta_{i',j'} \leq 0$, тогда путь π содержит не менее $\Delta_{i',j'} - \Delta_{i,j}$ операций вставки. Поэтому, учитывая соотношения (9.8), можно записать

$$c[i, j] - c[i', j'] \geq |\Delta_{i,j} - \Delta_{i',j'}|d^* = |(j - i) - (j' - i')|d^*. \quad (9.9)$$

В частности, полагая $[i', j'] = [0, 0]$, получаем неравенство $c[i, j] \geq |j - i|d^*$, справедливое для любой вершины $[i, j]$, лежащей на пути от $[0, 0]$ до $[n_1, n_2]$. Поскольку $c[i, j] \leq c[n_1, n_2]$ для любой такой вершины $[i, j]$ и имеет место равенство $d = d(\mathbf{x}_1, \mathbf{x}_2) = c[n_1, n_2]$, то отсюда следует, что

$$d/d^* \geq |j - i|, \quad (9.10)$$

поэтому $|\Delta_{i,j}| \leq \lfloor d/d^* \rfloor$. В нашем примере $\mathbf{x}_1 = \text{rests}$, $\mathbf{x}_2 = \text{stress}$ и $d^* = 1$, поэтому первоначальное множество диагоналей $\{-3, -2, -1, 0, 1, 2, 3\}$ графа G должно быть уменьшено по крайней мере вдвое, чтобы получить редуцированный граф, показанный на рис. 9.3.

Однако не все так просто. Рассмотрим произвольный путь на графе G от вершины $[0, 0]$ до вершины $[n_1, n_2]$, проходящий через некоторую промежуточную

вершину $[i, j]$. Применяя дважды неравенство (9.9), сначала к пути $[i, j] - [n_1, n_2]$, а затем к пути $[0, 0] - [i, j]$, и учитывая, что $c[0, 0] = 0$, получим усиление неравенства (9.10):

$$d/d^* \geq |j - i| + |(n_2 - n_1) - (j - i)|. \quad (9.11)$$

Рассмотрим два возможных случая.

- $j \leq i$ (в этом случае $\Delta_{i,j} \leq 0$)

В данной ситуации неравенство (9.11) переписывается как $d/d^* \geq (n_2 - n_1) + 2(i - j)$, поскольку мы положили, что $n_2 \geq n_1$. Поскольку здесь $i - j$ — неотрицательное целое число, последнее неравенство можно переписать в виде

$$q \geq i - j \geq 0, \quad (9.12)$$

где

$$q = \lfloor (d/d^* - (n_2 - n_1)) / 2 \rfloor.$$

Неравенство (9.12) справедливо для всех промежуточных вершин $[i, j]$, лежащих в графе G на пути $[0, 0] - [n_1, n_2]$.

- $j \geq i$ (в этом случае $\Delta_{i,j} \geq 0$)

Здесь надо исследовать две ситуации, которые зависят от того, будет ли выполняться неравенство $n_2 - n_1 \leq j - i$. Если $n_2 - n_1 > j - i$, то неравенство (9.11) запишется в виде $d/d^* \geq n_2 - n_1$. Отсюда для величины q , определяемой формулой (9.13), получаем оценку

$$q \geq 0. \quad (9.13)$$

С другой стороны, для ситуации $n_2 - n_1 \leq j - i$ находим, что

$$d/d^* \geq 2(j - i) - (n_2 - n_1),$$

откуда следует, что

$$q + (n_2 - n_1) \geq j - i \geq 0. \quad (9.14)$$

Итак, в результате наших небольших исследований получено три неравенства (9.12), (9.14) и (9.15). Неравенство (9.14) в данный момент не представляет для нас интереса, поскольку оно выражает тривиальное условие, что значение d/d^* должно быть равно не менее половины разности длин строк x_1 и x_2 . Неравенства (9.12) и (9.15) объединим как

$$-q \leq j - i \leq q + (n_2 - n_1). \quad (9.15)$$

Еще раз подчеркнем, что эти неравенства справедливы для любой промежуточной вершины $[i, j]$, лежащей в графе G на произвольном пути от вершины $[0, 0]$ до

вершины $[n_1, n_2]$. Поскольку $j - i = \Delta_{i,j}$, то неравенства (9.16) можно рассматривать как ограничения на промежуточные вершины $[i, j]$, т.е. эти вершины должны находиться в графе G между диагоналями $-q$ и $q + (n_2 - n_1)$. Если применить неравенства (9.16) к нашему примеру, где $\mathbf{x}_1 = \text{rests}$, $\mathbf{x}_2 = \text{stress}$ и $d^* = 1$, то получим, что $q = 1$, а эти неравенства запишутся как $-1 \leq j - i \leq 2$, что соответствует редуцированному графу на рис. 9.3, где все вершины располагаются на диагоналях $\{-1, 0, 1, 2\}$.

Чтобы создать алгоритм на основе неравенств (9.16), предварительно рассмотрим такую задачу: для некоторого заданного положительного целого числа D надо определить, будет ли выполняться неравенство $d(\mathbf{x}_1, \mathbf{x}_2) \leq D$. Для решения этой задачи на основе формулы (9.13) можно сначала вычислить пробное значение

$$Q = \lfloor (D/d^* - (n_2 - n_1)) / 2 \rfloor$$

и затем найти минимальную сумму весов на путях от вершины $[0,0]$ до вершины $[n_1, n_2]$, которые лежат между диагоналями, задаваемые неравенствами (9.16):

$$-Q \leq \Delta_{i,j} \leq Q + (n_2 - n_1).$$

Если эта минимальная сумма, которая “по идее” должна быть равна $c[n_1, n_2]$, больше значения D , то это означает, что не все вершины $[i,j]$ включены в рассчитываемую сумму для того, чтобы получить минимальное значение $c[n_1, n_2]$. В этом случае необходимо включить в рассмотрение дополнительные диагонали и, соответственно, вершины, которые лежат на них. А для этого нужно увеличить значение Q и, следовательно, значение D . Если же вычисленная сумма не превышает D , то это означает, что $c[n_1, n_2] \leq D$. Эти вычисления реализуем в виде функции TRIAL_DISTANCE (пробное расстояние), которая во всех случаях возвращает значение $c[n_1, n_2]$, но это значение будет считаться “правильным” только тогда, когда $c[n_1, n_2] \leq D$.

▷ *Вычисление $c[n_1, n_2]$ на основе заданного значения D*

function TRIAL_DISTANCE(D)

▷ *На основе формулы (9.13) вычисляется Q*

▷ *и затем интервал диагоналей*

$$Q \leftarrow \lfloor (D/d^* - (n_2 - n_1)) / 2 \rfloor$$

▷ *Вычисление $c[n_1, n_2]$*

for $i \leftarrow 1$ **to** n_1 **do**

for $j \leftarrow \max\{1, i - Q\}$ **to** $\min\{n_1, i + Q + (n_2 - n_1)\}$ **do**

$$c[i, j] \leftarrow \min\{c[i - 1, j] + d(\mathbf{x}_1[i], \varepsilon),$$

$$c[i, j - 1] + d(\varepsilon, \mathbf{x}_2[j]),$$

$$c[i - 1, j - 1] + d(\mathbf{x}_1[i], \mathbf{x}_2[j])\}$$

return $c[n_1, n_2]$

Теперь все, что остается сделать алгоритму Укконена–Майерса (представлен ниже как алгоритм 9.5.1), так это инициализировать массив c и затем предложить подходящие значения D для функции TRIAL_DISTANCE. По причинам, которые будут понятны после обсуждения сложности алгоритма УМ, начальное значение D берется равным 1, а затем удваивается до тех пор, пока не будет выполнено неравенство $c[n_1, n_2] \leq D$. На основе обсуждений, проведенных выше для обоснования алгоритма, мы можем утверждать, что алгоритм УМ выполняется корректно.

Алгоритм 9.5.1 (Алгоритм Укконена–Майерса)

```

▷ Вычисление  $d(\mathbf{x}_1, \mathbf{x}_2) = c[n_1, n_2]$ 
▷ Присвоение значений элементам нулевых столбца и строки массива  $c$ 
 $c[0, 0] \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n_1$  do
     $c[i, 0] \leftarrow c[i - 1, 0] + d(\mathbf{x}_1[i], \varepsilon)$ 
for  $j \leftarrow 1$  to  $n_2$  do
     $c[0, j] \leftarrow c[0, j - 1] + d(\varepsilon, \mathbf{x}_2[j])$ 
▷ Инициализация  $D$ 
 $D \leftarrow 1$ 
▷ Удвоение значения  $D$  до тех пор,
▷ пока не будет достигнуто неравенство  $D \geq \text{TRIAL\_DISTANCE}(D)$ 
while  $D < \text{TRIAL\_DISTANCE}(D)$  do
     $D \leftarrow 2D$ 
output  $c[n_1, n_2]$ 

```

Для оценки сложности нового алгоритма рассмотрим сначала функцию TRIAL_DISTANCE, в частности двойной цикл **for**. Эти циклы выполняются не более $2Q + (n_2 - n_1) + 1$ раз, поскольку, вследствие формулы (9.13), справедливо неравенство

$$2Q + (n_2 - n_1) + 1 \leq D/d^* + 1.$$

Таким образом, общее время выполнения функции TRIAL_DISTANCE имеет порядок $O(n_1 D)$. В алгоритме 9.5.1 эта функция выполняется $h + 1$ раз для значений $D = 1, 2, \dots, 2^h$, где h — наименьшее положительное целое, такое, что $d = d(\mathbf{x}_1, \mathbf{x}_2) \leq 2^h$. Следовательно, общее время выполнения функции TRIAL_DISTANCE в алгоритме 9.5.1 составит

$$O(1 + 2 + \dots + 2^h)n_1 = O(2^{h+1}n_1) = O(4n_1d) = O(n_1d).$$

Наши выкладки доказывают следующую теорему.

Теорема 9.5.1. Для заданных строк $\mathbf{x}_1 = \mathbf{x}_1[1..n_1]$ и $\mathbf{x}_2 = \mathbf{x}_2[1..n_2]$, $n_2 \geq n_1$, алгоритм 9.5.1 корректно вычисляет $d = d(\mathbf{x}_1, \mathbf{x}_2)$ за время порядка $O(n_1d)$ с использованием памяти объемом порядка $\Theta(n_1n_2)$. ■

Как показано в упр. 9.5.2, объем памяти, необходимый для выполнения алгоритма УМ, можно уменьшить до величины порядка $O(n_1d)$. В упр. 9.5.3 предложено модифицировать этот алгоритм таким образом, чтобы и время его выполнения имело порядок $O(n_1d)$. Однако этого можно добиться только тогда, когда известна верхняя граница d для расстояния $d(x_1, x_2)$.

Для специальных случаев расстояний преобразования и Левенштейна в работе [144] предложена “пошаговая” версия алгоритма УМ, которая позволяет вычислить $d(x_1, \lambda x_2)$ и соответствующий массив стоимостей на основе известных расстояния $d(x_1, \lambda x_2)$ и его массива стоимостей за время порядка $O(d)$, где d — снова известная верхняя граница для расстояния, а λ — произвольная буква. Если эта верхняя граница не известна, то время выполнения такого алгоритма возрастает до величины порядка $O(n_1 + n_2)$. Поскольку $d(x_1, x_2\lambda)$ можно приближенно вычислить на основе массива стоимостей для $d(x_1, x_2)$, алгоритм [144] можно использовать для вычисления искомого расстояния последовательно, начиная с $d(\varepsilon, \varepsilon)$ путем добавления к начальным строкам новых букв, на каждом шаге по одной букве (при этом не имеет значения, в начало или конец строк добавляются новые буквы). Если верхняя граница d для расстояния определяется на основе подобных вычислений, то время выполнения данного алгоритма имеет порядок $O(d(n_1 + n_2))$, т.е. такой же порядок, что и алгоритм УМ. Однако пошаговый алгоритм позволяет эффективно решать другие задачи приближенного сравнения с паттерном, в частности, задачи P1 и P2, сформулированные в разделе 13.4. Алгоритм [144] в вычислительном плане весьма сложен, поэтому более подробно мы его рассматривать не будем.

Более простую реализацию пошаговая идея нашла в алгоритме [133], который позволяет вычислить $d(x_1, \lambda x_2)$ или $d(x_1, x_2\lambda)$ за время порядка $O(n_1 + n_2)$ независимо от того, известна ли верхняя граница d для расстояния. Этот алгоритм предпочтительнее других, если верхняя граница d неизвестна или велика, например $d \in \Theta(n_1)$.

Наконец отметим, что в алгоритме УМ, подобно алгоритму Вагнера–Фишера, наибольшую общую подпоследовательность $LCS(x_1, x_2)$ можно восстановить, двигаясь по графу G в обратном порядке от вершины $[n_1, n_2]$ до вершины $[0, 0]$.

Упражнения 9.5

1. В теореме 9.5.1 утверждается, что алгоритм УМ выполняется за время $O(n_1d)$, хотя только для инициализации массива c необходимо время порядка $\Omega(n_2)$. Чтобы развеять все сомнения в корректности этой теоремы, докажите, что $n_2 \in O(n_1d)$.
2. В теореме 9.5.1 утверждается, что для выполнения алгоритма УМ необходим объем памяти порядка $\Theta(n_1n_2)$. Предложите усовершенствования в этом

алгоритме, которые позволили бы сократить необходимый объем памяти до величины порядка $O(n_1d)$.

3. Покажите, как можно модифицировать алгоритм УМ для вычисления $d(x_1, x_2)$ в случае, когда известна верхняя граница d этого расстояния.

9.6 Заключение

В этой главе мы рассмотрели четыре основных алгоритма, которые для заданных строк $x_1 = x_1[1..n_1]$ и $x_2 = x_2[1..n_2]$ вычисляют или расстояние $d(x_1, x_2)$, или наибольшую общую подпоследовательность $LCS(x_1, x_2)$.

- Алгоритм Вагнера–Фишера вычисляет расстояние $d(x_1, x_2)$ для произвольной функции расстояния d за время порядка $\Theta(n_1n_2)$ с использованием памяти объемом такого же порядка. Если $d = d_L$ (расстояние Левенштейна), то этот алгоритм позволяет также найти $LCS(x_1, x_2)$ в качестве “сопутствующего продукта” вычисления расстояния.
- Алгоритм Хешберга (алгоритм НЗ) вычисляет $LCS(x_1, x_2)$ (следовательно, и расстояние Левенштейна $d_L(x_1, x_2)$) за время порядка $\Theta(n_1n_2)$, но использует память объемом порядка $\Theta(n)$, где $n = \min\{n_1, n_2\}$.
- Алгоритм Ханта–Шиманского вычисляет $LCS(x_1, x_2)$ за время, которое обычно имеет порядок $\Theta(n \log n)$, с использованием памяти объемом порядка $\Theta(n)$, где $n = \max\{n_1, n_2\}$. Временная сложность этого алгоритма превышает $\Theta(n \log n)$ только тогда, когда строки x_1 и x_2 имеют количество совпадающих букв, превышающее величину порядка $\Theta(n)$ (например, когда размер алфавита мал по отношению к длине строк).
- Алгоритм Укконена–Майерса вычисляет расстояние $d(x_1, x_2)$ за время порядка $\Theta(dn)$ с использованием памяти объемом порядка $\Theta(n_1n_2)$, где $n = \min\{n_1, n_2\}$. Как и в алгоритме Вагнера–Фишера, этот алгоритм позволяет эффективно вычислить и $LCS(x_1, x_2)$, если используется функция расстояния Левенштейна d_L . Этот алгоритм более предпочтителен, по сравнению с другими алгоритмами, если известна верхняя граница d для расстояния $d(x_1, x_2)$ и эта граница не слишком велика (например, $d \in O(\log n)$ или даже $d \in O(n/\log n)$). Такая ситуация часто встречается на практике: обычно вычисление расстояния между строками представляет интерес только тогда, когда это расстояние относительно мало по сравнению с длиной строк.

Для расстояния преобразования d_E , как показано в работе [189], если входные строки выбираются случайным образом, то алгоритм УМ в среднем выполняется за время порядка $O(n + k^2)$, а в работах [189, 149] доказано, что с использованием дерева суффиксов этот алгоритм в самом худшем случае выполняется за время порядка $O(n \log \alpha + k^2)$, где α — размер алфавита.

Пошаговые версии алгоритма УМ применимы для расстояний преобразования и Левенштейна [144, 133].

Существуют другие алгоритмы, решающие те же задачи, что и рассмотрены в этой главе, или близкие к ним. В работе [173] показано, как “четыре русских” метода [20] можно применить для разбиения массива стоимостей на перекрывающиеся квадратные подматрицы так, чтобы время вычисления расстояния имело порядок $O(n^2/\log n)$. Однако эти алгоритмы становятся предпочтительнее алгоритма УМ только для очень больших строк (когда $n \geq 300\,000$). В статье [128] предложены алгоритмы для “тяжелых” общих подпоследовательностей, а в статье [204] — для больших “возрастающих” подпоследовательностей. Несколько недавних алгоритмов [63, 64, 125] предлагают алгоритмы для вычисления LCS, построенные на основе использования бинарных векторов [77, 7, 24], которые обсуждаются в разделе 7.4 и главе 10. В работе [63] показано, как можно реализовать алгоритмы Хешберга и Ханта–Шиманского с использованием таких векторов. В общем случае эти алгоритмы выполняются за время порядка $O(n_1 n_2 / w)$, где w — длина машинного слова. Такие алгоритмы показали себя очень быстрыми на практике. Также быстрым является алгоритм Эллисона и Дикса (Allison and Dix) [7], который использует бинарные векторы для вычисления *длины* LCS.

ГЛАВА 10

Приближенное сравнение с паттерном

Газета состоит из одинакового количества слов независимо от того, содержит она новости или нет.

— Генри Филдинг (1707–1754)

Количество различных алгоритмов приближенного сравнения с паттернами, предложенных в период с 80-х годов и до настоящего времени, сравнимо с количеством алгоритмов точного сравнения с паттернами (см. главы 7 и 8). На момент написания этой книги можно было без труда заметить, что для построения всего множества всех наиболее эффективных и востребованных алгоритмов приближенного сравнения используются всего две идеи, с которыми мы уже имели удовольствие познакомиться.

1. В разделе 7.4 был описан алгоритм Демелки–Бейза–Ятса–Гоннета (сокращенно, алгоритм ДБГ) для вычисления точного совпадения с паттерном. В этом алгоритме использовались бинарный вектор $s = s[1..m]$, который описывает текущее “состояние” вычислений, и “кодирующая” бинарная матрица $t = t[1..α, 1..m]$, ставящая в соответствие элементам паттерна $p[1..m]$ буквы $1..α$ индексированного алфавита. В этом случае можно эффективно реализовать бинарные операции на бинарных массивах s и t для пересчета вектора состояний s с целью поиска совпадения с паттерном следующей буквы $x[i + 1]$, входящей в заданную текстовую строку $x = x[1..n]$.

В разделах 10.2–10.4 мы покажем три применения этой идеи для решения задачи приближенного сравнения с паттерном.

2. В разделе 9.1 был введен двухмерный массив стоимостей c , элементы которого определялись равенством $c[i, j] = d(x_1[1..i], x_2[1..j])$, т.е. равны расстоянию между префиксами строк $x_1 = x_1[1..n_1]$ и $x_2 = x_2[1..n_2]$. Отношения базовой рекурсии между смежными элементами массива c дали возможность в главе 9 построить два алгоритма прямого вычисления расстояния между строками (алгоритмы Вагнера–Фишера и Укконена–Майерса) и два алгоритма вычисления наибольшей общей подпоследовательности LCS для двух строк (алгоритмы Хешберга и Ханта–Шиманского).

В разделе 10.1 мы немного изменим определение элементов массива c , что позволит адаптировать алгоритм Вагнера–Фишера для приближенного сравнения с паттерном. В разделах 10.3 и 10.4 рассмотрим несколько вариаций реализации и применения массива c и его возможного “сотрудничества” с бинарными массивами алгоритма ДБГ. Наконец, в разделе 10.5 дан краткий обзор теоретических исследований временной сложности недавних алгоритмов сравнения с паттерном, использующих функции расстояния преобразования и Левенштейна.

Отметим, что в работе [191] дан современный обзор алгоритмов приближенного сравнения с паттерном, включая описание методов, которые не рассматриваются в данной книге. Кроме того, в этой работе приведены результаты экспериментов по оценке относительной эффективности алгоритмов.

10.1 Алгоритм для произвольной функции расстояния

Мы начали главу 9 с обсуждения массива стоимостей c и его свойств. Здесь мы пойдем таким же путем. Переопределим массив стоимостей $c = c[0..n_1, 0..n_2]$ из главы 9 как массив $c = c[0..n, 1..m]$ с элементами, которые определяются соотношением

$$c[i, j] = \min_{0 \leq i' \leq i} d(x[i'..i], p[1..j]), \quad (10.1)$$

справедливым для всех $i \in 0..n$, $j \in 1..m$. Здесь, как и ранее, d — произвольная функция расстояния, удовлетворяющая условиям (2.2)–(2.5) метрики, включая функцию взвешенного расстояния (раздел 2.2). Минимум по i' в выражении (10.1) необходим вследствие того, что, как показано в упражнении 10.1.1, здесь необязательно выполнение равенства $i - i' + 1 = j$, т.е. не исключаются случаи, когда $|x[i'..i]| \neq |p[1..j]|$. Таким образом, $c[i, j]$ теперь равняется расстоянию между префиксом паттерна $p[1..j]$ и “ближайшим” суффиксом подстроки $x[1..i]$.

В частности, $c[i, m]$ теперь равно минимальному расстоянию между паттерном p и подстроками строки x , которые оканчиваются в позиции i .

Из соотношения (10.1) следует, что начальные значения для элементов нулевой строки массива c здесь вычисляются точно так же, как и в разделе 9.1:

- для всех $j \in 1..m$ $c[0, j] = d(\varepsilon, p[1..j])$ — суммарная стоимость операций вставки для префикса $p[1..j]$ паттерна p .

Однако начальные значения для первого столбца массива c вычисляются способом, отличным от того, как это делается в разделе 9.1, поскольку здесь фактически находится минимум расстояний между буквой $p[1]$ и подстроками $x[i'..i]$, а не вычисляется расстояние между $p[1]$ и единственным префиксом $x[1..i]$. Конечно, если $x[i] = p[1]$, тогда мы положим

$$c[i, 1] \leftarrow 0, \tag{10.2}$$

но если $x[i] \neq p[1]$, то в этом случае необходимо применить формулу

$$c[i, 1] \leftarrow \min\{c[i-1, 1] + d(x[i], \varepsilon), d(x[i], p[1])\}. \tag{10.3}$$

Отметим, что различия в инициализации массивов стоимостей c здесь и в разделе 9.1 — это *единственное* различие в вычислении этих массивов. Если значения в нулевой строке и первом столбце массива c вычислены корректно, тогда, поскольку для произвольной функции расстояния операции вставки и удаления можно упорядочить любым способом (упражнение 2.2.14), справедливы утверждения леммы 9.1.1, задающие отношения базисной рекурсии для вычисления элементов массива c . Поэтому, как и ранее, для всех $i \in 1..n$ и $j \in 2..m$ элементы $c[i, j]$ вычисляются по формуле

$$\begin{aligned} c[i, j] \leftarrow \min\{ & c[i-1, j] + d(x[i], \varepsilon), \\ & c[i, j-1] + d(\varepsilon, p[j]), \\ & c[i-1, j-1] + d(x[i], p[j])\}. \end{aligned} \tag{10.4}$$

Используем пример из главы 9, где вычислялось расстояние между строками *rests* и *stress*, для того чтобы показать различие между массивами c , вычисленным в главе 9 и вычисленным по формулам данной главы. Обозначим через c_{dist} массив стоимостей, вычисленный на основе значений расстояний (как в главе 9), а через c_{best} — массив, вычисленный на основе значений наилучшего (best) приближенного совпадения с паттерном всех суффиксов, которые оканчиваются в позиции i . Если для вычисления массива стоимостей используется расстояние Левенштейна, то получим массив c_{dist} , который представим в табл. 10.1. Массив c_{best} для $p = \textit{stress}$ и $x = \textit{rests}$ показан в табл. 10.2. Отметим, что все различия в этих массивах определяются только различными значениями в столбце, соответствующего $j = 1$.

Таблица 10.1. Массив c_{dist} для $d_L(\text{rests}, \text{stress})$

j	0	1	2	3	4	5	6
i	ε	s	t	r	e	s	s
0	ε	0	1	2	3	4	5
1	r	1	2	3	2	3	4
2	e	2	3	4	3	2	3
3	s	3	2	3	4	3	2
4	t	4	3	2	3	4	3
5	s	5	4	3	4	5	4

Таблица 10.2. Массив c_{best} для $p = \text{stress}$ и $x = \text{rests}$

j	1	2	3	4	5	6
i	s	t	r	e	s	s
0	ε	1	2	3	4	5
1	r	2	3	2	3	4
2	e	2	3	3	2	3
3	s	0	1	2	3	2
4	t	1	0	1	2	3
5	s	0	1	2	3	2

Для любого неотрицательного целого числа k по столбцу j массива c_{best} можно определить такую позицию i , для которой существует подстрока $u = x[i'..i]$, что $d(u, p[1..j]) \leq k$. В этом случае будем говорить, что в позиции i строки x **имеет место k -приближенное совпадение** с префиксом $p[1..j]$ паттерна p . Например, взяв $j = 5$ и $k = 2$ из табл. 10.2, находим, что в позициях $i = 3$ и $i = 5$ имеют место 2-приближенные совпадения с $p[1..5] = \text{stres}$. Аналогично для $j = 3$ и $k = 1$ в позиции $i = 4$ имеет место 1-приближенное совпадение с $p[1..3] = \text{str}$. Отметим, что совпадающие подстроки u не обязательно будут единственными.

Очевидно, что алгоритм Вагнера–Фишера (см. раздел 9.2), если в нем для инициализации массива стоимостей применить формулы (10.2) и (10.3), можно использовать для решения задачи приближенного сравнения с паттерном. Но, как отмечено выше, длина $i - i' + 1$ совпадающей подстроки u может не совпадать с длиной m паттерна. Поэтому, если необходимо определить хотя бы одну подстроку $u = x[i'..i]$, частично совпадающую с паттерном, то также необходимо указать и значение i' . Это нетрудно сделать, если знать последовательность перемещений v , введенных в разделе 9.2. Эти перемещения определяют последо-

вательность операторов редактирования, необходимых для преобразования подстроки u в паттерн p (или, что то же самое, для преобразования $p \rightarrow u$). Зная перемещения, можно вычислить значение i' . Как показано в упражнении 10.1.3, количество перемещений не превышает $2|p| = 2m$, тогда как количество k -приближенных совпадений в любой позиции i не превышает n . Поэтому общее время вычисления всех возможных значений i' имеет порядок $O(mn)$, т.е. не превышает время, необходимое для вычисления массива c . Следовательно, справедлива такая теорема.

Теорема 10.1.1. Для заданных текстовой строки $x[1..n]$, паттерна $p[1..m]$ и неотрицательного целого числа k модифицированный алгоритм Вагнера–Фишера вычисляет все k -приближенные совпадения строки x с паттерном p за время порядка $O(mn)$ с использованием памяти, объем которой имеет такой же порядок. ■

Еще раз подчеркнем, что модифицированный алгоритм Вагнера–Фишера вычисляет *все* k -приближенные совпадения строки x с любым префиксом $p[1..j]$ ($1 \leq j \leq m$) паттерна p .

Конечно, алгоритм, описанный в этом разделе, очень затратный как по используемому времени, так и по используемому объему памяти. Однако, как было показано в упражнении 9.3.3 и как будет показано в упражнении 10.1.4, объем используемой памяти можно легко уменьшить, если заметить, что в массиве c значения в строке i зависят только от значений в строке $i - 1$. Кроме того, данный алгоритм имеет одно несомненное достоинство: он применим для работы с любыми функциями расстояния, которые удовлетворяют аксиомам (2.2)–(2.5) метрики. В следующем разделе мы покажем, как путем наложения различного рода ограничений на область применения алгоритмов приближенного сравнения можно значительно уменьшить их временную сложность.

Упражнения 10.1

1. Для расстояния Левенштейна d_L приведите пример паттерна $p[1..m]$ и текстовой строки x , таких, что в формуле (10.1) элемент $c[i, m]$ вычислялся бы при значении i' , которое удовлетворяет соотношениям

а) $i + i' - 1 < m$;

б) $i + i' - 1 = m$;

в) $i + i' - 1 > m$.

Затем покажите, что все эти три соотношения могут иметь место тогда, когда стоимость некоторой операции подстановки превышает стоимость операций вставки и удаления.

2. Докажите, что для расстояния преобразования d_E без потери общности в выражении (10.1) можно считать, что $i + i' - 1 \leq m$.

3. Покажите, что максимальное число перемещений v , которые преобразуют совпадающую подстроку u в паттерн p , равно $2|p|$.
4. Запишите алгоритм (модификацию алгоритма Вагнера–Фишера, соответствующую теореме 10.1.1), который вычислял бы все целые i и соответствующие целые i' , такие, что $d(x[i'..i], p) \leq k$. Покажите, что для реализации этого алгоритма достаточно памяти объемом порядка $\Theta(\min\{mn, m^2\})$. Далее покажите, что если не требуется выводить значения i' , то объем памяти можно уменьшить до величины порядка $\Theta(m)$.

10.2 Алгоритм, вычисляющий несовпадения

В этом разделе мы покажем, как уже известный нам алгоритм ДБГ из раздела 7.4 можно легко модифицировать так, чтобы он вычислял k -приближенные совпадения, рассчитанные на основе расстояния Хемминга (см. раздел 2.2). В действительности, это только первое из нескольких возможных применений идеи алгоритма ДБГ к различным типам задач приближенного сравнения с паттерном. Здесь представим идею этого алгоритма в более общем контексте, чем мы делали это ранее. Вместе с тем в этом разделе мы будем применять только расстояние Хемминга, поэтому задачу, решаемую с помощью алгоритма данного раздела, назовем *задачей k -несовпадений*.

Далее будем считать, что задано некоторое целое число $k \geq 0$, а также заданы паттерн $p = p[1..m]$ и текстовая строка $x = x[1..n]$, которые определены на алфавите A . Как и в алгоритме ДБГ, здесь мы требуем, чтобы алфавит A был конечен, индексирован и известен заранее. В этом случае можно считать (см. упражнение 4.1.3), что $A = \{1, 2, \dots, \alpha\}$.

В исходном алгоритме ДБГ использовался *вектор состояний* $s = s[1..m]$, в котором элемент $s[j]$ указывал, есть ли совпадение текущей подстроки $x[i - j + 1..i]$ строки x с префиксом $p[1..j]$ паттерна p . Поскольку в разделе 7.4 мы рассматривали точное совпадение с паттерном, поэтому элементы $s[j]$ принимали только два значения: $s[j] = 0$, если имело место совпадение, и $s[j] = 1$ в случае несовпадения. Здесь мы рассматриваем k -приближенное сравнение с паттерном на основе расстояния Хемминга, поэтому естественно сейчас положить, что элементы $s[j]$ будут равны *количеству* несовпадающих позиций при сравнении $p[1..j]$ и текущей подстрокой строки x . Поскольку число несовпадающих позиций может быть от 0 до m , то для хранения этого числа необходимо не менее $\lceil \log_2(m + 1) \rceil$ бит на каждый элемент вектора s , а не один бит, как в алгоритме ДБГ. Обозначим через B число $\lceil \log_2(m + 1) \rceil$ — количество битов, требуемых для хранения каждого элемента вектора s (ниже мы покажем, как можно уменьшить это число).

Далее вместо формулы (7.13) для вычисления $s[j]$ в исходном алгоритме ДБГ будем использовать формулу

$$s[j] = d_H(x[i - j + 1..i], p[1..j]). \quad (10.5)$$

Эта формула применима для любого $j \in 1..m$, тогда как i последовательно принимает значения $1, 2, \dots, n$.

Напомним, что в алгоритме ДБГ также использовался массив $t = t[1..\alpha, 1..m]$, где для каждой буквы $h \in 1..\alpha$ и для каждой позиции $j \in 1..m$

$$t[h, j] = 0, \text{ если } p[j] = h, \text{ и } t[h, j] = 1 \text{ в противном случае.} \quad (10.6)$$

Здесь мы также используем массив t , но теперь каждый элемент этого массива будет не одиночным битом, а будет содержать $B \geq 1$ бит. При этом $t[h, j]$ будет также равняться 0 или 1, но перед ними будет стоять $B - 1$ значащих нулей. Эти дополнительные нули необходимы для того, чтобы можно было просто и эффективно вычислить сумму $s + t[h, 1..m]$, аналог операции логического сложения $s \vee t[h, 1..m]$ в алгоритме ДБГ. Как предложено показать в упражнении 10.2.2, массив t можно эффективно вычислить на предварительном этапе алгоритма.

Используя новые формулы (10.5) и (10.6) для вычисления вектора s и массива t , можем сформулировать и доказать результат, аналогичный лемме 7.4.1 (здесь операция логического сложения заменяется на обычную операцию сложения).

Лемма 10.2.1. Пусть $s^{(i)}$ — вектор состояний, соответствующий позиции i строки x . Начальные состояния определяются как $s^{(0)}[0..m] = 0^B 0^B \dots 0^B$. Тогда для всех $i \in 1..n$ и $j \in 1..m$ справедлива формула

$$s^{(i)}[j] = s^{(i-1)}[j - 1] + t[x[i], j]. \quad (10.7)$$

■

Доказательство в точности повторяет схему доказательства леммы 7.4.1, оставим его для упражнения 10.2.1. ■

Отметим, что возможны другие леммы, подобные леммам 10.2.1 и 7.4.1, если предложить оператор (+, \vee или какой-либо другой), который на основе значений для позиции $j - 1$ предыдущего состояния корректно пересчитывал бы значения для позиции j текущего состояния. Другими словами, леммы 10.2.1 и 7.4.1 можно рассматривать как частные случаи более общей леммы, где определен некий оператор (назовем его оператором *op*), удовлетворяющий этому условию. В дальнейшем при использовании идеи алгоритма ДБГ мы будем учитывать, что возможна такая общая лемма.

Подобно тому, как в разделе 7.4 на основе леммы 7.4.1 был создан алгоритм ДБГ, здесь на основе леммы 10.4.1 создадим аналогичный алгоритм для при-

ближенного сравнения с паттерном и назовем его алгоритмом ПДБГ (“приближенный” алгоритм ДБГ). В этом новом алгоритме на каждом шаге выполняются вычисления

$$s \leftarrow \text{rightshift}(s, B) + t[x[i], 1..m], \quad (10.8)$$

где

- сначала вектор состояний s с помощью процедуры *rightshift* сдвигается вправо на B бит (с учетом B нулей, входящих в значение $s[0]$), тем самым осуществляется переход от состояния $j - 1$ к состоянию j ($j \in 1..m$);
- затем для всех m позиций в сдвинутом векторе s выполняется пересчет путем прибавления соответствующих m элементов массива t , определяемых текущей буквой $x[i]$ строки x .

Реализация алгоритма ПДБГ очень похожа на реализацию алгоритма ДБГ, поэтому оставим ее для упражнения 10.2.3.

Чтобы оценить сложность нового алгоритма, сначала заметим, что для хранения вектора $s[1..m]$ требуется mB бит, а массива $t - \alpha mB$ бит. Так же как и в алгоритме ДБГ, предположим, что машинное слово имеет длину w бит. Тогда для хранения вектора $s[1..m]$ требуется $\lceil mB/w \rceil$ машинных слов, а массива $t - \alpha \lceil mB/w \rceil$ машинных слов. Далее предположим, что арифметические операции (такие, как сдвиг и сложение) на отдельных машинных словах могут выполняться за константное время. Тогда для каждого значения i операцию (10.8) можно выполнить за время порядка $\Theta(\lceil mB/w \rceil)$. На этом основании можем сформулировать аналог теоремы 7.4.3.

Теорема 10.2.2. Для заданных текстовой строки $x[1..n]$, паттерна $p[1..m]$ и неотрицательного целого числа k алгоритм ПДБГ вычисляет все k -приближенные совпадения строки x с паттерном p за время порядка $\Theta(\lceil mB/w \rceil n)$ с использованием памяти объемом порядка $\Theta(\alpha \lceil mB/w \rceil)$. ■

Поскольку $w -$ константа, а $B \in \Theta(\log m)$, то отсюда делаем заключение, что алгоритм ПДБГ имеет время выполнения порядка $\Theta(mn \log m)$. Но это еще не та эффективность, которую хотелось бы иметь. Конечно, если длина паттерна m мала настолько, что величина $\lceil mB/w \rceil$ ограничена небольшим целым числом, то в этом случае можно утверждать, что алгоритм ПДБГ на практике будет иметь время выполнения порядка $\Theta(n)$. Но чтобы быть уверенным в том, что величина mB/w действительно невелика, следует уменьшить B до величины порядка $\Theta(\log k)$.

Сделаем два почти очевидных наблюдения. Во-первых, мы можем считать, что всегда $k < m$, поскольку для $k \geq m$ задача k -приближенного сравнения тривиальна. Во-вторых, значение $s[j]$ нет необходимости вычислять точно, если это значение превышает k . Действительно, когда $s[j]$ впервые превышает k , то достаточно просто зафиксировать этот факт, так как несовпадение с паттерном $p[1..m]$ уже имеет место и не важно, будут ли зафиксированы еще другие совпадения

или несовпадения с суффиксами $p[j + 1..m]$ паттерна. Эти простые наблюдения позволяют утверждать, что для хранения необходимой информации достаточно $\lceil \log_2(k + 1) \rceil$ бит плюс один бит переполнения (фиксирующий, что значение $s[j]$ превышает k) для каждого $m + 1$ элемента вектора s и, вследствие этого, также для каждого элемента массива t . Модифицированный таким образом алгоритм ПДБГ удовлетворяет теореме 10.2.2 с константой $B = \lceil \log_2(k + 1) \rceil + 1$ и асимптотически имеет время выполнения порядка $\Theta(mn \log k)$. Как отмечалось в разделе 7.4, фактически $w \in \Omega(\log n)$, поэтому порядок времени выполнения алгоритма можно записать как $O(\lceil m \log k / \log n \rceil n)$. Отметим, что в этом случае для $k = 1$, $B = 1$ алгоритм ПДБГ ведет себя так же, как и алгоритм БДГ (что и ожидалось).

Для объяснения вычислительного механизма модифицированного алгоритма ПДБГ рассмотрим пример, взятый из работы [24] и представленный в табл. 10.3. Здесь рассматривается задача нахождения 2-приближенных совпадений паттерна $p = ababc$ со строкой $x = abdabababc$ на основе расстояния Хемминга. Поскольку здесь $k = 2$, необходимо $B = 3$ бит для каждого элемента вектора s (предполагая, что значения этих элементов не превысят 3) плюс один бит переполнения (который положим равным 1, если количество несовпадений превысит 3). Для удобства вычислений в момент времени, когда бит переполнения становится равным 1, его поместим в отдельный вектор, который обозначим s' . Вертикальные стрелки в табл. 10.3 указывают на две позиции в строке x , где достигается 2-приближенное совпадение с паттерном; точнее, это две позиции, где выполняется условие $s[5] + 2^2 s'[5] \leq 2$.

Таблица 10.3. Сравнение паттерна $p = ababc$ со строкой $x = abdabababc$, $k = 2$, с помощью модифицированного алгоритма ПДБГ

x	a	b	d	a	b	a	b	a	b	c
$t[x[i]]$	01011	10101	11111	01011	10101	01011	10101	01011	10101	11110
s	01011	10202	12131	02220	10323	02003	10301	02001	10301	12100
s'	01111	00111	00011	00001	00000	00010	00001	00010	00001	00010
								↑		↑

Бит переполнения в элементе $s[j]$ используется в следующей процедуре для проверки, превышает ли $s[j]$ значения k . Если бит переполнения равен 0, то он просто игнорируется. Если же он равен 1, то его надо переписать в вектор s' в ту же позицию j , а затем в векторе s его следует вернуть в состояние 0. После выполнения такой процедуры величина $s[j] + 2^{B-1} s'[j]$ будет нижней границей количества несовпадений в позиции j паттерна p (т.е. количество несовпадений в позиции j не меньше указанной величины). Описанную процедуру можно реализовать с помощью оператора логического умножения $op = \text{AND}$ с использованием шаблона $0^B 10^B 1 \dots 0^B 1$ для извлечения битов переполнения из вектора s и инверсного шаблона $1^B 01^B 0 \dots 1^B 0$ для возврата битов переполнения в нуле-

вое состояние. (Напомним, что бинарный логический оператор AND возвращает значение 1 только в том случае, если оба его операнда имеют значение 1.)

Из табл. 10.3 видно, что векторы s и s' на каждом шаге алгоритма должны сдвигаться вправо одновременно, что обеспечит синхронизацию позиций j в этих векторах. Это позволит проверить выполнение неравенства

$$s[m] + 2^{B-1} s'[m] \leq k,$$

для того чтобы определить, имеет ли место k несовпадений паттерна $p[1..m]$ в текущей позиции строки x .

Детали реализации модифицированного алгоритма ПДБГ оставим для упражнения 10.2.3. Отметим только, что объем памяти, необходимый для хранения векторов s и s' и соответствующих шаблонов, ограничен величиной порядка $O(\lceil mB/w \rceil)$, что и требуется для выполнения условий теоремы 10.2.2.

Рассмотрим пару небольших примеров, чтобы лучше разобраться в пространственных и временных требованиях алгоритма ПДБГ. Предположим, что длина машинного слова $w = 32$. Для аппроксимационной константы $k \leq 7$ значение B равно 4. Поэтому для любого паттерна длиной $m \leq 8$ величина $\lceil mB/w \rceil = 1$. В этом случае для хранения одного элемента вектора $s[1..8]$ (следовательно, и вектора $s'[1..8]$) достаточно одного машинного слова, а алгоритм ПДБГ будет выполняться за время порядка $\Theta(n)$. Если же $k \leq 3$, тогда значение $B = 3$, и для любого паттерна длиной $m \leq 10$ величина $\lceil mB/w \rceil$ снова равна 1. Продолжение подобного анализа оставим для упражнения 10.2.4.

Другие алгоритмы

Первым практическим алгоритмом для решения задачи k -несовпадений, по-видимому, был алгоритм, предложенный Ландау и Вишкиным (Landau, Vishkin) [146, 147]. Этот алгоритм использует двухмерный массив “несогласованности паттерна”, который вычисляется на предварительном этапе алгоритма. Затем алгоритм вычисляет массив “несогласованности текста”, на основе которого определяется k -приближенное совпадение паттерна и заданной текстовой строки. Включая предварительный этап, этот алгоритм выполняется за время порядка $O(k(n + m \log m))$. Как показано в упражнении 10.2.5, время выполнения алгоритма Ландау–Вишкина в самом худшем случае превышает верхнюю границу времени выполнения алгоритма ПДБГ на достаточно широком множестве входных данных. Гелилом и Джейнкарло (Galil, Giancarlo) [95] время выполнения алгоритма Ландау–Вишкина в самом худшем случае было улучшено до величины порядка $O(kn + m \log m)$ путем использования дерева суффиксов. Однако и этот алгоритм показал на практике медленную работу [24]. Другой алгоритм, использующий бинарные операции, описан в работе [23]. Этот алгоритм подсчитывает длины частичных совпадений с паттерном p для каждой позиции i строки x . Если в алфавите размером α все буквы в паттерне p и строке x могут появиться с одинаковой

вероятностью, то *математическое ожидание* времени выполнения этого алгоритма составит величину порядка $\Theta(mn/\alpha)$, что является некоторым улучшением по сравнению с алгоритмом ПДБГ в случае больших значений α (например, если $\alpha \in \Omega(m)$ или $\alpha \in \Omega(n)$).

Подход, использованный при построении алгоритма Бойера–Мура (раздел 7.2), также можно использовать для решения задачи k -несовпадений. В работе [218] описана обобщенная версия алгоритма Бойера–Мура–Хоспула (раздел 8.2), которая вычисляет все k -приближенные совпадения паттерна p и строки x за среднее время порядка $O\left(k\left(\frac{1}{m-k} + \frac{k}{\alpha}\right)n\right)$. Таким образом, этот алгоритм может конкурировать с алгоритмом ПДБГ в случае, если значение k мало по сравнению с m и α и значение m велико по сравнению с длиной машинного слова w . Представляет большой интерес предложенное Эль-Мабруком и Крочемором (El-Mabrouk, Crochemore) [82] объединение алгоритма ПДБГ и алгоритма Бойера–Мура для вычисления точного совпадения с паттерном. На практике для больших значений m и малых значений k этот гибридный алгоритм выполняется немного быстрее, чем алгоритм ПДБГ.

Даже на основе такого краткого знакомства с другими возможными алгоритмами для решения задачи k -несовпадений, можно сделать вывод, что на сегодняшний день алгоритм ПДБГ и его варианты, по-видимому, показывают наилучшие результаты в решении этой задачи.

Упражнения 10.2

1. Докажите лемму 10.2.1.
2. Запишите алгоритм, который вычислял бы массив t за время порядка $\Theta(\alpha m B/w)$.
3. Запишите исходный алгоритм ПДБГ (когда $B = \lceil \log_2(m+1) \rceil$) и его модификацию (когда $B = \lceil \log_2(k+1) \rceil + 1$).
4. Составьте таблицу, где бы для $k = 0, 3, 7, 15, 31, 63$ и $\lceil mB/w \rceil \in 1..10$ приводилась максимальная возможная длина паттерна m , с которым при этих условиях может работать алгоритм ПДБГ. Рассмотрите два варианта таблицы: когда длина машинного слова $w = 31$ и когда $w = 64$.
5. Для положительных целых чисел $n \geq m \geq k \geq 2$ докажите, что при выполнении неравенства $m \log_2 m/n \leq 2$ также справедливо неравенство

$$mn \log_2 k / \log_2 n < k(n + m \log_2 m).$$

Используя доказанное неравенство, сравните асимптотическую эффективность алгоритмов ПДБГ и Ландау–Вишкина.

10.3 Алгоритмы, вычисляющие разности

В этом разделе рассмотрим алгоритмы приближенного сравнения с паттерном, основанные на расстоянии преобразования. Назовем задачи, решаемые такими алгоритмами, *задачами k -разностей*. Напомним, что для расстояния преобразования стоимости всех однобуквенных операций удаления и вставки и операций подстановки положительны и равны между собой, т.е.

$$d_E(\varepsilon, \lambda) = d_E(\lambda, \varepsilon) = d_E(\lambda, \mu) > 0 \tag{10.9}$$

для любой буквы λ из алфавита A и для любых букв $\lambda \neq \mu$. Конечно, если $\lambda = \mu$, то $d_E(\lambda, \lambda) = 0$. Поскольку стоимость всех однобуквенных операций одинаковая, то без потери общности мы можем положить эту стоимость равной 1. На основании формул (10.2), (10.3) и леммы 9.1.1 нетрудно показать, что для расстояния преобразования массив стоимостей c можно представить в виде табл. 10.4. Здесь элементы δ_{i1} в столбце 1 вычисляются следующим образом (лемма 9.1.1):

$$\delta_{i1} = 0, \text{ если } x[i] = p[1], \text{ и } \delta_{i1} = 1, \text{ если } x[i] \neq p[1]. \tag{10.10}$$

Таблица 10.4. Массив стоимостей для расстояния преобразования

	j	0	1	2	...	m
i		ε	$p[1]$	$p[2]$...	$p[m]$
0	ε	0	1	2	...	m
1	$x[1]$	0	δ_{11}
2	$x[2]$	0	δ_{21}
	
n	$x[n]$	0	δ_{n1}

Теперь докажем важный результат, который первоначально в более общей форме был установлен в работе [173], а затем в форме, представленной ниже, приведен в статье [226].

Лемма 10.3.1. Для любого массива стоимостей $c = c[0..n, 0..m]$, вычисленного с использованием расстояния преобразования, справедливы соотношения ($i \in 1..n, j \in 1..m$):

- а) $c[i, j] - c[i, j - 1] \in \{-1, 0, 1\}$;
- б) $c[i, j] - c[i - 1, j] \in \{-1, 0, 1\}$;
- в) $c[i, j] - c[i - 1, j - 1] \in \{0, 1\}$.

Доказательство. Сначала заметим, что вследствие леммы 9.1.1 и равенств (10.9) значение любого элемента $c[i, j]$ массива c не может превышать более чем на 1 значения элементов $c[i - 1, j]$, $c[i, j - 1]$ и $c[i - 1, j - 1]$. Поскольку все значения массива c являются целыми числами, то остается только доказать, что разности в утверждении леммы не могут быть меньше -1 (соотношения a) и b) или 0 (соотношение c)).

Докажем соотношение a) по индукции, заметив в табл. 10.4, что это соотношение выполняется для $i = 1$. Предположим, что оно справедливо для некоторого $i \geq 1$ и для всех j . Теперь предположим, что для $i + 1$ и для некоторого наименьшего j соотношение a не выполняется, т.е. имеет место неравенство

$$c[i + 1, j] - c[i + 1, j + 1] < -1.$$

Но тогда $c[i + 1, j - 1] > c[i + 1, j] + 1$ и, следовательно,

$$\begin{aligned} c[i + 1, j] &= \min\{c[i, j], c[i, j - 1], c[i + 1, j - 1]\} + 1 = \\ &= \min\{c[i, j], c[i, j - 1]\} + 1. \end{aligned}$$

Отсюда вытекает, что

$$c[i + 1, j - 1] > \min\{c[i, j], c[i, j - 1]\} + 2.$$

С другой стороны, мы знаем, что для расстояния преобразования должно выполняться неравенство

$$c[i, j - 1] \geq c[i + 1, j - 1] - 1 \tag{10.11}$$

и поэтому должно выполняться неравенство

$$c[i + 1, j - 1] \geq c[i, j] + 2. \tag{10.12}$$

Из неравенств (10.11) и (10.12) вытекает, что

$$c[i - 1, j] - c[i, j] \geq c[i + 1, j - 1] - 1 - c[i, j] > 1.$$

Последнее неравенство противоречит лемме 9.1.1 для расстояния преобразования. Это противоречие и доказывает, что не существует таких значений j , для которых соотношение a не выполняется.

Соотношение b следует из соотношения a вследствие симметрии формулы (10.4) относительно i и j . Доказательство соотношения c оставим для упражнения 10.3.1. ■

Вооружившись этой леммой, теперь можем представить два алгоритма, решающих задачу k -разностей. Первый алгоритм предложен Укконеном [226]. В этом алгоритме скомбинированы результаты леммы 10.3.1 и подход, реализованный

в алгоритме Укконена–Майерса (раздел 9.5). Он выполняется в среднем за время порядка $O(kn)$. Второй алгоритм — это новый алгоритм, предложенный Майерсом [187]. Этот алгоритм основан на той же идее, что и алгоритм ДБГ (раздел 7.4), и использует в вычислениях бинарные операции. Его время выполнения в самом худшем случае имеет порядок $O(mn/w)$, т.е. не зависит от k . Идеи алгоритма Укконена можно применить для модификации алгоритма Майерса, после чего его время выполнения в среднем будет иметь порядок $O(kn/w)$. Алгоритм Ву–Менбера, третий алгоритм для решения задачи k -разностей, описан в разделе 10.4. Хотя этот алгоритм требует времени выполнения в самом худшем случае порядка $O(kmn/w)$, его область применения значительно шире, чем области применения алгоритмов Укконена и Майерса.

10.3.1 Алгоритм Укконена

Алгоритм Укконена–Майерса, используемый для вычисления расстояния между строками (раздел 9.5), в самом худшем случае выполняется за время порядка $O(dn)$, где d — фактическое расстояние между строками. Такой результат достигается путем уменьшения количества элементов массива c , участвующих в вычислениях (учитываются только те элементы массива c , которые располагаются на определенных диагоналях этого массива, а число таких диагоналей имеет порядок $O(d)$, подробнее см. раздел 9.5). Подобную стратегию использует алгоритм Укконена решения задачи k -разностей, причем время выполнения этого алгоритма в среднем составляет величину порядка $O(kn)$. Основой данного алгоритма является соотношение \wp леммы 10.3.1, которое фактически показывает, что необходимые для вычислений диагональные элементы массива c (вычисленного на основе расстояния преобразования) образуют неубывающую последовательность (не по значениям, а по местоположению). Если обозначить через j' самую правую позицию в строке i массива c , такую, что $c[i, j'] \leq k$, тогда для всех элементов в строке $i + 1$ в позициях $j = j' + 2, j' + 3, \dots, m$ справедливо неравенство $c[i + 1, j] > k$. Поэтому $m - j' - 1$ самых правых элементов, стоящих в строке $i + 1$ массива c , можно исключить из вычислений, поскольку они не могут определять k -приближенное совпадение с паттерном p . В общем случае получаем такое правило: во всех строках массива c , номер которых превышает i , можно не рассматривать элементы, стоящие на диагоналях, располагающихся правее диагонали $\Delta_{i, j'}$. Эта простая стратегия составила основу следующего алгоритма.

Алгоритм 10.3.1. (Алгоритм Укконена)

- ▷ *Вычисление всех k -приближенных совпадений*
- ▷ *паттерна $p[1..m]$ со строкой $x[1..n]$*
- инициализация элементов в нулевых строке и столбце массива c (см. табл. 10.4)
- ▷ *В строке 0 число k находится в позиции $j' = k$*

```

j' ← k
for i ← 1 to n do
    j'' ← 0    ▷ для случая, когда k = 0
    for j ← 1 to min{j' + 1, m} do
        вычисление c[i, j] на основе формулы (10.4)
        if c[i, j] ≤ k then j'' ← j
    ▷ Вывод результата
    if j'' = m then output i
    ▷ Сохранение самой правой позиции j', такой, что c[i, j'] ≤ k
    j' ← j''

```

Очевидно, что временная сложность алгоритма Укконена пропорциональна сумме n значений переменной j' , вычисленных во время выполнения алгоритма. Конечно, в самом худшем случае, когда все j' равны m (в этом случае в каждой строке $i \in 1..n$ $c[i, m] \leq k$), алгоритм Укконена выполняется за время порядка $\Theta(mn)$. Но даже если, например, условие $c[i, m/2] \leq k$ выполняется для числа строк порядка $\Theta(n)$, то и в этом случае алгоритм выполняется за время порядка $\Theta(mn)$.

Чтобы определить временную сложность алгоритма Укконена в среднем, необходимо оценить математическое ожидание $M[j']$ значения j' для каждой строки i массива c , соответствующего префиксу $\mathbf{x}[1..i]$ строки \mathbf{x} . Математическое ожидание $M[j']$ можно оценить путем вычисления математического ожидания значения расстояния преобразования между префиксом $\mathbf{p}[1..j]$ паттерна \mathbf{p} и случайно выбранной строки \mathbf{x} длиной не более j . В работе [226] утверждается: “вполне очевидно”, что $M[j'] \in O(k)$. Однако оказалось, что это утверждение “не очевидно” и оно не было доказано до 1992 года, когда появилась работа Чанга и Лэмпи (Chang, Lampe) [44] с доказательством этого факта.

Доказательство оценки $O(kn)$ временной сложности в среднем алгоритма Укконена мы начнем с доказательства двух лемм, первая из которых устанавливает одно простое свойство расстояния преобразования, а вторая лемма чисто техническая.

Лемма 10.3.2. Для всех целых $j \in 1..m$ и $i \in j..n$ справедливо неравенство

$$c[i, j] \geq d_E(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j])/2.$$

Доказательство. Пусть $d_E(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j]) = d$. Напомним, что существует такое целое число $i' \in i - j + 1..i$, что

$$c[i, j] = d_E(\mathbf{x}[i'..i], \mathbf{p}[1..j]) = d' \leq d.$$

На основе неравенства треугольника (2.5) получаем

$$d - d' \leq d_E(\mathbf{x}[i - j + 1..i], \mathbf{x}[i'..i]). \quad (10.13)$$

Очевидно следующее утверждение: расстояние преобразования между двумя строками длиной n_1 и n_2 (пусть для определенности $n_2 \geq n_1$) не меньше $n_2 - n_1$; минимум $n_2 - n_1$ этого расстояния достигается тогда, когда первая строка является суффиксом второй. На основании этого утверждения получаем неравенство

$$d_E(\mathbf{x}[i - j + 1..i], \mathbf{x}[i'..i]) \leq d_E(\mathbf{x}[i'..i], \mathbf{p}[1..j]) = d'. \quad (10.14)$$

Комбинируя неравенства (10.13) и (10.14), получаем требуемый результат. ■

Лемма 10.3.3. Пусть $p_{j,r}$ — вероятность события, что две случайные строки длиной $j > 0$ имеют общую подпоследовательность длиной $\lceil rj \rceil$, $r < 1$. Тогда для произвольного значения j и значения r , удовлетворяющего неравенствам $7/8 \leq r \leq 1$, существуют положительные константы s и t ($t < 1$), такие, что $p_{r,j} < st^j/j$.

Доказательство см. в [44]. ■

Чтобы использовать результаты последних лемм, введем величину

$$j'' = 2k/(1 - r),$$

где $7/8 \leq r \leq 1$, как того требует лемма 10.3.3. Отметим, что для всех целых $j \geq j''$ выполняется неравенство $j - 2k \geq rj$. Поскольку $j'' > 2$, математическое ожидание $M[j']$ ограничено сверху величиной

$$M[j'] < (j'' - 1) + \sum_{j \geq j''} (j \times P\{c[i, j] \leq k\}), \quad (10.15)$$

где $P\{c[i, j] \leq k\}$ — вероятность того, что $c[i, j] \leq k$. В этом выражении условие $c[i, j] \leq k$ имеет следствием (по лемме 10.3.2) неравенство

$$d_E(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j]) \leq 2k.$$

Однако, как показано в упражнении 10.3.7,

$$d_E(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j]) \geq j - |\text{LCS}(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j])|. \quad (10.16)$$

Поэтому из условия $c[i, j] \leq k$ также вытекает, что

$$\text{LCS}(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j]) \geq j - 2k,$$

где, если выполняется условие $j \geq j''$, $j - 2k \geq rj$. Следовательно, при выполнении неравенства $j \geq j''$ и на основании леммы 10.3.3 при соответствующем выборе констант s и t имеем

$$\begin{aligned} P\{c[i, j] \leq k\} &\leq P\{\text{LCS}(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j]) \geq rj\} \leq \\ &\leq P\{\text{LCS}(\mathbf{x}[i - j + 1..i], \mathbf{p}[1..j]) \geq \lceil rj \rceil\} < st^j/j. \end{aligned} \quad (10.17)$$

Подставив (10.17) в выражение (10.15), получаем

$$M[j'] < (j'' - 1) + \sum_{j \geq j''} (jst^j/j) = (j'' - 1) + O(1) \in O(k).$$

Отсюда следует, что алгоритм Укконена выполняется в среднем за время порядка $O(kn)$. Также заметим, что, как и в алгоритме Вагнера–Фишера, для реализации алгоритма Укконена достаточно хранить в памяти только две строки массива c . Таким образом, мы доказали следующую теорему.

Теорема 10.3.4. Алгоритм 10.3.1 решает задачу k -разностей для паттерна $p[1..m]$ и текстовой строки $x[1..n]$ в среднем за время порядка $O(kn)$ с использованием памяти объемом порядка $\Theta(m)$. ■

Приведенная верхняя граница времени выполнения в среднем алгоритма Укконена не кажется вполне удовлетворительной, однако на практике алгоритм выполняется достаточно быстро. В упражнениях 10.3.4 и 10.3.5 показано, как можно изменить этот алгоритм, чтобы он работал и с использованием расстояния Левенштейна.

Через три года после появления первоначального варианта алгоритма Укконена, он был усовершенствован на основе применения дерева суффиксов так, что его время выполнения в самом худшем случае имело порядок $O(k^2n)$ с использованием памяти объемом порядка $O(n)$ [149]. Позднее были предложены варианты этого алгоритма с временем выполнения в самом худшем случае порядка $O(kn)$ и с использованием памяти объемом порядка $O(m^2)$ [148, 98]. В 90-х годах проводилось дальнейшее совершенствование данного алгоритма [44, 222, 232], кульминацией этого процесса стал алгоритм Майерса [187], который мы рассмотрим в следующем разделе.

10.3.2 Алгоритм Майерса

В этом разделе представим алгоритм, который для достижения очень высокой эффективности (в самом худшем случае время вычислений не зависит от аппроксимационной константы k) использует бинарные векторы оригинальным способом, основываясь на особых свойствах расстояния преобразования (лемма 10.3.1). Как и в других алгоритмах, использующих бинарные векторы, будем предполагать, что алфавит индексирован (это требование обычно выполняется на практике), т.е. будем считать, что имеем алфавит $A = \{1, 2, \dots, \alpha\}$.

Матрицы разностей

Сначала покажем, как в алгоритме Майерса массив стоимостей c , вычисленный для расстояния преобразования, можно заменить матрицей разностей Δ , точнее, двумя матрицами разностей Δ и Δ' . Начнем с определения векторов строк

$\Delta_i = \Delta_i[1..m]$ *горизонтальных разностей* ($i = 1, 2, \dots, n$), где для каждого $j \in 1..m$

$$\Delta_i[j] = c[i, j] - c[i, j - 1], \quad (10.18)$$

и определения векторов строк $\Delta'_i = \Delta'_i[1..m]$ *вертикальных разностей* ($i = 1, 2, \dots, n$), где для каждого $j \in 1..m$

$$\Delta'_i[j] = c[i, j] - c[i - 1, j]. \quad (10.19)$$

Напомним (см. табл. 10.4), что нам известны значения элементов массива стоимостей, которые стоят в нулевых строке и столбце:

$$c[0, j] = j \text{ для всех } j \in 0..m, \quad (10.20)$$

$$c[i, 0] = 0 \text{ для всех } i \in 0..n. \quad (10.21)$$

Теперь на основании формул (10.18) и (10.21) мы можем вычислить все элементы *матрицы разностей* $\Delta = \Delta[1..n, 1..m]$:

$$\Delta[i, j] = \Delta_i[j] \text{ для всех } i \in 1..n \text{ и } j \in 0..m. \quad (10.22)$$

Аналогично на основании формул (10.19) и (10.20) вычисляются элементы матрицы Δ' :

$$\Delta'[i, j] = \Delta'_i[j] \text{ для всех } i \in 1..n \text{ и } j \in 0..m. \quad (10.23)$$

Для удобства дальнейших вычислений расширим определение матриц Δ и Δ' , включив в них нулевые строки и нулевые столбцы, положив при этом

$$\Delta_0[j] = \Delta'_0[j] = 0 \text{ для всех } j \in 0..m, \quad (10.24)$$

$$\Delta_i[0] = \Delta'_i[0] = 0 \text{ для всех } i \in 1..n. \quad (10.25)$$

Из приведенных формул видно, что векторы строк (10.18) и (10.19) можно вычислить последовательно $\Delta_{i-1} \rightarrow \Delta_i$ и $\Delta'_{i-1} \rightarrow \Delta'_i$ ($i = 1, 2, \dots, n$). Поэтому таким же способом можно вычислить и все элементы (10.22) и (10.23) матриц Δ и Δ' . Однако напомним, что наша цель заключается в том, чтобы для каждого значения i вычислить $c[i, m]$ и проверить выполнение неравенства $c[i, m] \leq k$. Для этого на основании формулы (10.19) можем записать

$$c[i, m] = c[0, m] + \sum_{i'=1}^i \Delta'_{i'}[m] = m + \sum_{i'=1}^i \Delta'_{i'}[m]. \quad (10.26)$$

Таким образом, вычисление массива c можно заменить вычислением матриц разностей Δ и Δ' (точнее, вычислением векторов строк Δ_i и Δ'_i).

Преимуществом такого способа вычислений является то, что в силу леммы 10.3.1 нам надо знать значения только тех элементов матриц Δ и Δ' , которые

стоят на диагоналях $\{-1, 0, 1\}$. Если мы найдем способ эффективного вычисления этих значений (а значит, и $c[i, m]$), то тем самым полностью используем специфические свойства расстояния преобразования, формализованные в лемме 10.3.1, что, в свою очередь, приведет к уменьшению общего времени решения задачи k -разностей. Эта идея реализована в алгоритме Майерса.

Чтобы лучше понять отношения между массивом стоимостей c и соответствующими матрицами разностей Δ и Δ' в случае расстояния преобразования, предлагаем сравнить табл. 10.5 и 10.6, где они приведены. (Напомним, что массив c для этого примера в случае расстояния Левенштейна приведен в табл. 10.2.)

Таблица 10.5. Массив стоимостей c для $x = \text{rests}$ и $p = \text{stress}$ на основе расстояния преобразования

j	0	1	2	3	4	5	6	
i	ε	s	t	r	e	s	s	
0	ε	0	1	2	3	4	5	6
1	r	0	1	2	2	3	4	5
2	e	0	1	2	3	2	3	4
3	s	0	0	1	2	3	2	3
4	t	0	1	0	1	2	3	3
5	s	0	0	1	1	2	2	3

Таблица 10.6. Матрицы разностей для $x = \text{rests}$ и $p = \text{stress}$

j	1	2	3	4	5	6
i	s	t	r	e	s	s
1	r	1	1	0	1	1
2	e	1	1	1	-1	1
3	s	0	1	1	1	-1
4	t	1	-1	1	1	0
5	s	0	1	0	1	0

Матрица Δ (горизонтальные разности)

j	1	2	3	4	5	6	
i	s	t	r	e	s	s	
1	r	0	0	0	-1	-1	-1
2	e	0	0	0	1	-1	-1
3	s	0	-1	-1	-1	1	-1
4	t	0	1	-1	-1	-1	0
5	s	0	-1	1	0	0	0

Матрица Δ' (вертикальные разности)

Двухшаговая рекурсия

Теперь получим формулы для вычисления векторов строк Δ_i и Δ'_i . Сначала рассмотрим Δ_i . Применяя формулу (10.18) и базовую рекурсию (10.4), получаем

следующую формулу, справедливую для всех $i \in 1..n$ и $j \in 0..m$.

$$\begin{aligned} \Delta_i[j] &= c[i, j] - c[i, j - 1] = \\ &= \min\{c[i - 1, j] + 1, c[i, j - 1] + 1, c[i - 1, j - 1] + d(x[i], p[j])\} - \\ &\quad - c[i, j - 1]. \end{aligned}$$

Как и в алгоритмах ДБГ и ПДБГ, будем считать, что значение $d(x[i], p[j])$ можно получить как значение элемента, стоящего в позиции $[x[i], j]$ бинарного массива $t = t[1..n, 1..m]$ (определение этого массива см. в (10.6)). Напомним, что массив t вычисляется на предварительном этапе алгоритма за время порядка $\Theta(\alpha \lceil m/w \rceil + m)$ (теорема 7.4.4). Далее для краткости записи будем использовать обозначение $t_{ij} = t[x[i], j]$. На основании формул (10.18), (10.19), (10.24) и (10.25) нетрудно получить формулу для вычисления горизонтальных разностей

$$\Delta_i[j] = \min \{ \Delta_{i-1}[j], \Delta'_i[j - 1], t_{ij} - 1 \} + (1 - \Delta'_i[j - 1]), \quad (10.27)$$

справедливую для всех $i \in 1..n$ и $j \in 0..m$. Аналогично получаем формулу для вычисления вертикальных разностей:

$$\Delta'_i[j] = \min \{ \Delta_{i-1}[j], \Delta'_i[j - 1], t_{ij} - 1 \} + (1 - \Delta_{i-1}[j]), \quad (10.28)$$

Таким образом, элементы $\Delta[i, j]$ и $\Delta'[i, j]$ вычисляются на основе только элементов $\Delta[i - 1, j]$ и $\Delta'[i, j - 1]$ и предварительно вычисленных значений t_{ij} . Формулы (10.27) и (10.28) составляют двухшаговую рекурсию, которая является основой алгоритма Майерса.

Использование бинарных операций

Чтобы реализовать вычисления, определяемые формулами (10.27) и (10.28), с помощью бинарных операций надо сначала представить переменные, участвующие в этих формулах, в бинарном виде. Для этого введем четыре новых бинарных вектора, которые будут представлять положительные (P) и отрицательные (M) значения векторов Δ_i и Δ'_i . В частности для векторов Δ_i определим бинарные векторы $P_i[1..m]$ и $M_i[1..m]$ следующим образом ($j \in 1..m$):

$$\begin{aligned} P_i[j] &= 1, \text{ если } \Delta_i[j] = 1, \text{ и } P_i[j] = 0 \text{ в противном случае;} \\ M_i[j] &= 1, \text{ если } \Delta_i[j] = -1, \text{ и } M_i[j] = 0 \text{ в противном случае.} \end{aligned}$$

Аналогично для Δ'_i определяются бинарные векторы $P'_i[1..m]$ и $M'_i[1..m]$:

$$\begin{aligned} P'_i[j] &= 1, \text{ если } \Delta'_i[j] = 1, \text{ и } P'_i[j] = 0 \text{ в противном случае;} \\ M'_i[j] &= 1, \text{ если } \Delta'_i[j] = -1, \text{ и } M'_i[j] = 0 \text{ в противном случае.} \end{aligned}$$

Из этих определений следует, что $\Delta_i[j] = 0$ только тогда, когда $P_i[j] = M_i[j] = 0$, и $\Delta'_i = 0$ только в том случае, если $P'_i[j] = M'_i[j] = 0$. Поэтому мы имеем возможность во всех ситуациях заменить вычисление векторов Δ_i и Δ'_i вычислением векторов P_i , M_i , P'_i и M'_i .

Рассмотрим процесс вычисления векторов горизонтальных разностей в соответствии с формулами (10.27). Покажем, что эти вычисления можно заменить вычислениями с использованием только бинарных операций, выполняемых над отдельными битами. Для этого можно применить следующую последовательность операций.

$$\begin{aligned} K[j] &\leftarrow \overline{t_{ij}} \vee M_{i-1}[j], \\ P_i[j] &\leftarrow M'_i[j-1] \vee \overline{(K[j] \vee P'_i[j-1])}, \\ M_i[j] &\leftarrow P'_i[j-1] \wedge K[j], \end{aligned} \quad (10.29)$$

где

- $K = K[1..m]$ — бинарный вектор, введенный для удобства вычислений,
- знак \vee обозначает бинарную операцию логического сложения OR (результат этой операции равен 0 только в том случае, если *оба* операнда равны 0),
- знак \wedge обозначает бинарную операцию логического умножения AND (результат этой операции равен 0 только в том случае, если *хотя бы один* из операндов равен 0),
- обозначение $\overline{\langle \text{выражение} \rangle}$ указывает на операцию логического отрицания, примененную к логическому выражению $\langle \text{выражение} \rangle$.

Чтобы удостовериться в правильности последних формул, сначала заметим, что в формуле (10.27) участвует всего три переменные: две из них ($\Delta_{i-1}[j]$ и $\Delta'_i[j-1]$) могут принимать только значения $\{-1, 0, 1\}$, а третья (t_{ij}) является по определению бинарной. Поэтому всего может быть $3 \times 3 \times 2 = 18$ комбинаций значений переменных, для которых следует выполнить вычисления по формуле (10.27). Рассмотрим одну из возможных комбинаций переменных, например такую: $\Delta_{i-1}[j] = 1$, $\Delta'_i[j-1] = -1$, $t_{ij} = 1$. Для использования в формулах (10.29) вычисляем значения $M_{i-1}[j] = 0$, $M'_i[j-1] = 1$, $P'_i[j-1] = 0$ и $\overline{t_{ij}} = 0$. Теперь, применяя формулу (10.27), находим, что $\Delta_i[j] = 1$, а на основании формул (10.29) имеем $K[j] = 0$, $P_i[j] = 1$ и $M_i[j] = 1$. Очевидно, результаты вычисления по формулам (10.27) и (10.29) совпадают. Для доказательства правильности формул (10.29) осталось проверить еще 17 комбинаций значений переменных, что мы оставим для упражнения 10.3.8.

Вследствие определенной симметрии горизонтальных и вертикальных разностей, нетрудно получить аналог формул (10.29) для векторов вертикальных

разностей (10.28):

$$\begin{aligned}
 K'[j] &\leftarrow \overline{t_{ij}} \vee M'_i[j - 1], \\
 P'_i[j] &\leftarrow M_{i-1}[j] \vee \overline{(K'[j] \vee P_{i-1}[j])}, \\
 M'_i[j] &\leftarrow P_{i-1}[j] \wedge K'[j],
 \end{aligned}
 \tag{10.30}$$

где $K' = K'[1..m]$ — бинарный вектор, введенный для удобства вычислений.

Логика “клеточных” вычислений

Если значения $K[j]$ вычислены, то на формулы (10.29) можно смотреть как на формулы вычисления горизонтальных разностей, записанных в ячейке (“клетке”) $[i, j]$ массива c , использующие значения вертикальных разностей, записанных в ячейке $[i, j - 1]$ того же массива c . Аналогично, если известны значения $K'[j]$, тогда формулы (10.30) можно применить для вычисления вертикальных разностей, записанных в ячейке $[i, j]$ массива c , используя при этом только значения горизонтальных разностей, записанных в ячейке $[i - 1, j]$. Эта зависимость между горизонтальными и вертикальными разностями показана графически на рис. 10.1, где в объединенной таблице представлены элементы как матрицы горизонтальных разностей Δ (эти элементы обозначены буквой h), так и матрицы вертикальных разностей Δ' (эти элементы обозначены буквой v).

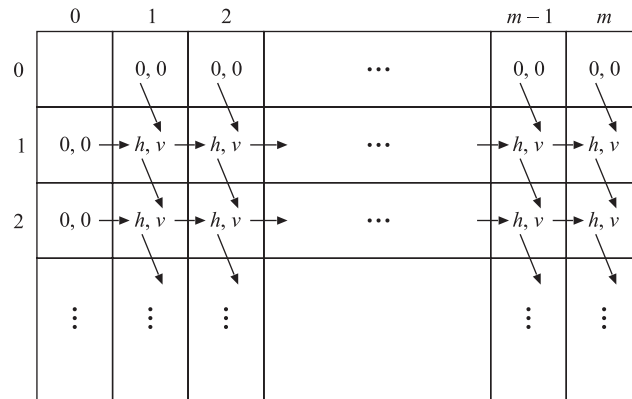


Рис. 10.1. Схема вычисления горизонтальных и вертикальных разностей

Оставив на время в стороне вопрос о вычислении векторов K и K' , здесь отметим, что теперь можно вычислить все элементы матриц Δ и Δ' строка за строкой, используя в качестве начальных значений нулевые значения из нулевой строки и столбца (см. формулы (10.24) и (10.25)). Так, на i -м шаге вычислений

- на основе значений $\Delta_{i-1}[j]$ вычисляются значения $\Delta'_i[j]$ (т.е. на основе горизонтальных разностей вычисляются вертикальные),
- на основе значений $\Delta'_i[j-1]$ вычисляются значения $\Delta_i[j]$ (здесь на основе вертикальных разностей вычисляются горизонтальные).

Если ввести переменную $cost$ (стоимость) с начальным значением $cost \leftarrow m$, то одновременно с описанными выше вычислениями на i -м шаге можно вычислить значение $c[i, m]$:

$$cost \leftarrow cost + \Delta'_{i-1}[m]. \quad (10.31)$$

Точно так же, как и в алгоритмах ДБГ и ПДБГ, в алгоритме Майерса в повышении эффективности вычислений определяющую роль играет тот факт, что бинарные операции над бинарными векторами можно выполнять совместно, а не по отдельности. Стандартные операции над машинными словами могут выполняться за константное время с помощью одной машинной команды сразу применительно ко всем w бит, составляющим машинное слово (обычно w равно 32 или 64). Так, используя стандартную операцию *rightshift* (сдвиг вправо), можно получить доступ сразу ко всем нужным $j-1$ элементам бинарного вектора без непосредственного указания адресов этих элементов и не используя ссылку на j -й элемент. Более того, если $m > w$, то в этом случае бинарный вектор будет храниться в $\lceil m/w \rceil$ смежных машинных словах, и стандартные бинарные операции над ними можно реализовать как операцию над одним словом.

После такого “лирического” отступления, попытаемся записать с помощью стандартных бинарных операций вычисления, задаваемые формулами (10.29)–(10.31). (Далее будем использовать обозначение $t_i = t[x[i], 1..m]$.)

1. $K' \leftarrow \bar{t}_i \vee \text{rightshift}(M'_i, 1)$.
2. $P'_i \leftarrow M_{i-1} \vee \overline{(K' \vee P_{i-1})}$.
3. $M'_i \leftarrow P_{i-1} \wedge K'$.
4. $cost \leftarrow cost + P'_i[m] - M'_i[m]$. (10.32)
5. $K \leftarrow \bar{t}_i \vee M_{i-1}$.
6. $P_i \leftarrow \text{rightshift}(M'_i, 1) \vee \overline{(K \vee \text{rightshift}(P'_i, 1))}$.
7. $M_i \leftarrow \text{rightshift}(P'_i, 1) \wedge K$.

Второй и третий операторы из набора (10.32) вычисляют значения P' и M' на основании значений K' (вычисленных с помощью первого оператора) и предыдущих значений P и M . Это позволяет четвертому оператору вычислить значение $cost$. Шестой и седьмой операторы вычисляют текущие значения P и M на основании значений K (вычисленных с помощью пятого оператора) и сдвигов текущих значений P' и M' . Очевидно, что для реализации этих вычислений необходимо хранить в памяти только переменные K, P, M, K', P' и M' , каждая из которых

является бинарным вектором длиной m , при этом операторы 1–3 и 5–7 выполняются за время порядка $\Theta(\lceil m/w \rceil)$, а 4-й оператор — за константное время. Таким образом, для реализации этих вычислений необходима память объемом порядка $\Theta(\lceil m/w \rceil)$ машинных слов, а общее время выполнения всех n шагов вычислений составит величину порядка $\Theta(\lceil m/w \rceil n)$. (Мы предполагали, что оператор *rightshift* после сдвига ставит в позицию $j = 1$ значение 0.)

К сожалению, набор операторов (10.32) работать не будет! Проблема заключается в том, что вычисление значений \mathbf{K}' в первом операторе зависит от вектора M'_i , который вычисляется только в третьем операторе. Однако в работе [187] с помощью достаточно сложных выкладок показано, что \mathbf{K}' можно вычислить по формуле

$$\mathbf{K}' \leftarrow \bar{t}_i \vee ((\bar{t}_i \wedge P_{i-1}) + P_{i-1}) \hat{\ } P_{i-1}, \quad (10.33)$$

где символ $\hat{\ }$ обозначает логическую операцию исключающего ИЛИ (эта операция возвращает значение 0 только в том случае, если операнды имеют одинаковые значения). Мы не будем приводить доказательство формулы (10.33), однако заменим этой формулой 1-й оператор в наборе (10.32).

Мы полностью описали алгоритм Майерса, осталось его записать в явном виде — оставим это для упражнения 10.3.10. Сформулируем основной результат данного подраздела.

Теорема 10.3.5. Алгоритм Майерса, основанный на формулах (10.32) и (10.33), корректно вычисляет все k -приближенные совпадения паттерна $\mathbf{p} = \mathbf{p}[1..m]$ и текстовой строки $\mathbf{x} = \mathbf{x}[1..n]$ на основе расстояния преобразования за время порядка $\Theta(\lceil m/w \rceil n)$ с использованием памяти объемом порядка $\Theta(\lceil m/w \rceil)$. ■

Как отмечалось в разделе 7.4, мы можем положить, что $w \in \Omega(\log n)$. Поэтому верхние границы теоремы 10.3.5 можно записать в форме $O(kn/\log n)$ и $O(mn/\log n)$ соответственно.

Усовершенствование алгоритма Майерса

В предыдущем разделе мы рассмотрели алгоритм Укконена; ожидаемое время вычисления достигает величины порядка $O(kn)$ благодаря тому, что диагональные элементы в массиве \mathbf{c} составляют неубывающую последовательность (лемма 10.3.1). Здесь существенно, что для каждой строки i массива \mathbf{c} известно значение j' , такое, что $\mathbf{c}[i, j'] \leq k$. Это ограничивает вычисляемые элементы в строке $i + 1$ столбцом $j' + 1$. Поскольку математическое ожидание числа j' имеет порядок $O(k)$, то и математическое ожидание времени обработки каждой строки массива \mathbf{c} также будет иметь такой же порядок.

В работе [187] показано, что такую же стратегию можно применить и в алгоритме Майерса, ограничивая величиной порядка $O(k)$ ожидаемое число битов,

которые будут обрабатываться в каждой строке матриц разностей. Далее такой усовершенствованный алгоритм Майерса для краткости будем называть алгоритмом ММ (модифицированный Майерса). Для этого алгоритма имеет место следующая теорема.

Теорема 10.3.6. Алгоритм ММ вычисляет все k -приближенные совпадения паттерна $p = p[1..m]$ и текстовой строки $x = x[1..n]$ на основе расстояния преобразования в среднем за время порядка $\Theta(\lceil k/w \rceil n)$ и в самом худшем случае за время порядка $\Theta(\lceil m/w \rceil n)$. ■

Эффективность алгоритмов решения задачи k -разностей

В отличие от алгоритмов непосредственного использования массива стоимостей c (такой алгоритм описан в разделе 10.1), самый ранний нетривиальный алгоритм решения задачи k -разностей был предложен Ландау и Вишкиным в 1986 году [150]. Этот алгоритм неоднократно усовершенствовался, в особенности Чангом и Лоулером (Chang, Lawler) [45, 46]. В результате появилось целое семейство алгоритмов, которые получили название “фильтрационные”. Время выполнения этих алгоритмов в самом худшем случае составляет $O(kn + m)$ и не меньше в среднем величины порядка $O((kn \log m)/m)$ для значений $k < m/(\log m + O(1))$. Как видно, эти алгоритмы могут конкурировать с алгоритмом ММ.

Другой подход, основанный на специальном определении расстояния между строками [81, 224], был разработан Укконеном [224]. Этот подход позволяет разрабатывать эффективные алгоритмы k -приближенного сравнения с паттерном. Здесь сначала вычисляется множество подстрок строки x , которое гарантированно содержит все решения задачи k -разностей, а затем используется эффективный алгоритм для удаления из этого множества тех подстрок, которые не содержат решения данной задачи.

В работе [187] приведены результаты экспериментального сравнения алгоритма ММ с другими современными алгоритмами решения задачи k -разностей. В частности, было проведено сравнение с алгоритмом [44], основанным на подходе Укконена, и с алгоритмами [232, 26], основанными на использовании бинарных векторов, методика работы с которыми получила дальнейшее развитие в работах [24, 233]. (Интересный алгоритм [26] будет описан в следующей главе (подраздел 11.2.4) как основной алгоритм для приближенного сравнения с множественными паттернами.) На очень широком наборе тестовых данных алгоритм ММ показал более быструю работу по сравнению с другими алгоритмами, что подтвердило на практике теоретические утверждения теоремы 10.3.6.

Таким образом, на сегодняшний день алгоритм ММ является лучшим выбором для приближенного сравнения с паттерном, если используется функция расстояния преобразования. В следующем разделе рассмотрим еще один алгоритм, который не является таким быстрым на практике и в теории, как алгоритм

ММ для расстояния преобразования, однако имеет значительно более широкую область применения.

Упражнения 10.3

1. Докажите утверждение *c*) леммы 10.3.1.
2. Докажите корректность алгоритма Укконена.
3. Охарактеризуйте ситуации, в которых алгоритм Укконена требует времени выполнения порядка $\Theta(mn)$.
4. Сформулируйте и докажите аналог леммы 10.3.1 для расстояния Левенштейна.
5. На основе предыдущего упражнения разработайте аналог алгоритма Укконена для расстояния Левенштейна.
6. Справедлива ли лемма 10.3.2 для других типов функций расстояния? Расширьте рамки этой леммы и докажите ее расширенный вариант, явно сформулировав предположения, на которых она будет основана.
7. Докажите неравенство (10.16).
8. Докажите корректность формул (10.27) и (10.28).
9. Проверьте эквивалентность формул (10.27) и (10.29).
10. Запишите в явном виде алгоритм Майерса.
11. Обсудите проблемы, которые могут возникнуть в алгоритме Майерса, если вместо расстояния преобразования использовать расстояние Левенштейна.

10.4 Быстрый и гибкий алгоритм Ву и Менбера

В разделе 10.2 описано обобщение алгоритма ДБГ (см. раздел 7.4), которое эффективно решает задачу k -несовпадений. Здесь мы рассмотрим другой алгоритм (технически более сложный), который расширяет подход алгоритма ДБГ, что обеспечивает эффективное решение задач приближенного сравнения с паттерном для широкого круга различных функций расстояния. Этот алгоритм решает не только задачу k -разностей, описанную в предыдущем разделе, но и задачу сравнения с регулярными выражениями (раздел 2.2). Описываемый в этом разделе алгоритм предложен Ву и Менбером (Wu, Manber) [233]; будем для краткости называть его алгоритмом ВМ. Алгоритм ВМ для большинства типов расстояний в самом худшем случае выполняется за время порядка $O(k\lceil m/w \rceil n)$, где w — длина машинного слова, а k ($k \geq 1$) — константа (граница) аппроксимации [233]. Хотя этот алгоритм не такой быстрый, как алгоритм ММ, он обладает большей гибкостью.

Как и большинство алгоритмов приближенного сравнения с паттерном, алгоритм ВМ использует соотношения рекурсии (10.4) для неотрицательной матрицы стоимостей $c = c[0..n, 1..m]$ (раздел 9.1), где для всех $i \in 1..n$ и $j \in 2..m$

$$c[i, j] \leftarrow \min\{c[i-1, j] + d(x[i], \varepsilon), \\ c[i, j-1] + d(\varepsilon, p[j]), \\ c[i-1, j-1] + d(x[i], p[j])\}.$$

Эта формула показывает, что $c[i, j]$ вычисляется как минимум сумм трех значений соседних элементов массива c и стоимостей операций соответственно удаления (буквы $x[i]$), вставки (буквы $p[j]$) и постановки (буквы $p[j]$ вместо буквы $x[i]$). В разделе 9.1 мы видели, что алгоритм Винера–Фишера использует эту формулу для непосредственного вычисления массива c . Однако в алгоритме ВМ, как и в алгоритме ДБГ, для более эффективного вычисления данного массива используются бинарные векторы.

Будем использовать бинарный вектор $s = s[1..m]$ и бинарный массив $t = t[1..n, 1..m]$, определенные в разделе 7.4. Напомним, что для заданного значения i и для всех $j \in 1..m$ $s[j] = 0$ только тогда, когда $x[i-j+1..i] = p[1..j]$, тогда как $t[h, j] = 0$ только в том случае, если $p[j] = h$, где h — буква алфавита. Конечно, здесь, как и в алгоритме ДБГ, на алфавит накладываются ограничения: алфавит A должен быть конечен, известен заранее и индексирован (см. раздел 4.1). В этом случае без уменьшения общности его можно представить в форме $\{1, 2, \dots, \alpha\}$.

В алгоритме ВМ вместо одного бинарного вектора s используется $k+1$ подобных бинарных векторов, которые для $q \in 0..k$ определяются следующим образом:

$$s_q[j] = 0, \text{ если } c[i, j] \leq q, \text{ и } s_q[j] = 1 \text{ в противном случае.} \quad (10.34)$$

Отметим, что $s_0 = s$, где s — определенный выше “стандартный” бинарный вектор. Также заметим, что для фиксированного i и любого целого $q \in 0..k-1$ выполняется импликация

$$s_q[j] = 0 \Rightarrow s_{q+1}[j] = 0. \quad (10.35)$$

Для сокращения записи введем обозначения $I = d(x[i], \varepsilon)$, $D = d(\varepsilon, p[j])$ и $S = d(x[i], p[j])$. Таким образом, I , D и S — неотрицательные целые числа, причем $S = 0$ только в том случае, если $t[x[i], j] = 0$. В этих обозначениях формулу (10.4) можно переписать так:

$$c[i, j] \leftarrow \min\{c[i-1, j] + I, c[i, j-1] + D, c[i-1, j-1] + S\}. \quad (10.36)$$

Также обозначим через $s_q^{(i)}$ ($q = 0, 1, \dots, k$) $k+1$ бинарные векторы, соответствующие строкам i массива c . Фактически, что необходимо сделать в алгоритме ВМ,

так это для каждого i ($i = 1, 2, \dots, n$) вычислить все векторы $\mathbf{s}_q^{(i)}$ ($q = 0, 1, \dots, k$) на основе вычисленных на предыдущем шаге векторов $\mathbf{s}_q^{(i-1)}$. Как и в алгоритме ДБГ, здесь вычисления производятся *на месте*, поэтому требуется хранить в памяти только $k + 1$ бинарных векторов, а не все $(k + 1)n$ векторов.

Рассмотрим значение $\mathbf{s}_q^{(i)}[j]$ для некоторого $q \in 0..k$. Мы хотим, чтобы $\mathbf{s}_q^{(i)}[j]$ равнялось 0, если выполняется неравенство $c[i, j] \leq q$, а в противном случае $\mathbf{s}_q^{(i)}[j]$ должно равняться 1. Из формулы (10.36) следует, что $\mathbf{s}_q^{(i)}[j] = 0$ только в том случае, если выполняется одно из следующих условий.

1. $c[i - 1, j] \leq q - I \Leftrightarrow +\mathbf{s}_{q-I}^{(i-1)}[j] = 0$.
2. $c[i, j - 1] \leq q - D \Leftrightarrow +\mathbf{s}_{q-D}^{(i)}[j - 1] = 0$.
3. $c[i - 1, j - 1] \leq q - S \Leftrightarrow +\mathbf{s}_{q-S}^{(i-1)}[j - 1] = 0$.

Отметим, что третье условие выполняется только тогда, когда $S = 0$ (т.е. когда $t[x[i], j] = 0$). Учитывая данные условия, значение $\mathbf{s}_q^{(i)}[j]$ можно вычислить по формуле

$$\begin{aligned} \mathbf{s}_q^{(i)}[j] \leftarrow & \mathbf{s}_{q-I}^{(i-1)}[j] \wedge \mathbf{s}_{q-D}^{(i)}[j - 1] \wedge \mathbf{s}_{q-S}^{(i-1)}[j - 1] \\ & \wedge \left(\mathbf{s}_q^{(i-1)}[j - 1] \vee t[x[i], j] \right), \end{aligned} \quad (10.37)$$

где символ \wedge обозначает операцию логического умножения AND, а символ \vee — операцию логического сложения OR.

Напомним (см. раздел 7.4), что бинарный вектор $\mathbf{s}^{(i)}$ формируется путем сдвига вправо вектора $\mathbf{s}^{(i-1)}$, т.е. для всех $j \in 2..m$ $\mathbf{s}^{(i)}[j] = \mathbf{s}^{(i-1)}[j - 1]$, а в самую левую позицию нового вектора помещается 0 (поэтому всегда $\mathbf{s}^{(i)}[1] = 0$). Используя форму записи из раздела 7.4 и учитывая, что для всех значений q $\mathbf{s}_q^{(0)}[0..m] = 01^m$, формулу (10.37) можно записать следующим образом ($i = 1, 2, \dots, n$).

$$\begin{aligned} \mathbf{s}_q^{(i)} \leftarrow & \mathbf{s}_{q-I}^{(i-1)} \wedge \mathit{rightshift}(\mathbf{s}_{q-D}^{(i)}, 1) \\ & \wedge \mathit{rightshift}(\mathbf{s}_{q-S}^{(i-1)}, 1) \wedge \left(\mathit{rightshift}(\mathbf{s}_q^{(i-1)}, 1) \vee t[x[i], 1..m] \right). \end{aligned} \quad (10.38)$$

После каждого вычисления формулы (10.38) необходимо проверить текущее значение $\mathbf{s}_k[m] = \mathbf{s}_k^{(i)}[m]$, чтобы определить (в соответствии с определением (10.34)), выполняется ли неравенство $c[i, m] \leq k$, т.е. достигается ли k -приближенное совпадение паттерна \mathbf{p} с подстрокой $\mathbf{x}[1..i]$. На основе формулы (10.38) и строится алгоритм ВМ, являясь полным аналогом алгоритма ДБГ. Явную запись алгоритма ВМ оставим для (простого) упражнения 10.4.1.

Отметим, что вычисления по формуле (10.38) зависят от фиксированных целочисленных значений I , D и S , которые не превосходят значения q . Если $I = D = S = 1$, то формула (10.38) вычисляется для расстояния преобразования.

Если же $I = D = 1$ и $S = 2$, то формула (10.38) работает для расстояния Левенштейна. В табл. 10.7 приведены массив t и векторы $s_0^{(i)}$ и $s_1^{(i)}$, вычисленные для расстояния Левенштейна для задачи 1-приближенного сравнения паттерна $p = str$ и строки $x = rests$ (массив стоимостей для этого примера приведен в табл. 10.2). Чтобы формулу (10.38) использовать для других функций расстояния, надо просто в соответствии с этой функцией изменить значения констант I , D и S .

Таблица 10.7. Приближенное сравнение с паттерном на основе расстояния Левенштейна: $A = \{e, r, s, t\}$, $p = str$, $x = rests$, $k = 1$

	s	t	r		s_0	s_1
e	1	1	1		r	111
r	1	1	0		e	111
s	0	1	1		s	011
t	1	0	1		t	001
	Массив t				s	011
				Бинарные векторы		

Для того чтобы найти все k -приближенные совпадения паттерна и заданной строки, в алгоритме ВМ вычисления по формуле (10.38) должны выполняться $(k + 1)n$ раз. Каждое вычисление по этой формуле требует выполнения трех сдвигов вправо, трех операций AND и одной операции OR над каждым машинным словом, составляющих каждый вектор s_q ($q = 1, 2, \dots, k$). Эти операции над одним машинным словом выполняются за константное время. Если, как и ранее, через w обозначить количество битов в машинном слове и принять допущения о предварительном вычислении и хранении массива t (теорема 7.4.4), то мы можем сформулировать аналог теоремы 7.4.3.

Теорема 10.4.1. Алгоритм ВМ вычисляет все k -приближенные совпадения паттерна $p = p[1..m]$ и текстовой строки $x = x[1..n]$ за время порядка $\Theta(k \lceil m/w \rceil n)$ с использованием памяти объемом порядка $\Theta((\alpha + k) \lceil m/w \rceil)$. ■

Как и для всех алгоритмов, использующих методологию алгоритма ДБГ, в случае малых значений m (и даже в случае только ограниченных значений m) для получения асимптотической оценки сложности алгоритма выражение $\lceil m/w \rceil$ можно включить в константу пропорциональности, что для алгоритма ВМ дает асимптотическую оценку временной сложности в виде $\Theta(kn)$. (Дальнейшие исследования в этом направлении оставим для упражнения 10.4.4.) Напомним также замечание, сделанное в конце раздела 7.4, в соответствии с которым можно положить $w \in \Omega(\log n)$. Тогда для алгоритма ВМ получаем временную границу $O(kmn / \log n)$.

Итак, в этом разделе описан эффективный и очень гибкий алгоритм для решения задачи приближенного сравнения с паттерном. Этот алгоритм легко адаптируется ко всем функциям расстояния, для которых стоимости операций вставки, удаления и подстановки фиксированы и выражаются целыми числами. С этой точки зрения алгоритм ВМ имеет бóльшую область применимости (более гибкий), чем алгоритмы решения задач k -несовпадений и k -разностей, описанные в разделах 10.2 и 10.3. Однако этот алгоритм имеет свои ограничения: его нельзя использовать, если стоимость операций редактирования переменная или принимает нецелые значения, как это бывает в случае весовых матриц (см. раздел 2.2). Но, как будет показано далее, алгоритм ВМ относительно просто и без потери эффективности можно модифицировать для решения практически любых задач сравнения с паттерном. Отметим, что такой расширенный алгоритм ВМ реализован в виде пакета *agrep* на языке Unix; его можно найти в Internet по адресу <http://webglimpse.net>.

Сравнение множеств букв

Во многих приложениях задачи сравнения с паттерном возникает ситуация, когда совпадение с паттерном считается достигнутым в некоторой позиции j , если в этой позиции стоит буква из некоторого заранее определенного множества букв. Например, в строке можно искать вхождение подстроки “ref”, которая предшествует натуральному числу из интервала от 1 до 7, либо можно искать вхождение подстроки “ion”, которой предшествует одна из букв “c”, “s” или “t”. В этих ситуациях можно определить паттерны, подобные $p = ref[1 - 7]$ или $p = [c, s, t]ion$, где множество букв для сравнения идентифицируется метасимволами в виде квадратных скобок []. Алгоритм ВМ может работать с такими паттернами, для чего необходимо внести только небольшие изменения на предварительном этапе вычисления массива t . Все, что необходимо сделать для корректной работы алгоритма ВМ, так это положить

$$t[1, 4] \leftarrow t[2, 4] \leftarrow \dots \leftarrow t[7, 4] \leftarrow 0$$

для паттерна $p = ref[1 - 7]$, что позволит распознать цифры 1–7 в 4-й позиции паттерна, или

$$t[c, 1] \leftarrow t[s, 1] \leftarrow t[t, 1] \leftarrow 0$$

для паттерна $p = [c, s, t]ion$, что позволит распознать одну из букв c, s, t в 1-й позиции паттерна. Никаких других изменений в алгоритме ВМ делать не требуется.

Сравнение с символами замещения

Напомним (см. раздел 2.2), что существует два общепринятых символа замещения:

- символ • замещает одну любую букву из алфавита A ;

символ $*$ замещает любую строку из множества A^* .

Использование символа \bullet — это случай использования специального множества букв для сравнения, который рассмотрен в предыдущем подразделе и который требует только небольших изменений в вычислении массива t на предварительном этапе алгоритма ВМ.

Использование символа замещения $*$ для строк уже требует внесения некоторых изменений в сам алгоритм ВМ. Рассмотрим паттерн вида

$$p = p[1..j_1] * p[j_1 + 1..j_2] * \dots * p[j_{r-1} + 1..j_r] * p[j_r + 1..m],$$

где метасимвол $*$ вставлен в паттерн после позиций j_1, j_2, \dots, j_r , но этот метасимвол не рассматривается как буква паттерна p . Поэтому, чтобы указать местоположение метасимволов $*$, вводится новый бинарный вектор $t^* = t^*[1..m]$, в котором $t^*[j] = 0$ только в том случае, если $j \in \{j_1, j_2, \dots, j_r\}$.

Если для некоторых $h \in 1..r$, $i \in 1..n$ и $q \in 0..k$ выполняется $s_q^{(i-1)}[j_h] = 0$, то в этом случае имеет место q -приближенное совпадение $p[1..j_h]$ (включая все вхождения $*$ слева от позиции j_h) с подстрокой $x[i'..j-1]$ строки x . Поскольку после $p[1..j_h]$ следует символ $*$, каждая подстрока $x[i'..i'']$, где $i'' > i-1$, также имеет q -приближенное совпадение с $p[1..j_h]$. Другими словами, если для некоторой позиции $i-1$ строки x выполняется равенство $s_q^{(i-1)}[j_h] = 0$, то для всех позиций $i'' > i-1$ должно выполняться аналогичное равенство $s_q^{(i'')}[j_h] = 0$. Чтобы добиться такого выполнения, надо после каждого вычисления по формуле (10.38) сделать присвоение

$$s_q^{(i)} \leftarrow s_q^{(i)} \wedge (s_q^{(i-1)} \vee t^*). \quad (10.39)$$

В этом случае, если в результате вычисления формулы (10.38) получили $s_q^{(i)} \leftarrow 1$, то присвоение $s_q^{(i)}[j] \leftarrow 0$ будет иметь место только тогда, когда $t^*[j] = 0$ и $s_q^{(i-1)}[j] = 0$. Поэтому формула (10.39) будет присваивать нулевые значения только в соответствии с позициями j_1, j_2, \dots, j_r паттерна и для позиций, следующих после позиции $i-1$ строки x .

Хотя вычисление по формуле (10.39) требует выполнение дополнительных операций AND и OR для каждого значения q и каждого значения i , это не влияет на асимптотическую оценку сложности алгоритма.

Комбинирование точного и приближенного совпадений

Возможны ситуации, когда при сравнении строки с паттерном необходимо найти точное совпадение с одной частью паттерна, а с другой частью паттерна достаточно приближенного совпадения. Например, в последовательностях ДНК надо найти фрагмент вида CGATGAGAGCGC, в котором первые четыре буквы (CGAT) и три последние (CGC) будут точно соответствовать этим буквам, а для

средней части, состоящей из пяти букв, достаточно 2-приближенного совпадения с GAGAG. Осуществить такое сравнение достаточно просто, поскольку точное совпадение означает, что в определенных позициях ($j = 1 - 4$ и $j=10-12$) паттерна не разрешается выполнять операции удаления, вставки и подстановки, тогда как для позиций приближенного совпадения следует положить $s_q[j] \leftarrow 1$ для всех $q \in 0..k$. Чтобы сделать такое, введем новый бинарный вектор $e = e[1..m]$, в котором $e[j] = 1$ только в том случае, если в позиции j требуется точное совпадение. В нашем примере $e \leftarrow 11110000111$. Этот вектор используется для присваивания $s_q^{(i)}[j] \leftarrow 1$, если не выполняются условия точного совпадения $s_q^{(i-1)}[j-1] = 0$ и $x[i] = p[j]$. С использованием вектора e формулу (10.38) следует изменить следующим образом:

$$s_q^{(i)} \leftarrow ((s_{q-I}^{(i-1)} \wedge \text{rightshift}(s_{q-D}^{(i)}, 1) \wedge \text{rightshift}(s_{q-S}^{(i-1)}, 1)) \vee e) \wedge \text{rightshift}(s_q^{(i-1)}, 1) \vee t[x[i], 1..m]. \quad (10.40)$$

Хотя формула (10.40) добавляет одну операцию OR к операциям, выполняемым по формуле (10.38), асимптотическая оценка сложности алгоритма не изменяется.

Сравнение с минимальным k

Представляет интерес задача нахождения *минимального* значения k^* , для которого имеет место одно или несколько k^* -приближенных совпадений паттерна и заданной строки. Чтобы вычислить такое k^* , сначала можно выполнить точное сравнение ($k = 0$) с помощью, например, алгоритма ДБГ. Если точного совпадения нет, то затем выполняется алгоритм ВМ для $k = 1, 3, \dots, 2^B - 1$ до тех пор, пока для некоторого значения $B \geq 1$ не будет найдено, по крайней мере, одно k^* -приближенное совпадение с паттерном, при этом $k^* \in 2^{B-1}..2^B - 1$. Каждое выполнение алгоритма ВМ требует времени, пропорционального используемому значению k . Поэтому общее время, требуемое для B выполнений алгоритма ВМ, пропорционально сумме

$$\Sigma = 1 + 3 + .. + (2^B - 1) < 2^{B+1}.$$

Поскольку $k^* \geq 2^{B-1}$, то $\Sigma/k^* < 4$; следовательно, время вычисления k^* имеет порядок $O(4k^*[m/w]n)$.

Другие применения алгоритма ВМ

Как мы увидим в следующей главе, алгоритм ВМ можно применить для сравнения с регулярными выражениями, а после небольшой модернизации — для эффективного сравнения с множественными паттернами. И конечно, эти применения алгоритма можно комбинировать с рассмотренными выше применениями. Например, алгоритм ВМ можно использовать для поиска совпадения с несколькими паттернами, каждый из которых может содержать символы замещения, или

для поиска комбинированного точного и приближенного совпадения, а также для нахождения минимального значения k . Таким образом, если паттерн не является экстремально длинным, то эти примеры еще раз свидетельствуют о чрезвычайной гибкости алгоритма ВМ!

В работе [231] предложены определенные усовершенствования алгоритма ВМ. В разделе 11.2.4 будет описан другой алгоритм [26], который предлагает эффективный подход к решению задачи приближенного сравнения с множественными паттернами.

Упражнения 10.4

1. Запишите алгоритм ВМ (используя формулу (10.38)) и докажите его корректность. Удостоверьтесь, что алгоритм работает корректно, если одно из значений I , D или S превышает значение q . Объясните, почему при любом i необходимо проверить только значение $s_k[m]$, а не все значения $s_q[m]$, $q = 1, 2, \dots, k$.
2. Как изменить формулу (10.38), чтобы алгоритм ВМ вычислял расстояние Хемминга? Сравните пространственную и временную сложности измененного алгоритма ВМ с аналогичными показателями алгоритма ДБГ.
3. Для расстояния преобразования рассчитайте аналог табл. 10.7 для $k = 1$ и $k = 2$.
4. Постройте таблицу, показывающую зависимость максимальной длины паттерна m , который может обрабатывать алгоритм ВМ, от величин $k \in 1..10$ и $B \in 1..10$ для машинных слов длиной $w = 32$ и $w = 64$ (B — количество машинных слов, необходимых для хранения одного бинарного вектора s_q).

10.5 Сложность алгоритмов приближенного сравнения с паттерном

Как и в случае точного сравнения с паттерном (раздел 8.8), здесь рассмотрим теоретические границы сложности алгоритмов приближенного сравнения с паттерном. Поскольку для малых значений аппроксимационной константы k приближенное сравнение паттерна и строки близко к их точному сравнению, естественно предположить, что должны существовать алгоритмы с линейным временем выполнения в самом худшем случае и почти линейным временем в среднем — это временные ограничения алгоритмов точного сравнения с паттерном. Действительно, существуют алгоритмы решения задачи k -разностей, асимптотические оценки времени выполнения которых близки к этим требованиям. В табл. 10.8 для сравнения приведены оценки времени выполнения двух уже изученных алгоритмов (алгоритм Майерса и алгоритм ММ), одного алгоритма, который будет рассмотрен

в следующей главе (алгоритм BYN), и двух недавно разработанных алгоритмов, которые подробно в этой книге не рассматриваются.

Таблица 10.8. Сложность алгоритмов приближенного сравнения с паттерном

Алгоритм	Временная граница в самом худшем случае	Временная граница в среднем
Майерса (подразд. 10.3.2)	$O(mn/\log n)$	$O(mn/\log n)$
ММ (подразд. 10.3.2)	$O(mn/\log n)$	$O(kn/\log n)$
BYN (подразд. 11.2.4)	$O(kn(m-k)/\log n)$	
Чанга и Лоулера [45]	$\Theta(kn+m)$	$O(kn \log m/m)$
Коула и Харихарана [51]	$O(n(k^4/m+1)+m)$	

На основе этой таблицы можно сделать такие заключения.

- Если значение m мало по сравнению со значением n (например, $m \in O(\log n)$), алгоритмы Майерса и ММ в самом худшем случае выполняются за линейное и почти линейное по n время. Если же принять $k \in O(\log n)$, то алгоритм ММ выполняется в среднем за время порядка $O(n)$.
- Для алгоритма BYN временная граница выполнения в самом худшем случае в общем случае превышает аналогичные границы для алгоритмов Майерса и ММ. Разработан вариант алгоритма BYN, среднее время выполнения которого имеет порядок $O(n)$, если значение отношения k/m “среднего размера”. Однако достижение устойчивого линейного поведения в среднем для этого алгоритма является пока нерешенной задачей.
- В алгоритме Чанга–Лоулера линейное в среднем время выполнения порядка $O(n)$ гарантировано, если выполняется условие $k < m/(\log m + O(1))$, тогда как в алгоритме ММ соответствующее условие имеет вид $k \in O(\log n)$. Поэтому при выполнении условия

$$m \in \Theta(\log m \log n) \tag{10.41}$$

верхние границы по k времени выполнения обоих алгоритмов эквивалентны с некоторым преимуществом, по крайней мере теоретическим, алгоритма ММ для случая значения m , меньшего, чем требует условие (10.41), и с преимуществом алгоритма Чанга–Лоулера для больших значений m .

- Алгоритм Коула–Харихарана имеет линейное время выполнения в самом худшем случае при выполнении условия $k \in O(m^{1/4})$. В работе [51] сделано предположение, что это условие можно ослабить, потребовав $k \in O(m^{1/3})$.

В статье [191] проведен более детальный сравнительный анализ эффективности алгоритмов приближенного сравнения с паттерном. Заинтересованных читателей мы направляем к этой работе.

ГЛАВА 11

Регулярные выражения и множественные паттерны

Я ненавижу неверные слова и тщательно, упрямо и угрюмо ищу те слова, которые соответствуют данному предмету.

— Уолтер Севидж Лэндор (1775–1864),
Воображаемые беседы

Я осознаю, что многие мои коллеги будут удивлены и даже шокированы, что только в главе 11 почти на 300-й странице книги я ввожу в рассмотрение конечные автоматы, которые являются центральной идеей как компьютерных наук в целом, так и вычислительных процессов обработки строковых последовательностей в частности.

В свою защиту я могу сказать, что в книге основное внимание уделяется алгоритмам, и при их описании я использую “минималистский” подход: объяснение идей, которые лежат в основе алгоритмов, дано очень кратко с минимумом обращений к теоретическим изысканиям. Правильно это или нет — судить читателю, но мне кажется, что объяснение основ алгоритмов будет более ясным и понятным без привлечения идеи конечного автомата.

Однако теперь ситуация изменилась: алгоритмы сравнения с регулярными выражениями и с множественными паттернами невозможно описать без использования понятия конечного автомата. В этой главе мы встретимся с различными типами конечных автоматов, каждый из которых решает свои задачи. Чтобы сделать дальнейший материал книги более понятным, представим очень краткое и не

формальное введение в теорию конечных автоматов. Желая ознакомиться с теорией конечных автоматов в более полном объеме предлагаем обратиться к книге [159].

Самое простое определение конечного автомата — это представить его в виде помеченного ориентированного графа, в котором

- вершины графа представляют “состояния” — состояния конечного автомата, состояния процесса вычисления;
- метки вершин идентифицируют состояния;
- дуги обозначают переход из одного состояния в другое;
- метки дуг — отдельные буквы или строки, относящиеся к соответствующему переходу из одного состояния в другое.

Из этого определения нетрудно заметить, что мы уже встречались в этой книге с конечными автоматами — это дерево суффиксов! Чтобы сделать это замечание очевидным, на рис. 11.1 показано дерево суффиксов, которое ранее было представлено на рис. 5.2. (Отличие этих рисунков состоит только в том, что здесь из рис. 5.2 удалены связи, внутренние вершины помечены как состояния s_1, s_2, \dots, s_{11} и нисходящие ребра заменены ориентированными дугами.)

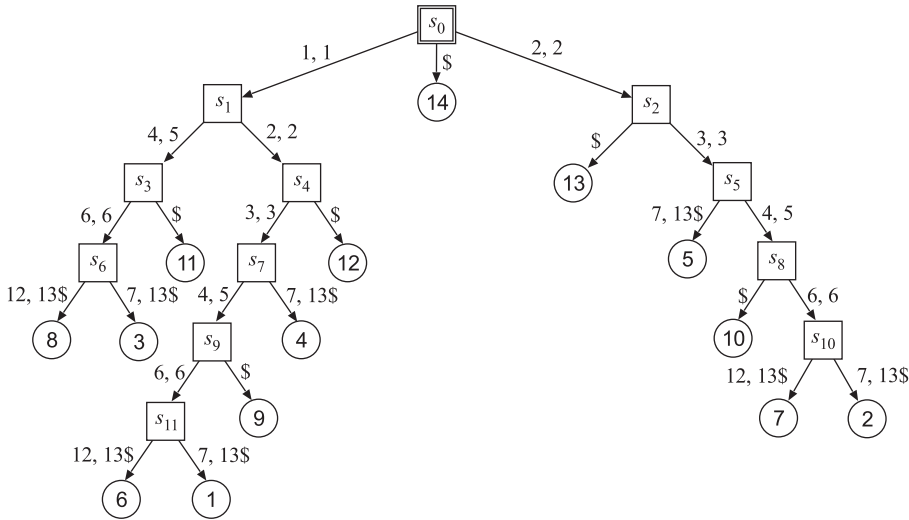


Рис. 11.1. Дерево суффиксов для $f_6 = abaababaabaab\$$ в виде конечного автомата

Если мы подставим входную строку $p = ab\$$ в вершину (состояние) s_0 , тогда первая буква $p[1] = a$ будет равна $f_6[1] = a$, что определяет переход в состояние s_1 . В этом состоянии $p[2..3] = f_6[4..5] = ab$, поэтому далее осуществляется переход в состояние s_3 , где наконец $p[4] = \$$ с переходом в конечное состояние

11, которое показывает, что $f_6[11..14] = p$. В данном случае очевидно, что произвольный паттерн p , поступающий на вход “суффиксного” автомата, достигнет конечного состояния t только в том случае, если p является суффиксом $f_7[t..14]$ строки $abaababaabaab\$$. В противном случае паттерн p “застрянет” в одной из внутренних вершин, поскольку в нем нет буквы или подстроки, совпадающих с метками дуг, выходящих из этой вершины.

Дерево суффиксов любой строки x (и даже набора строк) можно интерпретировать и исследовать подобным способом — таким образом, мы видим, что конечные автоматы действительно предлагают новый графический способ представления вычислений, выполняемых над строковыми последовательностями. В этой главе мы изучим несколько подобных типов конечных автоматов, которые будем называть *суффиксными автоматами*.

Однако следует отметить, что суффиксные автоматы решают задачу, которая в некотором смысле обратная задаче, рассматриваемой в этой главе: суффиксные автоматы легко решают задачу сравнения одного паттерна с некоторым набором текстовых строк путем включения всех этих строк в одно дерево суффиксов. В этой главе решается задача одновременного сравнения одной текстовой строки с набором паттернов (с множественным паттерном). Как мы увидим далее, для решения последней задачи можно построить конечный автомат, который включал бы все паттерны.

В разделе 11.1 сначала будут рассмотрены недетерминированные конечные автоматы (сокращенно, NDFA — от non-deterministic finite automata), с помощью которых решим задачу сравнения с регулярными выражениями. Затем покажем, как на основе NDFA построить детерминированные конечные автоматы (сокращенно, DFA) для более эффективного решения той же задачи. В последнем подразделе раздела 11.1 рассмотрим вариант алгоритма BM (см. раздел 10.4), решающий задачу k -приближенного сравнения с регулярными выражениями. Здесь этот алгоритм будет представлен в виде модели NDFA, использующей бинарные векторы. На этом основании будет сделан вывод, что алгоритм ДБГ (раздел 7.4) и его многочисленные варианты и модификации также можно представить в виде конечных автоматов.

Раздел 11.2 начнем с описания конечных автоматов, решающих задачу точного сравнения с множественными паттернами, причем основой для решения сначала будет служить алгоритм КМП (раздел 7.1), а затем — алгоритм BM (раздел 7.2). Далее представим два алгоритма для приближенного сравнения с множественным паттерном, один из которых основан на алгоритме BM, а второй предложен Бейза-Ятсом и Наварро (Baeza-Yates, Navarro).

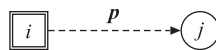
11.1 Алгоритмы для регулярных выражений

В этом разделе рассмотрим три алгоритма, два из которых “классические” (разработаны до 1970 года) и один относительно новый. Эти алгоритмы вычисляют все совпадения всех подстрок строки x с паттерном p , который представлен в виде регулярного выражения. Как указывалось в разделе 2.2, сравнение с регулярными выражениями можно считать экстремальным случаем приближенного сравнения.

11.1.1 Недетерминированные конечные автоматы

Прежде чем использовать конечные автоматы (применительно к регулярным выражениям), необходимо сначала научиться их строить. Для построения таких автоматов опишем метод, восходящий к Томпсону (Thompson) [119], основанный на четырех правилах, применяемых рекурсивно. Мы будем строить *недетерминированный конечный автомат* (N DFA). “Недетерминированность” здесь проявляется в том, что из некоторых состояний (вершин) возможно несколько переходов (выходных дуг), помеченных пустой строкой ε , и при этом определение выходной дуги неоднозначно. Как мы увидим далее, NDFA, подобно суффиксному автомату, описанному выше, имеет единственное *начальное состояние* или *источник* и, в отличие от суффиксного автомата, только одну конечную вершину, соответствующую *поглощающему состоянию* или *стоку*. Описываемый здесь NDFA имеет дуги, помеченные или одной буквой или пустой строкой, что напоминает некомпактное дерево суффиксов, описанное в разделе 2.1.

Будем представлять автомат NDFA(p), соответствующий регулярному выражению p , в виде



Здесь i и j — метки начального и поглощающего состояний соответственно, а пунктирная линия обозначает возможные пути от i до j .

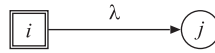
Пусть регулярное выражение p определено на алфавите $A = \{\lambda_1, \lambda_2, \dots, \lambda_\alpha\}$. Построение автомата NDFA начнем с вычисления

$$p^* \leftarrow (\lambda_1|\lambda_2|\dots|\lambda_\alpha)^*p.$$

Это позволит сравнивать паттерн p и входную строку x начиная с любой позиции этой строки, поскольку сравнение паттерна p с любой позицией строки x эквивалентно сравнению p^* со строкой x , начинаемого с первой позиции этой строки.

Далее предполагаем, что p^* сканируется по одной позиции слева направо. Автомат NDFA(p) строится по следующим правилам.

П.1. Для произвольной буквы $\lambda \in A$ (не метасимвола) автомат $\text{N DFA}(\lambda)$ определяется следующим образом:

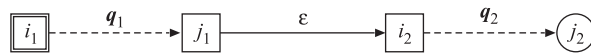


Очевидно, что определенный таким образом автомат $\text{N DFA}(\lambda)$ допустим только для одиночной буквы, но не для какой-либо строки, состоящей из двух и более букв.

П.2. Предположим, что для некоторых регулярных выражений q_1 и q_2 уже определены автоматы $\text{N DFA}(q_1)$ и $\text{N DFA}(q_2)$:

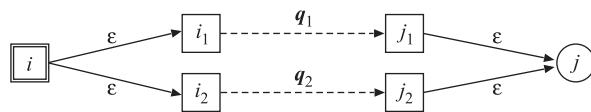


Тогда $\text{N DFA}(q_1 q_2)$ строится так:



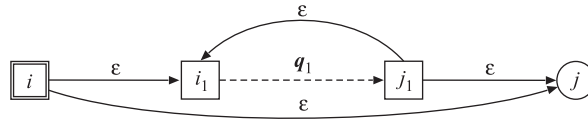
Здесь пустая метка ϵ обозначает переход из одного состояния в другое без инициации этого перехода каким-либо входным воздействием. В данном случае такой “пустой” переход переопределяет поглощающее состояние j_1 автомата $\text{N DFA}(q_1)$ в начальное состояние j_2 автомата $\text{N DFA}(q_2)$, при этом оба состояния j_1 и j_2 становятся внутренними состояниями автомата $\text{N DFA}(q_1 q_2)$. Отметим, что автомат $\text{N DFA}(q_1 q_2)$ определен только для строк вида $q_1 \epsilon q_2 = q_1 q_2$.

П.3. Пусть для регулярных выражений q_1 и q_2 уже определены автоматы $\text{N DFA}(q_1)$ и $\text{N DFA}(q_2)$. Тогда автомат $\text{N DFA}(q_1 | q_2)$ строится так:



Таким образом, поглощающие состояния автомата $\text{N DFA}(q_1)$ и начальные состояния автомата $\text{N DFA}(q_2)$ становятся обычными внутренними состояниями автомата $\text{N DFA}(q_1 | q_2)$. Отметим, что здесь пустой строкой ϵ помечены как выходные дуги (переходы) нового начального состояния i , так и входные дуги нового поглощающего состояния j . Автомат $\text{N DFA}(q_1 | q_2)$ доступен только для строк q_1 и q_2 , но не доступен ни для каких-либо других строк.

П.4. Предположим, что для некоторого регулярного выражения q_1 определен автомат $NFA(q_1)$. Тогда автомат $NFA(q_1^*)$ строится так:



В силу пустых переходов из состояния j_1 в состояние i_1 и из состояния j_1 в состояние j автомат $NFA(q_1^*)$ доступен как для q_1 , так и для любой кратной строки q_1^k , $k \geq 2$. Также в силу пустых переходов из состояния i в состояние j автомат $NFA(q_1^*)$ доступен и для пустых строк. Но для любых других строк автомат не доступен. Здесь также начальное состояние i_1 и поглощающее состояние j_1 исходного автомата $NFA(q_1)$ становятся обычными внутренними вершинами автомата $NFA(q_1^*)$.

Эти правила позволяют построить NFA для выражения p^* в тех случаях, когда в этом выражении встречаются отдельные буквы, конкатенации букв и регулярных выражений, а также метасимволы $|$ и $*$. Однако эти правила не помогут в случае вложенных операторов (определяемых с помощью метасимволов “скобки”), например, для такого выражения:

$$p = (q_1 | (q_2^* | q_3))^* | (q_4 | q_5)^*, \tag{11.1}$$

где q_i ($i = 1, 2, \dots, 5$) — регулярные выражения (возможно, обычные строки).

Предполагая, что буквы алфавита можно распознать за константное время, выражение, подобное (11.1), можно распознать и проанализировать за время, пропорциональное количеству символов, входящих в *синтаксическое дерево*. На рис. 11.2 показано такое дерево, соответствующее выражению (11.1). Отметим, что в этом дереве каждое выражение q_i , входящее в p , соответствует листу (конечному узлу). Синтаксическое дерево можно преобразовать в автомат $NFA(p)$ за время, пропорциональное $|p|$. Подробности этого процесса можно найти в [118].

Уже на основе сформулированных выше правил построения автомата $NFA(p)$ можно сделать некоторые выводы о свойствах такого автомата.

- Как показано в упражнении 11.1.2, если m — количество букв в выражении p , тогда $|p| \leq 6m$.
- Поскольку после применения любого из правил П.1–П.4 в NFA остается только одно начальное и одно поглощающее состояние, то и в $NFA(p)$ также имеется единственное начальное и единственное поглощающее состояние.
- Каждая вершина, не являющаяся поглощающим состоянием, имеет или одну исходящую дугу, помеченную буквой, или не более двух исходящих дуг,

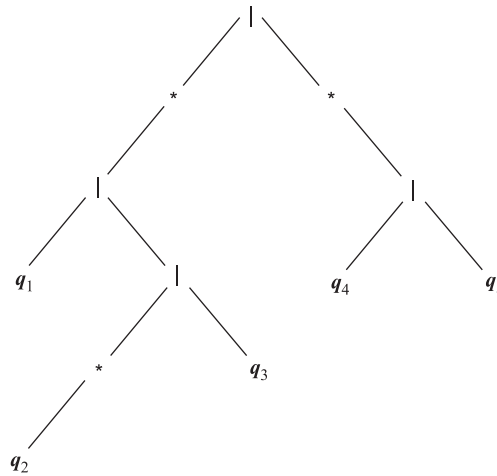


Рис. 11.2. Синтаксическое дерево для выражения (11.1)

помеченных пустой строкой ε . Аналогично каждая вершина, не являющаяся начальным состоянием, имеет или одну входящую дугу, помеченную буквой, или не более двух входящих дуг, помеченных ε . Отсюда вытекают еще два вывода, которые в дальнейшем будут полезны при построении алгоритма сравнения с регулярными выражениями:

- вершина, не являющаяся поглощающим состоянием, имеет единственную исходящую “непустую” дугу (т.е. дугу, помеченную не ε) тогда и только тогда, когда она не имеет исходящих “пустых” дуг;
 - каждая “непустая” дуга ведет к вершине, которая либо соответствует поглощающему состоянию, либо имеет “пустые” исходящие дуги.
- Пусть $\text{N DFA}(p) = (V, A)$, где V — множество вершин, а A — множество дуг. Когда, как показано в упражнении 11.1.3,

$$|V| \leq 8m, |A| \leq 13m, \tag{11.2}$$

где, как и ранее, m — количество отдельных букв в выражении p . Отсюда следует, что $\text{N DFA}(p)$ можно вычислить за время порядка $\Theta(m)$.

На рис. 11.3 показан типичный N DFA , в данном случае построенный для паттерна (регулярного выражения) $p = (a | b)^* ab^* a$, где префикс $(a | b)^*$, как отмечалось выше, гарантирует совпадение с буквами, стоящими в любой позиции любой строки x , определенной на алфавите $\{a, b\}$.

Теперь, когда показано, что автомат $\text{N DFA}(p)$ можно построить за время, пропорциональное количеству букв в регулярном выражении p , расскажем, как

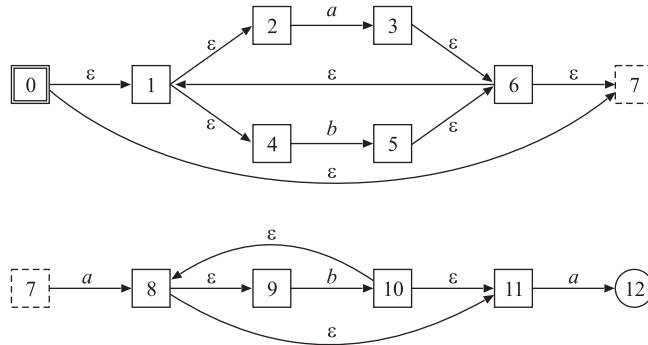


Рис. 11.3. Пример NFA

применить такой автомат. Допустим, на вход автомата поочередно поступают буквы текстовой строки $x = x[1..n]$. Для каждой вновь поступающей буквы строки x необходимо определить возможные переходы из текущего активного множества состояний (множества вершин) автомата $NFA(p)$, при этом совпадение с паттерном будет зарегистрировано только в том случае, если эти переходы приведут к поглощающему состоянию.

Для выполнения таких вычислений введем функцию *trans* (от transition — переход), которая для заданного множества состояний S и заданной буквы $\lambda \in A$ (текущая входная буква) вычислит новое множество состояний, достижимых путем следования по дугам, помеченным буквой λ и выходящим из вершин множества S . Более точно, состояние j автомата $NFA(p)$ назовем λ -*достижимым* из состояния i , если выполняются следующие условия:

- состояние (вершина) i имеет выходную дугу, помеченную буквой λ ;
- в $NFA(p)$ существует ориентированный путь от вершины i до вершины j , дуги которого (пути) помечены как $\lambda \epsilon^k$ для некоторого целого $k \geq 0$;
- состояние j не имеет выходных дуг, помеченных ϵ .

Данные условия полностью определяют функцию $trans(S, \lambda)$. Расширим это определение для того, чтобы включить два особых случая.

- В случае, когда $S = \emptyset$ (пустое множество), функция $trans(S, \lambda)$ интерпретируется как $trans(\{0\}, \lambda)$, где 0 — метка начального состояния автомата $NFA(p)$.
- Обозначение $trans(S, \epsilon)$ будем понимать как множество всех состояний, достижимых из множества S , если следовать по выходящим из этого множества дугам, помеченным пустой строкой ϵ .

Для автомата, показанного на рис. 11.3, приведем для примера несколько значений функции *trans*.

$$\begin{aligned} trans(\emptyset, a) &= trans(\{3\}, a) = trans(\{12\}, \varepsilon) = \emptyset. \\ trans(\{2\}, a) &= trans(\emptyset, \varepsilon) = trans(\{5\}, \varepsilon) = \{2, 4, 7\}. \end{aligned}$$

Приведенный ниже алгоритм 11.1.1 использует функцию *trans* для вычисления в данной строке x всех позиций i , для которых имеет место совпадение подстроки $x[i'..i]$ с непустым регулярным выражением p . Доказательство корректности этого алгоритма оставим для упражнения 11.1.8, доказательство базируется на корректности реализации функции *trans* и на свойствах автомата N DFA.

Алгоритм 11.1.1. (Алгоритм N DFA)

```

▷ Нахождение всех подстрок строки  $x$ ,
▷ совпадающих с заданным непустым регулярным выражением  $p$ 
построение N DFA( $p$ ) с поглощающим состоянием  $w$ 
 $S \leftarrow trans(\emptyset, \varepsilon)$ 
for  $i \leftarrow 1$  to  $n$  do
     $S \leftarrow trans(S, x[i])$ 
    if  $w \in S$  then
        output  $i$ 
     $S \leftarrow trans(S, \varepsilon)$ 
    if  $w \in S$  then
        output  $i$ 
    
```

Для того чтобы определить сложность алгоритма N DFA, очевидно, сначала надо определить сложность функции *trans*. Эта функция выполняется $2n + 1$ раз, поочередно используя в качестве аргументов ε и $x[i]$. Функцию *trans* можно реализовать с использованием двух очередей, требующих памяти порядка $O(|V|)$ и бинарного отображения, требующего памяти порядка $\Theta(|V|)$, где $|V| \in \Theta(m)$ — количество состояний (вершин) автомата N DFA(p), m — количество букв в выражении p . Записав множество состояний N DFA в одну очередь, нетрудно вычислить переходы (не более двух переходов) из этих состояний и записать состояния, полученные в результате этих переходов, в другую очередь. Бинарное отображение необходимо для того, чтобы все состояния, записанные во вторую очередь, были различны. Поскольку во второй очереди надо хранить не более $|V|$ состояний и поскольку вычисляется не более двух переходов для каждого состояния, то функция *trans* будет выполняться за время порядка $O(|V|) = O(m)$. Напомним, что автомат N DFA(p) строится за время порядка $\Theta(m)$. Таким образом, доказан следующий основной результат данного подраздела.

Теорема 11.1.1. Для непустого регулярного выражения p , содержащего m букв, и для заданной строки $x = x[1..n]$ алгоритм 11.1.1 вычисляет все подстроки строки x , совпадающие с паттерном p , за время порядка $O(mn)$ с использованием памяти объемом $\Theta(m)$. ■

С помощью “четырёх русских” методов [20] можно уменьшить время выполнения алгоритма N DFA до величины порядка $O(mn/\log n)$, но за счет увеличения объема используемой памяти до величины порядка $O(mn/\log n)$ [188].

11.1.2 Детерминированные конечные автоматы

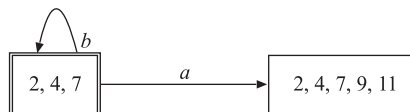
В этом подразделе мы рассмотрим способ построения детерминированного конечного автомата (сокращенно, DFA — от deterministic finite automaton) для заданного регулярного выражения p , т.е. такого конечного автомата, где не будет “пустых” переходов, а из каждого состояния можно будет перейти только в одно состояние, соответствующее какой-либо букве алфавита A . Такой автомат DFA(p) можно использовать, как мы увидим ниже, для вычисления всех совпадений с паттерном p заданной строки x , причем в большинстве случаев это можно сделать более эффективно, чем посредством автомата N DFA(p). Далее будем предполагать, что алфавит A индексирован (раздел 4.1).

Построение DFA(p) для заданного регулярного выражения p основано на простой идее: исключить из автомата N DFA(p) пустые переходы путем объединения вершин, которые можно достичь посредством пустых переходов.

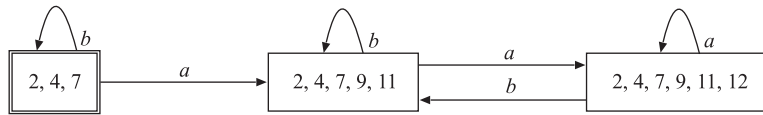
Для примера рассмотрим N DFA, показанный на рис. 11.3. Начиная с начального состояния 0 следуем по всем дугам, помеченным ϵ , до тех пор, пока не встретится вершина, из которой выходит “непустая” дуга — в данном случае это вершины 2, 4 и 7. Таким образом, начальной вершиной автомата DFA($(a | b)^*ab^*a$) будет



Теперь найдем в N DFA все состояния, достижимые из состояний 2, 4, 7 путем a -перехода и последующих пустых переходов. Это состояния 2, 4, 7, 9 и 11, которые объединим в новую вершину:



Здесь также показано, что b -переход из вершин 2, 4, 7 приводит к тому же множеству вершин. Продолжая этот процесс, получим автомат DFA($(a | b)^*ab^*a$) следующего вида.



Такие построения можно выполнить на основе любого N DFA, получив в результате DFA, в котором каждой вершине соответствует уникальное множество состояний, достижимых из ранее определенного множества состояний для новой входной буквы. Процесс построения завершается, когда у любой вершины не будет исходящих дуг, требующих определения нового множества состояний, отличного от ранее определенных множеств. Любая вершина DFA, включающая поглощающее состояние автомата N DFA, также называется поглощающим состоянием. Нетрудно показать, что описанный процесс построения DFA должен обязательно завершиться и что этот DFA должен иметь хотя бы одно поглощающее состояние. На основе рассмотренного примера можно сделать предположения, что для любого регулярного выражения p

- количество вершин в $DFA(p)$ имеет порядок $O(m)$, где m — количество букв в выражении p ;
- в $DFA(p)$ существует в точности одно поглощающее состояние.

Однако, как показано в упражнении 11.1.10, оба эти предположения неверны! Рассмотрим, например, регулярное выражение

$$p = (a | b)^* a (a | b)^{(m-3)/2},$$

описывающее все строки, определенные на алфавите $\{a, b\}$, в которых после буквы a стоят в точности $m' = (m - 3)/2$ произвольных букв (включая букву a). Для $m' = 2$ паттерн p будет иметь одно совпадение со строкой $x = abbab$, два совпадения со строкой $x = abaab$ и три совпадения с $x = aaabb$. Чтобы подсчитать количество букв, следующих за последней буквой a , в автомате DFA необходимо создать $\Theta(2^{m/2})$ различных вершин, многие из которых будут содержать поглощающее состояние автомата N DFA, и вследствие этого такие вершины DFA сами будут поглощающими состояниями. Как показано в упражнении 11.1.23, в самом худшем случае количество вершин в DFA может иметь порядок $\Theta(2^m)$.

Описанный процесс построения DFA не труден для понимания, однако его непросто реализовать программно. Вот одна из трудностей, возникающая при попытке практической реализации этого процесса. Каждый новый переход в автомате DFA определяет вершину, которая помечается набором целых чисел (эти числа принадлежат некоему множеству из $m + 1$ чисел). Необходимо определить, встречался ли текущий набор чисел ранее (как метка уже созданных вершин DFA) или нет, и если не встречался, то его надо сохранить. Нетрудно присвоить каждой вершине, помеченной набором чисел $\{s_1, s_2, \dots, s_t\}$, $1 \leq t \leq m + 1$, уникальный целочисленный идентификатор вида

$$S = 2^{s_1} + 2^{s_2} + \dots + 2^{s_t},$$

задающий бинарное отображение меток и требующий для хранения $O(m)$ бит. Таким образом, на каждом шаге построения DFA идентификатор S новой вершины надо сравнивать с $O(2^m)$ идентификаторами ранее созданных вершин.

Будем предполагать, что идентификаторы S принимают значения из интервала от 1 до 2^{m+1} , — всегда можно так перенумеровать метки автомата N DFA (в общем случае, за время, пропорциональное m), чтобы это условие выполнялось. Теперь можно определить массив $I[1..2^{m+1}]$, в котором любой элемент $I[S]$ или принимал бы значение NULL, или был бы указателем на вершину с идентификатором S автомата DFA. Итак, для хранения массива I требуется память объемом $\Theta(2^m)$, каждый элемент $I[S]$ этого массива можно изменить за константное время, а поиск нужного идентификатора в массиве I можно выполнить за время порядка $O(2^m)$. Заметим также, что для каждой вершины DFA необходима дополнительная память, пропорциональная α (размеру алфавита), поскольку каждая вершина имеет не более α исходящих дуг. Для заполнения такой памяти необходимо время порядка $O(\alpha)$ в расчете на одну вершину. Таким образом, мы доказали следующую теорему.

Теорема 11.1.2. Для заданного регулярного выражения p , содержащего m букв, автомат DFA(p) можно построить за время порядка $O(2^m \alpha)$ с использованием памяти объемом $O(2^m \alpha)$. ■

Как показано в [175], автомат DFA можно построить непосредственно на основе синтаксического дерева, однако такой способ не уменьшает границы сложности процесса построения DFA, приведенные в теореме 11.1.2.

Когда DFA(p) построен, нахождение всех подстрок строки $x = x[1..n]$, совпадающих с паттерном p , считается уже тривиальной задачей. Назовем алгоритмом DFA алгоритм, решающий эту задачу на основе построенного автомата DFA. В этом алгоритме должна использоваться функция, подобная функции *trans*, которая бы выполняла переходы в DFA и возвращала текущую позицию строки x , для которой достигнуто поглощающее состояние. Для индексированного алфавита любой переход можно выполнить за константное время, для неиндексированного алфавита для этого необходимо время порядка $O(\log \alpha)$ (как для дерева суффиксов). Имеет место следующая теорема.

Теорема 11.1.3. На основе автомата DFA(p), построенного для непустого регулярного выражения p , для заданной строки $x = x[1..n]$ (p и x определены на алфавите A) можно вычислить все подстроки строки x , совпадающие с паттерном p , за время порядка $O(n \log \alpha)$ (порядка $O(n)$, если алфавит A индексирован) с использованием памяти фиксированного объема. ■

Приведенную оценку эффективности для алгоритма DFA в случае индексированного алфавита можно доказать с помощью *переходной матрицы* $T =$

$= T[1..M, 1..\alpha]$, где M — количество вершин в $DFA(p)$. Введение этой матрицы требует переопределения вершин числами от 1 до M вместо идентификаторов S , введенных при построении DFA. Каждый элемент $T[j, h]$ переходной матрицы равен метке (номеру) вершины DFA(p), которую можно достигнуть из вершины j при предъявлении на вход автомата буквы $h \in 1..\alpha$.

Даже для неиндексированного алфавита алгоритм DFA находит отдельные совпадения с регулярными выражениями чрезвычайно быстро, особенно в ситуациях, когда с одним паттерном надо сравнить несколько текстовых строк (поскольку построенный автомат DFA можно использовать повторно).

Однако наши дифирамбы в адрес алгоритма DFA несколько не обоснованы: если в процесс сравнения с паттерном включить процесс построения DFA, то общее время выполнения такого алгоритма составит время порядка $O(2^m \alpha + n \log \alpha)$ с использованием памяти объемом $O(2^m \alpha)$ — весьма сомнительное улучшение аналогичных оценок $O(mn)$ и $O(m)$ алгоритма N DFA. Но есть и хорошие новости: случай, когда количество вершин в DFA имеет порядок $O(2^m)$, является патологическим. В большинстве случаев (как в рассмотренном выше примере) количество вершин в DFA имеет порядок $O(m)$. Поэтому временные и пространственные границы для алгоритма DFA можно уменьшить до величин порядка $O(m\alpha + n)$ и $O(n\alpha)$ соответственно. Кроме того, в [6] предложено несколько подходов, применение которых уменьшает временные и пространственные показатели вычислений на основе DFA. Один из таких подходов основан на “отложенных вычислениях переходов”, когда возможные переходы вычисляются не в процессе построения DFA, а, по мере необходимости, в процессе сравнения с паттерном. Такой подход, согласно [2], уменьшает общее время построения и использования DFA до величины порядка $O(m + n)$ даже в самом худшем случае.

11.1.3 Модификация алгоритма ВМ

В этом подразделе мы покажем, как для решения задачи сравнения с регулярным выражением можно применить методику использования бинарных векторов алгоритма ВМ (раздел 10.4), где эта методика позволила эффективно вычислить все k -приближенные совпадения текстовой строки x и паттерна p . Для применения этой методики в данном случае положим, что в N DFA(p) каждое состояние представимо одним битом¹ и если вершины помечены числами от 0 до m' , то каждой вершине соответствует бинарный вектор, содержащий $m' + 1$ бит. Напомним (подраздел 11.1.1), что если паттерн p содержит m букв, то $m' \in \Theta(m)$. Здесь, как и в алгоритме ВМ, должны выполняться некоторые условия: алфавит A должен быть конечным, известным заранее и индексированным. В этом случае можно считать, что $A = \{1, 2, \dots, \alpha\}$.

¹Этот бит должен только показать, активно ли данное состояние. — *Примеч. ред.*

В новой модификации алгоритма ВМ (назовем модифицированный алгоритм алгоритмом ВМ*) построим бинарные векторы $s_q[1..m']$, $q = 0, 1, \dots, k$, которые будут почти полным аналогом векторов (10.34). Единственное отличие состоит в том, что теперь j -й элемент в s_q определяет не j -ю позицию в p , а j -ю метку в $\text{N DFA}(p)$. Конечно, если ищется точное совпадение (тогда $k = 0$), то в этом случае будет только один бинарный вектор s_0 .

Аналогично переопределим бинарный массив $t = t[1..\alpha, 0..m']$, положив, что $t[h, j] = 0$ только тогда, когда дуга, направленная из вершины j в вершину $j + 1$, имеет меткой букву h . Теперь для любой позиции i строки x выражение

$$s_q^{(i-1)}[j - 1] \vee t[x[i], j] \tag{11.3}$$

в формуле (10.37) будет вычислено. Возвращаясь к рис. 11.3, мы видим, что выражение (11.3) корректно вычисляется для вершин 2, 3 и буквы a , для вершин 4, 5 и буквы b и т.д. Более того, из этого рисунка видно, что эту корректность для непустых переходов можно обеспечить путем присваивания последовательных меток конечным вершинам непустых дуг.

Однако для обеспечения корректности вычислений для всех вершин необходимо учесть еще две возможные ситуации.

1. Если в p входит метасимвол $|$, тогда в $\text{N DFA}(p)$ должна существовать по крайней мере одна пара вершин с последовательными метками $(j, j + 1)$, между которыми *нет* перехода (правило 3 построения N DFA). На рис. 11.3 есть такая пара вершин — это вершины (3, 4). Такую ситуацию легко формализовать, если ввести специальный шаблон, запрещающий для подобного случая переход из j в $j + 1$.
2. В $\text{N DFA}(p)$, как правило, много пар вершин (j_1, j_2) , соединенных пустыми дугами (т.е. дугами, помеченными пустой строкой ϵ), причем среди них обязательно есть пары вершин с последовательными метками, когда $j_2 = j_1 + 1$. Но вхождение в p метасимволов $|$ и $*$ обязательно порождает дуги (j_1, j_2) , для которых $j_2 \neq j_1 + 1$, просто как результат применения правил 3 и 4 построения N DFA . На рис. 11.3 это дуги (1, 4), (3, 6), (8, 11) и (10, 8).

Ниже опишем механизм обращения с пустыми дугами независимо от того, соединяют они последовательные вершины или нет.

Существование пустой дуги (j_1, j_2) просто показывает, что если имеет место приближенное совпадение, соответствующее вершине j_1 , то это же утверждение справедливо и для вершины j_2 . Напомним (см. подраздел 11.1.1), что в N DFA существует два типа вершин:

- вершины, из которых исходит одна или две пустые дуги (первый тип вершин);
- вершины, из которых исходит одна непустая дуга или которые вообще не имеют исходящих дуг (второй тип вершин).

Нам надо все вершины j_1 первого типа заменить вершинами j_2 второго типа. Для этого можно проследить пути от вершин первого типа, проходящие по пустым дугам и заканчивающиеся в вершинах второго типа, — правила построения NDFA предоставляют такую возможность. Этот процесс можно представить в виде функции $f : J_1 \rightarrow J_2$, отображающей множество J_1 вершин первого типа на множество J_2 вершин второго типа. Поскольку элементами множеств J_1 и J_2 являются несовпадающие целые числа из интервала от 0 до m' , то функцию f рационально представить как отображение бинарных векторов длиной $m' + 1$ на множество подобных бинарных векторов.

Чтобы реализовать такое отображение эффективно, в алгоритме ВМ* на предварительном этапе строятся массивы T_r , которые представляют все возможные переходы $J_1 \rightarrow J_2$. Каждый такой массив соответствует одному байту (восьми последовательным битам) в некотором векторе s_q . Всего $\lceil m'/8 \rceil$ таких массивов T_r . Тогда $r = 0, 1, \dots, \lceil m'/8 \rceil - 1$, что соответствует позициям

$$8r..8r + 7 = 0..7, 8..15, \dots, 8(\lceil m'/8 \rceil - 1)..m'$$

векторов s_q (вершинам NDFA(p)). Каждый такой байт представляет целое число b из интервала от 0 до 255. Для каждого значения $b \in 0..255$ в некоторой позиции $j_1 \in 0..7$ этого байта стоит 0, показывая, что вершина $j_1 + 8r$ является текущей активной вершиной NDFA(p). Другими словами, значение 0 в позиции j_1 указывает на то, что вершине $j_1 + 8r$ соответствует q -приближенное совпадение некоторой подстроки строки x и регулярного выражения p . Тогда в бинарном векторе $T_r[b]$ длиной $m' + 1$ в каждой позиции $j_2 \in 0..m' + 1$ стоит значение 0 только в том случае, если в NDFA(p) вершина j_2 является вершиной второго типа, достижимой из одной из вершин $j_1 + 8r$ посредством пути, состоящего только из пустых дуг. Каждый массив T_r является двумерным массивом, содержащим $2^8(m' + 1)$ байт; его можно вычислить за время порядка $\Theta(m')$. Тогда общее время вычисления всех $m'/8$ массивов T_r составляет величину порядка $\Theta((m')^2/8)$.

Обозначим через $s_{q,r}$ значение r -го байта вектора s_q . С помощью T_r для любого $r \in 0..\lceil m'/8 \rceil$ можно вычислить выражение

$$T_r[s_{q,r}] \wedge s_q, \tag{11.4}$$

где \wedge обозначает операцию логического умножения AND. Результатом вычисления (11.4) всегда будет значение 0, кроме случая, когда оба операнда в (11.4) равны 1. Поэтому активная в настоящий момент вершина $j_1 + 8r$ будет оставаться активной до тех пор, пока все вершины j_2 второго типа, достижимые из вершины $j_1 + 8r$, также будут активными. Поскольку операцию AND можно выполнить сразу над всем машинным словом, то $\lceil m'/8 \rceil$ вычислений выражения (11.4) на каждом шаге алгоритма требуют времени порядка $\Theta((m')^2/8w)$, где, как и ранее, w — длина машинного слова. Всего в процессе вычислений по формуле (10.37)

в алгоритме ВМ выражение (11.4) будет вычисляться kn раз. Отметим, что совпадение с паттерном будет зарегистрировано, если будет выполняться равенство $s_k[m'] = 0$.

В статье [233] доказано, что время выполнения алгоритма ВМ* имеет порядок

$$O(k \lceil 2m' / \log_2 n \rceil \lceil m' / w \rceil n).$$

Если значение m' не слишком велико, а значение n достаточно велико, тогда можно считать, что время выполнения алгоритма ВМ* имеет порядок $O(kn)$ с дополнительным временем порядка $\Theta((m')^2)$ для вычислений, выполняемых на предварительном этапе алгоритма.

Упражнения 11.1

1. Постройте $\text{N DFA}(p)$ для регулярного выражения (11.1).
2. Пусть в регулярное выражение p входит m букв. Покажите, что в этом выражении не может быть
 - более $2m$ скобок,
 - более m метасимволов $|$,
 - более $2m$ метасимволов $*$.
 Затем докажите, что количество $|p|$ символов в выражении p удовлетворяет неравенству $|p| \leq 6m$.
3. Основываясь на предыдущем упражнении, докажите неравенства (11.2). Можно ли изменить правила П.1.–П.4. так, чтобы уменьшить правые части этих неравенств?
4. Для автомата N DFA , показанного на рис. 11.3, докажите равенство $\text{trans}(\{5\}, \varepsilon) = \{2, 4, 7\}$. Правильно ли, что в последнее множество не включены состояния 1 и 6?
5. На рис. 11.3 видно, что N DFA может содержать циклы. Несмотря на это докажите, что $s \notin \text{trans}(s, \varepsilon)$ для любого состояния s , не являющегося начальным. Справедливо ли такое же утверждение для начального состояния?
6. Запишите алгоритм для функции $\text{trans}(S, \lambda)$ и докажите его корректность. Обратите особое внимание на случай, когда вторым аргументом функции будет пустая строка ε !
7. Покажите на примере, что если выражение $(\lambda_1 | \lambda_2 | \dots | \lambda_\alpha)^*$ не предшествует регулярному выражению p , то алгоритм 11.1.1, работающий с N DFA , построенным без учета этого выражения, может не обнаружить совпадения некоторых подстрок строки x с паттерном p .
8. Докажите корректность алгоритма 11.1.1.

9. Докажите, что в синтаксическом дереве, соответствующем некоторому регулярному выражению, понятия “внутренний узел” и “узел, имеющий метку в виде метасимвола”, эквивалентны.
10. Постройте синтаксические деревья и автоматы N DFA и DFA для следующих регулярных выражений.
 - $p = (a | b)^* aba$;
 - $p = ab^{m-1}$;
 - $p = a(a | b)(a | b)$;
 - $p = a(a | b | c)(a | b | c)$.

Основываясь на последних двух примерах, подсчитайте количество вершин в DFA(p), если $p = \lambda_1(\lambda_1 | \lambda_2 | \dots | \lambda_\alpha)^{(m-1)/\alpha}$. Сколько поглощающих состояний в этом DFA?

11. Докажите, что процесс преобразования N DFA в DFA должен обязательно завершиться.
12. Докажите, что в любом автомате DFA имеется хотя бы одно поглощающее состояние.
13. Покажите, что метки вершин в DFA совпадают с метками только тех вершин N DFA, которые имеют непустые исходящие дуги. Это утверждение не относится к вершинам DFA, соответствующим поглощающим состояниям. На основе этого утверждения докажите, что максимально возможное число вершин в DFA равно 2^{m+1} .
14. Запишите алгоритм преобразования N DFA в DFA.
15. Запишите алгоритм DFA нахождения всех совпадений подстрок строки x с регулярным выражением p с временем выполнения порядка $\Theta(n)$ для алфавита, который
 - упорядочен, но не индексирован;
 - индексирован.

Для этого измените функцию *trans* так, чтобы она для заданного идентификатора S вершины и буквы λ вычисляла следующую вершину $trans(S, \lambda)$.

16. Запишите алгоритм выполнения предварительного этапа алгоритма ВМ*.

11.2 Алгоритмы сравнения с множественными паттернами

В этом разделе рассмотрим решение задачи сравнения заданной текстовой строки $x = x[1..n]$ с множественным паттерном $P = \{p_1, p_2, \dots, p_r\}$. Поскольку

множественный паттерн P можно записать в форме $p_1 \mid p_2 \mid \dots \mid p_r$, то данную задачу можно рассматривать как задачу сравнения с регулярными выражениями специального типа; алгоритмы для решения этой задачи представлены в предыдущем разделе. В этом разделе мы опишем более эффективные алгоритмы решения данной задачи, которые будут учитывать специфику множественных паттернов.

Алгоритмы для точного сравнения с паттерном описаны в подразделах 11.2.1 и 11.2.2. Это модификации алгоритма КМП (алгоритм Кнута–Морриса–Пратта, раздел 7.1) и алгоритма БМ (алгоритм Бойера–Мура, раздел 7.2). В двух последних подразделах описаны алгоритмы для приближенного сравнения с множественными паттернами: в подразделе 11.2.3 представлена модификация алгоритма ВМ (алгоритм Ву–Менбера, раздел 10.4), а в подразделе 11.2.4 — новый (в этой книге) алгоритм Бейза–Ятса–Наварро (Baeza-Yates-Navarro, сокращенно — алгоритм ВУН). Новая модификация алгоритма ВМ применима к функциям расстояния любых типов, а алгоритм ВУН — только для расстояний преобразования, Левенштейна и Хемминга.

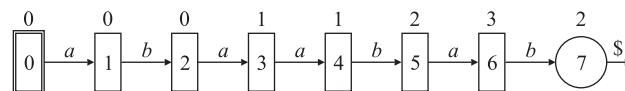
11.2.1 Автомат Ахо–Корасика: модификация алгоритма КМП

Теперь, когда мы познакомились с конечными автоматами, читателю должно быть ясно, что алгоритм КМП можно представить в виде конечного автомата. Действительно, в алгоритме КМП сначала сравнивается очередной входной символ, и если такого совпадения не наблюдается, тогда выполняется сравнение с символами наибольшей грани, соответствующей текущей входной позиции. Другими словами, мы уже имеем конечный автомат, построенный по правилам алгоритма КМП построения массива граней.

Чтобы понять “природу” конечного автомата Ахо–Корасика (сокращенно, АСФА — Aho-Corasick finite automaton) [3], рассмотрим небольшой пример. Предположим сначала, что имеем простой паттерн $p = abaabab$, который сравнивается с текстовой строкой $x = x[1..n]$. Для данного паттерна p получаем следующий массив граней (раздел 1.3):

$$\begin{array}{cccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\
 p = & a & b & a & a & b & a & b \\
 \beta = & 0 & 0 & 1 & 1 & 2 & 3 & 2
 \end{array}$$

Тогда автомат АСФА(p) будет иметь вид



В ACFA(p) метки *внутри* вершин обозначают позиции символов в паттерне p , метки на дугах — символы паттерна, определяющие переход к следующей вершине (состоянию). Вспоминая определение деревьев суффиксов в разделе 2.1, здесь мы видим, что ACFA(p) является ничем иным, как некомпактным синтаксическим деревом для паттерна p . Добавим в это дерево метки *над* вершинами, которые соответствуют элементам массива граней и указывают вершину (состояние), в которую должен перейти автомат в случае *несовпадения* входного символа с меткой исходящей дуги. (С этой точки зрения массив граней можно назвать “функцией несовпадения”!) Таким образом, эти метки в действительности определяют “обратные” дуги (переходы) от текущей вершины (состояния) в одну из предыдущих вершин, что соответствует “текущему” значению массива граней. В случае несовпадения в начальной вершине (вершине 0) для определенности положим “возврат” в эту же вершину.

Как и в случае деревьев суффиксов, примем соглашение, что в ACFA(p) исходящие дуги конечных вершин помечены специальным символом \$, который не может совпадать ни с одной буквой входной строки. Это обеспечит корректное поведение автомата после установления факта совпадения входной строки и паттерна.

В общем случае, поскольку массив граней можно вычислить за линейное время (теоремы 1.3.2 и 1.3.3), автомат ACFA(p) для любого одиночного паттерна $p = p[1..m]$ можно построить за время порядка $\Theta(m)$ с использованием памяти объемом $\Theta(m)$. Когда такой автомат построен, алгоритм ACFA (представленный ниже как алгоритм 11.2.1) без труда найдет все совпадения паттерна p и строки x . Алгоритм ACFA можно рассматривать как вариант алгоритма КМП, однако на самом деле он является зеркальным отражением алгоритма 1.3.1 вычисления массива граней.

Алгоритм 11.2.1. (Алгоритм ACFA)

```

▷ Вычисление позиций  $i$  в строке  $x$ , таких, что подстрока  $x[i'..i]$  в ACFA
▷ определяет путь от состояния 0 до поглощающего состояния
 $s \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n$  do
    while  $s > 0$  and not  $match(s, x[i])$  do
         $s \leftarrow \beta[s]$ 
    if  $match(s, x[i])$  then
         $s \leftarrow s + inc(s, x[i])$ 
        if  $accept(s)$  then output  $(i, s)$ 

```

Выполнение алгоритма ACFA зависит от четырех типов вычислений.

- Массив граней β вычисляется на предварительном этапе построения автомата NDFA. Поэтому для любого состояния s значение $\beta[s]$ можно получить за константное время. Напомним, что мы положили $\beta[0] = 0$.
- В случае обнаружения совпадения текущее значение s должно возрасти на величину, определяемую буквой $x[i]$, для которой зафиксировано совпадение. В нашем простом примере шаг возрастания значения s всегда равен 1, однако для множественных паттернов шаг уже не будет постоянным. Значение шага будет храниться в метке исходящей из вершины s дуги, которая соответствует букве $x[i]$. В этом случае значение шага можно получить за константное время с помощью функции *inc*.
- Каждое состояние содержит бит, который указывает, является ли данное состояние поглощающим. Значение этого бита возвращает функция *accept*.
- Функция *match* возвращает значение истины тогда, когда имеет место совпадение метки на дуге, исходящей из вершины s , и буквы $x[i]$. В нашем примере с одним паттерном эта функция выполняет одно буквенное сравнение и поэтому вычисляется за константное время. Однако в случае множественных паттернов, как мы увидим далее, вычисление функции *match* усложняется.

Расширим приведенный выше пример и рассмотрим множественный паттерн, состоящий из трех простых паттернов:

$$P = \{p_1, p_2, p_3\} = \{abaabab, abab, ababcabab\}. \tag{11.5}$$

Применяя к этому паттерну ту же методологию построения некомпактного синтаксического дерева, получим автомат ACFA(P), показанный на рис. 11.4. Отметим, что общее количество *всех* дуг в ACFA(P) (т.е. включая как явно показанные на рисунке дуги, так и “обратные” дуги, определяемые метками над вершинами) не превышает величины $2m + 1$, где m — сумма длин частных паттернов в множественном паттерне P .

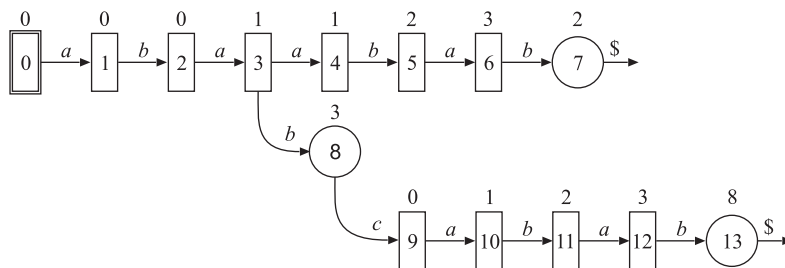


Рис. 11.4. Автомат ACFA(P) для паттерна (11.5)

На этом примере очевидна необходимость внесения изменений в процесс вычисления обычного массива граней для того, чтобы пересчитать перемещение некоторых позиций в паттернах с учетом номеров (меток) состояний автомата $ACFA(P)$. Так позиции 4 в паттернах $abab$ и $abaabab$ соответствует состояние 8, позиции 9 в паттерне $ababcabab$ — состояние 13, а $\beta[13] = 8$. Оставим подробности этого процесса для упражнения 11.2.2. В этом же упражнении рассмотрен случай, когда один паттерн является гранью другого паттерна, как и в паттерне (11.5).

Имея массив β , корректно вычисленный для автомата $ACFA(P)$, нетрудно показать, что алгоритм 11.2.1 корректно определяет все подстроки строки x , совпадающие с паттернами множественного паттерна P . Цикл **while** гарантирует, что в случае обнаружения несовпадения с некоторым префиксом $p[1..j]$ некоторого паттерна p далее сравнения с гранями подстроки $p[1..j - 1]$ выполняются в убывающем порядке длин этих граней, точно так же как в алгоритме КМП. Таким образом, корректность алгоритма $ACFA$ для множественных паттернов следует непосредственно из корректности алгоритма КМП. Более того, если функцию $match$ можно вычислить за константное время, тогда время выполнения алгоритма 11.2.1 будет пропорциональным величине n . Это утверждение вытекает из того факта, что общее количество итераций цикла **while**, просуммированное по всем n этапам алгоритма, не может превышать величины $n - 1$. Такой же факт имеет место для цикла **while** процесса вычисления массива граней (см. раздел 1.3)! Таким образом, предполагая, что функция $match$ вычисляется за константное время, линейность времени выполнения алгоритма 11.2.1 доказывается на основе тех же аргументов, что и утверждение о линейности времени выполнения алгоритма вычисления массива граней (теорема 1.3.3).

Что же на самом деле можно сказать о времени вычисления функции $match$? В общем случае автомат $ACFA(P)$ является (слегка приукрашенным) некомпактным синтаксическим деревом, построенным на множестве $P = \{p_1, p_2, \dots, p_r\}$ различных строк длиной m_1, m_2, \dots, m_r соответственно. (Общую длину этих строк обозначим m , т.е. $m = m_1 + m_2 + \dots + m_r$.) Если алфавит имеет размер α , то в $ACFA(P)$ может быть несколько вершин, имеющих $O(\alpha)$ исходящих дуг. Поэтому для корректного выбора исходящей дуги в каждой вершине необходимо хранить данные объемом $O(\alpha)$. Тогда выбрать правильную дугу можно за константное время, если алфавит индексирован, либо за время порядка $O(\log \alpha)$, если алфавит упорядочен, но не индексирован. Подобная ситуация для деревьев суффиксов исследована в обсуждении 2.1.1.

Даже когда частные паттерны длинные, а их количество относительно невелико, процесс корректного выбора дуги может быть трудоемким. Например, предположим, что $r = \alpha$ и паттерны p_h имеют префиксы λ_h^2 , $h = 1, 2, \dots, \alpha$. Тогда в $ACFA(P)$ вершина 0 имеет α исходящих дуг. Пусть в строке $x = x[1..n]$ нет квадратов отдельных букв, а любая буква алфавита встречается n/α раз. При сравнении этой строки с паттерном P вершина 0 будет посещаться n раз, и при

каждом ее посещении необходимо выбрать одну дугу среди α дуг. Если алфавит не индексирован, то для определения нужной дуги необходимо использовать или дерево поиска, или массив соответствия. Таким образом, только обработка вершины 0 потребует времени порядка $O(n \log \alpha)$.

Наши рассуждения обобщим в виде следующей теоремы.

Теорема 11.2.1. Пусть $P = \{p_1, p_2, \dots, p_r\}$ — множество паттернов общей длиной m , $x = x[1..n]$ — заданная текстовая строка, строка x и паттерны P определены на алфавите размером α . Тогда для хранения автомата ACFA(P) необходима память объемом $\Theta(\alpha m)$ и

- для индексированного алфавита A автомат ACFA(P) строится за время порядка $\Theta(m)$, а алгоритм 11.2.1 выполняется за время порядка $\Theta(n)$;
- для упорядоченного, но не индексированного алфавита A автомат ACFA(P) строится за время порядка $\Theta(m \log \alpha)$, а алгоритм 11.2.1 выполняется за время порядка $\Theta(n \log \alpha)$. ■

Таким образом, автомат Ахо–Корасика позволяет сравнить текстовую строку $x[1..n]$ с множеством паттернов общей длиной m за время порядка $\Theta(m + n)$ для индексированных алфавитов и за время порядка $\Theta((m + n) \log \alpha)$ для упорядоченных алфавитов. Другими словами, алгоритм КМП можно распространить на случай множественных паттернов практически без потери эффективности вычислений.

Если в алгоритме ACFA использовать ориентированные ациклические графы слов (см. раздел 5.3), тогда этот алгоритм в самом худшем случае потребует выполнения не более $2n$ буквенных сравнений. Такая модификация алгоритма ACFA наиболее эффективна тогда, когда размер алфавита небольшой, а все частные паттерны длинные.

11.2.2 Автомат Комменца–Вальтера: модификация алгоритма БМ

Так же как алгоритм ACFA является модификацией алгоритма КМП для случая множественных паттернов, так и алгоритм CWFA (Commentz–Walter Finite Automata — конечный автомат Комменца–Вальтера) является модифицированным вариантом алгоритма БМ [57]. В алгоритме CWFA существенно использование тех же массивов δ_1 и δ_2 , что и в алгоритме БМ (раздел 7.2). Различие в использовании этих массивов заключается только в том, что в алгоритме CWFA сдвиги в этих массивах вычисляются как соответствующие минимумы по позициям во *всех* паттернах множественного паттерна P , а не как минимумы по позициям одного паттерна. Чтобы сделать это различие более ясным, рассмотрим автомат CWFA(P) для множественного паттерна (11.5), введенного в предыдущем под-

разделе. Как показано на рис. 11.5, состояния в автомате CWFA соответствуют позициям в строках, записанным в обратном порядке. Такой конечный автомат соответствует алгоритму БМ, в котором сравнение букв текстовой строки с буквами паттерна происходит при перемещении по паттерну справа налево. Напомним, что в алгоритме БМ несовпадение букв в процессе сравнения требует вычисления нового значения i и последующего сравнения буквы $x[i]$ с самым правым символом паттерна — на рис. 11.5 это соответствуют возврату в вершину 0 после каждого несовпадения, установленного в процессе сравнения.

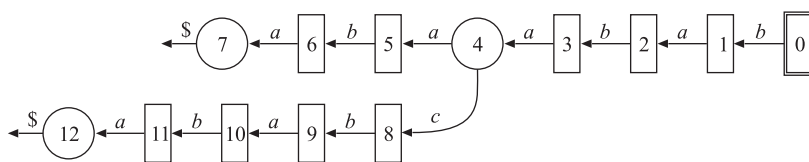


Рис. 11.5. Автомат CWFA(P) для паттерна (11.5)

Предположим, что на вход автомата CWFA(P), показанного на рис. 11.5, подается текстовая строка $x = \dots \underline{bab} \dots$. (Здесь подчеркивание указывает на буквы, совпадающие с паттерном.) Для определенности положим, что $x[i' + 1..i' + 3] = bab$. Тогда на дуге, исходящей из вершины 3, будет иметь место несовпадение $x[i'] = b \neq a$. Так как имеет место несовпадение, в стандартном алгоритме БМ вычисляются δ_2 -сдвиги, отдельно для каждого паттерна

$$abab, abaabab, ababcabab \tag{11.6}$$

из множественного паттерна P . Получим следующие значения для новой позиции i : $i = i' + 5$, $i = i' + 8$, $i = i' + 9$, каждое из которых соответствует “своему” паттерну из (11.6). В алгоритме CWFA выбирается минимальный сдвиг, т.е. новой позицией будет $i' + 5$. Аналогично входная буква c в вершине 0 определяет δ_1 -сдвиги $i = i' + 4$, $i = i' + 6$, $i = i' + 4$, соответствующие паттернам (11.6), среди которых выбирается минимальное значение $i = i' + 4$.

Рассмотрим общий случай, когда множественный паттерн $P = \{p_1, p_2, \dots, p_r\}$ задан как множество r частных паттернов длиной соответственно m_1, m_2, \dots, m_r , общая длина этих паттернов равна $m = m_1 + m_2 + \dots + m_r$. Значения массивов δ_1 и δ_2 вычисляются отдельно для каждого паттерна из P , затем определяется минимум по r паттернам. Это означает, что для каждого состояния s находится минимум по не более чем r паттернам, которые имеют общий префикс со строкой, определяемой путем, ведущим к s . Поскольку значения массивов δ_1 и δ_2 для одного паттерна можно вычислить за время, пропорциональное длине паттерна (см. упражнение 7.2.4 и теорему 7.2.2), то сдвиги для множественного паттерна P

можно вычислить за время порядка $O(m)$. Технические детали этого процесса мы опускаем.

Далее, в алгоритме БМ процедуру $(i, j) \leftarrow hctam(i, m)$ необходимо заменить новой процедурой $(i, s) \leftarrow hctam(i, 0)$, которая, увеличивая значение i , осуществляет перемещение по пути от вершины 0 до некоторой вершины s , в которой будет иметь место несовпадение (такое несовпадение обязательно будет хотя бы потому, что исходящие дуги конечных вершин помечены специальным символом $\$$, которого по условию нет в текстовых строках). Функции $hctam$ мы также делегируем операцию вывода пары значений (i, s) в том случае, когда будет достигнуто поглощающее состояние s (т.е. когда будет достигнуто состояние, для которого $accept(s) = 0$). Таким образом, алгоритм CWFA (представлен ниже как алгоритм 11.2.2) формально почти полностью идентичен алгоритму БМ.

Алгоритм 11.2.2 (Алгоритм CWFA)

▷ Вычисление всех позиций i в строке x , таких, что
 ▷ подстрока $x[i..i']$ в автомате $CWFA(P)$ определяет путь
 ▷ от состояния 0 до некоторого поглощающего состояния,
 ▷ при этом $P = \{p_1, p_2, \dots, p_r\}$
 $i \leftarrow \min\{m_1, m_2, \dots, m_r\}$
while $i \leq n$ **do**
 $(i, s) \leftarrow hctam(i, 0)$
 $i \leftarrow i + \max\{\delta_1[x[i], s], \delta_2[s]\}$

Отметим, что поскольку δ_2 -сдвиг положителен для любого паттерна p , то значение i возрастает на каждом шаге алгоритма не менее, чем на 1. Это гарантирует завершение выполнения алгоритма 11.2.2.

Подобно автомату ACFA(P), автомат CWFA(P) требует для хранения памяти объемом $\Theta(\alpha m)$ и может быть построен за время $\Theta(m)$, если алфавит индексирован, или за время $\Theta(m \log \alpha)$, если алфавит только упорядочен. Как и алгоритм БМ, алгоритм CWFA имеет сложности при работе с периодическими паттернами, поэтому в самом худшем случае его время выполнения может иметь порядок $\Theta(mn)$. Как показано в [2], алгоритм CWFA в среднем работает быстрее, чем алгоритм ACFA, в том случае, если количество r частных паттернов невелико, и медленнее в случае большого значения r .

11.2.3 Аппроксимирующие паттерны: модификация алгоритма БМ

Рассмотрим задачу k -приближенного сравнения заданной текстовой строки $x = x[1..n]$ с множеством паттернов $P = \{p_1, p_2, \dots, p_r\}$, которые имеют длины

соответственно m_1, m_2, \dots, m_r . В этом подразделе опишем модификацию алгоритма Ву–Менбера (алгоритм ВМ), который эффективно решает эту задачу [233].

Основная идея модификации алгоритма ВМ — обрабатывать множественный паттерн $p = p_1 p_2 \dots p_r$ как единый паттерн длиной $m = m_1 + m_2 + \dots + m_r$, при этом, конечно, необходимы средства от возможных ошибок, связанных с границами паттернов.

Напомним, что в алгоритме ВМ для каждого $q \in 0..k$ и $i \in 1..n$ вычисляется бинарный вектор $s_q^{(i)}$, где для всех $j \in 1..m$

$$s_q^{(i)}[j] = 0 \Leftrightarrow c[i, j] \leq q,$$

а $c[i, j]$ — элемент матрицы стоимостей (раздел 10.1), вычисляемый для подстроки $x[i'..i]$ и префикса $p[1..j]$ паттерна. В алгоритме ВМ также используется бинарный массив t , где для каждого $h \in 0..\alpha$ и $j \in 1..m$

$$t[h, j] = 0 \Leftrightarrow p[j] = h.$$

Для того чтобы определить границы между паттернами, в модифицированном алгоритме применим шаблон

$$M = 01^{m_1-1}01^{m_2-1} \dots 01^{m_r-1}$$

длиной m , который будет использоваться для вычисления

$$t'[h, 1..m] = t[h, 1..m] \vee M,$$

где, как и ранее (раздел 10.4), \vee обозначает операцию логического сложения OR. Тогда $t'[h, j] = 0$ только в том случае, когда j будет первой позицией одного из паттернов p_1, p_2, \dots, p_r и при этом $p[j] = h$.

Первый бинарный вектор $s_0^{(i)}[j]$ вычисляется по стандартным формулам (10.37):

$$s_q^{(i)}[j] \leftarrow s_{q-I}^{(i-1)}[j] \wedge s_{q-D}^{(i)}[j-1] \wedge s_{q-S}^{(i-1)}[j-1] \wedge (s_q^{(i-1)}[j-1] \vee t[x[i], j]).$$

При $q = 0$ эти вычисления корректны для всех позиций j , за исключением тех позиций j_0 , с которых начинаются частные паттерны. Чтобы избежать ошибок в этих позициях, необходимо положить $s_0^{(i)}[j_0] \leftarrow 0$, если $x[i] = p[j_0]$, и $s_0^{(i)}[j_0] \leftarrow 1$ в противном случае. Добиться такого результата можно с помощью вычисления $s_0^{(i)}[j] \wedge t'[x[i], j]$ для всех $j \in 1..m$, где \wedge обозначает операцию логического умножения AND. В результате таких вычислений битовые значения в позициях, не совпадающих с начальными позициями паттернов, не изменятся, а значения в начальных позициях паттернов изменятся (с 1 на 0) только в том случае, если $t'[x[i], j_0] = 0$.

При $q = 1$, чтобы избежать подобных ошибок в процессе вычисления $s_1^{(i)}[j]$, необходимо положить $s_1^{(i)}[j_0] \leftarrow 0$ для начальных позиций j_0 всех паттернов из P . (Здесь просматривается аналогия со сдвигом 0 в $s_1^{(i)}[1]$, выполняемым процедурой *rightshift* в формуле (10.38) исходного алгоритма ВМ.) Итак, после вычислений по формуле (10.37) для всех $j \in 1..m$ следует выполнить дополнительные вычисления

$$s_1^{(i)}[j] \leftarrow s_1^{(i)}[j] \wedge M_1,$$

где $M_1 = M$. Эти дополнительные вычисления гарантируют, что первые биты для каждого частного паттерна будут равны 0, тогда как значения остальных битов останутся вычисленными по формуле (10.37).

В общем случае после вычисления по формуле (10.37) необходимо выполнить дополнительные вычисления

$$s_q^{(i)}[j] \leftarrow s_q^{(i)}[j] \wedge M_q,$$

где $M_q = 0^q 1^{m_1 - q} 0^q 1^{m_2 - q} \dots 0^q 1^{m_r - q}$.

Дополнительные вычисления, необходимые для коррекции векторов $s_q^{(i)}$, можно выполнить за константное время, и поскольку эти дополнительные вычисления заключаются в выполнении логических операций, то их можно реализовать применительно к целым машинным словам, а не к отдельным битам, как показано здесь. Следовательно, сложность модифицированного алгоритма ВМ не увеличилась, и мы имеем право сформулировать теорему, которая является аналогом теоремы 10.4.1.

Теорема 11.2.2. Пусть заданы текстовая строка $x[1..n]$ и множество паттернов $P = \{p_1, p_2, \dots, p_r\}$, которые определены на индексированном алфавите размером α . Обозначим $m = |p_1| + |p_2| + \dots + |p_r|$. Тогда алгоритм ВМ вычисляет все k -приближенные совпадения множественного паттерна P и строки x за время порядка $\Theta(k \lceil m/w \rceil n)$ с использованием памяти объемом порядка $\Theta((\alpha + k) \lceil m/w \rceil)$, где w — длина машинного слова. ■

В работе [186] описан алгоритм сравнения с множественным аппроксимирующим паттерном, в котором скомбинированы бинарное отображение и методы хеширования. Этот алгоритм, по-видимому, будет очень быстрым, однако он применим только для расстояния преобразования и только для аппроксимационной константы k , равной 1.

11.2.4 Аппроксимирующие паттерны: алгоритм Бейза-Ятса-Наварро

Мы закончим этот раздел и данную главу описанием еще одного алгоритма, использующего бинарные операции, который первоначально предлагался автора-

ми Бейзом-Ятсом и Наварро (Baeza-Yates, Navarro) для решения задачи приближенного сравнения с одним паттерном [26] и только позже был модифицирован для приближенного сравнения с множественными паттернами [27]. Этот алгоритм (назовем его алгоритмом BYN — от Baeza-Yates–Navarro), примененный к множественным паттернам, является примером “фильтрационного” алгоритма, о котором упоминалось в подразделе 10.3.2. В таких алгоритмах сначала очень быстро определяется *сверхмножество* подстрок текстовой строки, которые приблизительно совпадают хотя бы с одним паттерном, затем из этого множества с помощью также очень быстрой технологии отфильтровываются все ложные решения. Подобно всем алгоритмам, основанным на бинарных операциях, в алгоритме BYN предполагается, что текстовые строки и паттерны определены на индексированном алгоритме $1, 2, \dots, \alpha$ (раздел 4.1).

Описание алгоритма BYN начнем со случая приближенного сравнения с одним паттерном, где будет использоваться расстояние преобразования (т.е. будет решаться задача k -разностей, описанная в разделе 10.3). Затем покажем, как этот алгоритм можно применить к множественным паттернам. Важно отметить, что методология алгоритма BYN не ограничивается только расстоянием преобразования. Как будет показано, этот алгоритм можно представить в виде модели конечного автомата, в котором будут “видны” операции вставки, удаления и подстановки. Поэтому данный алгоритм после соответствующей “обработки” этих операций будет применим и к расстояниям Левенштейна и Хемминга.

Автомат $\text{BYNFA}(p, k)$

Алгоритм BYN для k -приближенного сравнения заданной строки x и одного паттерна p можно представить в виде конечного автомата “нестандартного” вида — назовем этот автомат $\text{BYNFA}(p, k)$. Как показано на рис. 11.6, автомат BYNFA для паттерна $p = \text{this}$ и аппроксимационной константы $k = 2$ представим в виде прямоугольной решетки с $k + 1$ строкой и $m + 1$ столбцом (m — длина паттерна). Без потери общности можно считать, что $m \geq k$. Строка h ($h \in 0..k$) решетки соответствует h -разностям между входной текстовой строкой и паттерном (подсчитанным на основе расстояния преобразования). Столбец j ($j \in 0..m$) определяет приближенное сравнение с префиксом $p[1..j]$ паттерна. Пересечения (h, j) строк и столбцов решетки интерпретируются как состояния автомата. Если состояние (h, j) *активно*, тогда имеет место h -приближенное совпадение между $p[1..j]$ и подстрокой $x[1..i]$ (i — текущая позиция в строке x). Поэтому первоначально все состояния в столбце $j = 0$ будут активными и, кроме того, все состояния (h, j) , для которых $0 \leq j \leq h \leq k$, также будут активны, что соответствует h -несовпадениям между первыми k позициями паттерна и пустыми префиксами строки x . Первоначально активные состояния на рис. 11.6 помечены меткой 0.

Как и в других конечных автоматах, изученных нами, здесь дуги определяют переходы из одного состояния в другое, определяемое текущей (i -й) входной

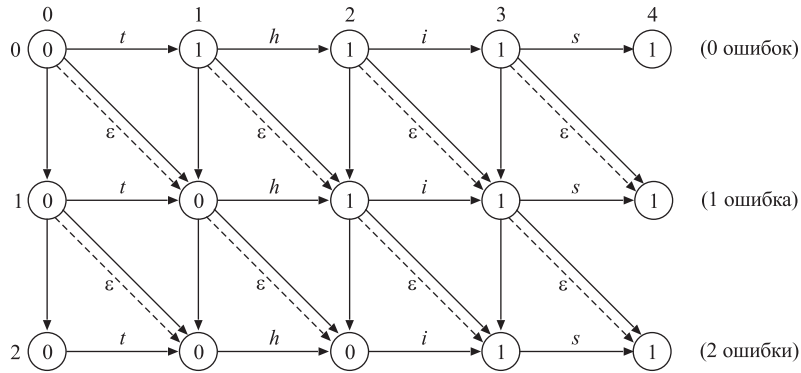


Рис. 11.6. Автомат BYNFA(*this*, 2)

буквой строки x . Состояния будут активными или неактивными в соответствии со следующими правилами.

- Первоначально активное состояние остается таковым, если из этого состояния начинается сравнение с любой позицией текстовой строки.
- Любое состояние, достижимое из другого активного состояния в соответствии с текущей буквой строки x , становится активным.
- Любое состояние, не являющееся начальным и не достижимое в соответствии с текущей буквой строки x , становится неактивным.

В автомате BYNFA все переходы между состояниями можно отнести к одному из четырех типов. Один тип переходов оставляет без изменения уровень ошибки h , а три других типа увеличивают этот уровень на единицу.

1. Горизонтальный переход $(h, j - 1) \rightarrow (h, j)$ соответствует совпадению буквы $x[i]$ и буквы $p[j]$, которой помечен этот переход. Значение h остается неизменным, поскольку новая буква $x[i]$ не увеличивает текущую разность между x и p , вычисляемую на основе расстояния преобразования.
2. Вертикальный переход $(h - 1, j) \rightarrow (h, j)$ становится возможным тогда, когда для достижения совпадения текущую букву $x[i]$ необходимо вставить в паттерн p в позицию j . При этом позиция j паттерна не изменяется, а значение разности h увеличивается на единицу.
3. Диагональный переход $(h - 1, j - 1) \rightarrow (h, j)$ по сплошной линии соответствует подстановке вместо $p[j]$ буквы $x[i]$, при этом позиции в паттерне p и строке x сдвигаются на одну позицию вправо, а значение разности h увеличивается на единицу.
4. Диагональный переход $(h - 1, j - 1) \rightarrow (h, j)$ по пунктирной линии (помеченной ϵ) становится возможным тогда, когда для достижения совпадения

текущую букву $p[j]$ можно удалить из паттерна. Это “пустой” или автоматический переход, поскольку здесь удаляется буква из паттерна без каких-либо изменений в текстовой строке x .

Отметим, что пустые диагональные переходы гарантируют активность состояний $(h + q, j + q)$ для всех целых $q \in 0..k - h$, если будет активно состояние (h, j) . И наоборот, состояние (h, j) становится активным, если активно какое-либо состояние $(h - q, m - q)$, $0 \leq q \leq k$. Таким образом, активность состояния (h, m) указывает на то, что имеет место h -приближенное совпадение $p[1..j]$ с некоторой подстрокой строки x , заканчивающейся в текущей позиции i .

Чтобы лучше понять операции, выполняемые автоматом BYNFA, рассмотрим два варианта, когда на вход автомата $\text{BYNFA}(\text{this}, 2)$ поступают строки $x = \text{hash}$ и $x = \text{whist}$. В табл. 11.1 и 11.2 представлены результаты обработки этих строк; здесь показаны состояния, которые становятся активными после обработки каждой входной буквы (но сюда не включены начальные активные состояния). Отметим, что пара чисел (i, h) будет выходным результатом тогда, когда существует наименьшее целое число h , такое, что для входной буквы $x[i]$ становится активным состояние (h, m) .

Таблица 11.1. Обработка автоматом $\text{BYNFA}(\text{this}, 2)$ строки $x = \text{hash}$

i	$x[i]$	Активные состояния	Совпадения (i, h)
1	h	$(1, 2), (2, 3)$	
2	a	$(2, 3)$	
3	s	$(2, 4)$	$(3, 2)$
4	h	$(1, 2), (2, 3)$	

Таблица 11.2. Обработка автоматом $\text{BYNFA}(\text{this}, 2)$ строки $x = \text{whist}$

i	$x[i]$	Активные состояния	Совпадения (i, h)
1	w		
2	h	$(1, 2), (2, 3)$	
3	i	$(1, 3), (2, 3), (2, 4)$	$(3, 2)$
4	s	$(1, 4), (2, 4)$	$(4, 1)$
5	t	$(0, 1), (1, 2), (2, 3)$	

Компактное представление автомата BYNFA

Нетрудно заметить, что для организации работы автомата $\text{BYNFA}(p, k)$ необходимо отслеживать изменение его состояний. Для этих целей введем бинарный

массив $A = A[0..k, 0..m]$, в котором для всех $h \in 0..k$ и $j \in 0..m$ $A[h, j] = 0$ только в том случае, если состояние (h, j) активно, и $A[h, j] = 1$ в противном случае. Сначала положим $A[h, j] = 0$ для всех h и j , таких, что $0 \leq j \leq h \leq k$. Обозначим через $A^{(0)}$ начальное состояние автомата BYNFA, а через $A^{(i)}$ ($i = 1, 2, \dots, n$) — последующие состояния этого автомата, соответствующие входным буквам $x[1], x[2], \dots, x[n]$.

Если известны состояния BYNFA для некоторого значения $i - 1$, тогда для следующего значения i и для всех h и j , таких, что $0 \leq h \leq j \leq m$, новые состояния пересчитываются по формуле

$$A^{(i)}[h, j] \leftarrow (A^{(i-1)}[h, j - 1] \vee t[x[i], j]) \wedge A^{(i-1)}[h - 1, j] \wedge A^{(i-1)}[h - 1, j - 1] \wedge A^{(i)}[h - 1, j - 1], \quad (11.7)$$

где

- бинарный массив $t = t[1..\alpha, 1..m]$ вычисляется заранее точно так, как в алгоритмах ДБГ (раздел 7.4) или ВМ (подраздел 11.2.3), т.е. для каждой буквы $r \in 1..\alpha$ и для всех $j \in 1..m$ $t[r, j] = 0$ только в том случае, если $p[j] = r$;
- символ \vee обозначает операцию логического сложения OR (результатом этой операции будет значение 0 только тогда, когда *оба* операнда равны 0);
- символ \wedge обозначает операцию логического умножения AND (результатом этой операции будет значение 0 тогда, когда *хотя бы один* операнд равен 0).

Отметим, что для корректного пересчета $A^{(i)}[h, j]$ на основе значения $A^{(i)}[h - 1, j - 1]$, стоящего на той же диагонали решетки, значения h и j должны идти в возрастающем порядке.

Чтобы при пересчете состояний автомата BYNFA можно было использовать логические операции, выполняемые параллельно над всеми битами машинного слова, необходимо выполнить некоторые условия. Так для хранения состояний используем массив A , состоящий из $(k + 1)(m + 1)$ бит. Но, как показано в упражнении 11.2.9, для каждой позиции i строки x в этом массиве необходимо пересчитывать только $(k + 1)(2m - k)/2 \in \Theta(mk)$ бит. Кроме того, как показано выше, наличие пустых переходов позволяет установить активные состояния, просто перемещаясь по диагоналям решетки от некоторого активного состояния. Поэтому для пересчета состояний на диагоналях необходимо хранить в памяти для каждой диагонали только номер строки, где находится активное состояние.

Отметим также, что благодаря эффекту пустых переходов никакое состояние (h, j) , $j > k + h$, не может стать активным, если не будет активным состояние (k, m) . Следовательно, пересчет состояний на диагоналях можно начинать с состояний $(0, j)$, $j = 1, 2, \dots, m - k$. Такое упрощение позволяет определить все позиции i строки x , где имеет место приближенное совпадение, однако не позволяет определить соответствующий уровень ошибки h , как это сделано в табл. 11.1 и 11.2.

Для работы с диагональными состояниями определим одномерный массив $D = D[0..m - k + 1]$, в котором j -й элемент соответствует диагонали, начинающейся в состоянии $(0, j)$. Значением элемента $D[j]$ является целое число из интервала $0..k + 1$. Если $D[j] = k + 1$, то это означает, что на j -й диагонали нет активных состояний. Во всех других случаях значение $D[j]$ равно наименьшему номеру строки, где есть активное состояние на j -й диагонали. Обозначим через $D^{(0)}$ первоначальное состояние массива D : здесь $D^{(0)}[0] = D^{(0)}[m - k + 1] = 0$, а для остальных элементов $j \in 1..m - k$ $D^{(0)}[j] = k + 1$. На основании формулы (11.7) получаем следующее рекуррентное соотношение для значений массива $D^{(i)}$, соответствующего позиции i строки x ($j \in 1..m - k$):

$$D^{(i)}[j] \leftarrow \min\{D^{(i-1)}[j + 1] + 1, D^{(i-1)}[j] + 1, g(i - 1, j - 1)\}, \quad (11.8)$$

где $g(i - 1, j - 1)$ — минимальный номер строки h_{\min} на диагонали $j - 1$, такой, что

- $D^{(i-1)}[j - 1] \leq h_{\min} \leq k$ (другими словами, $A[h_{\min}, j - 1]$ активно для $x[i - 1]$),
- $t[x[i], j - 1 + h_{\min}] = 0$ (это означает, что $x[i] = p[j - 1 + h_{\min}]$).

Если такого номера h_{\min} не существует, тогда $g(i - 1, j - 1) = k + 1$. В правой части формулы (11.8) первое выражение соответствует вертикальному переходу в BYNFA(p, k) (вставке буквы $x[i]$), второе выражение — диагональному переходу (подстановка $x[i]$ вместо $p[j]$), третье — горизонтальному переходу (имеет место совпадение буквы $x[i]$ и некоторой буквы паттерна, стоящей в позиции, соответствующей активному состоянию на диагонали $j - 1$). Таким образом, формула (11.8) определяет все активные состояния каждой диагонали j ($j \in 1..m$) на основе входной буквы $x[i]$ и на знании наименьших (по номеру строки) активных состояний на диагоналях $j + 1$, j и $j - 1$, соответствующих букве $x[i - 1]$.

Отметим, что h -приближенное совпадение паттерна p с некоторой подстрокой $x[i'..i]$ строки x будет иметь место только тогда, когда $D^{(i)}[m - k] \leq k$ (здесь, конечно, $h \in 0..k$).

Алгоритм BYN для одного паттерна

Обычное компьютерное представление неотрицательных целых значений элементов $D[j]$, которые по величине не превышают $k + 1$, требует $\lceil \log_2(k + 2) \rceil$ бит, а всего для массива $D[1..m - k]$ — $\lceil \log_2(k + 2) \rceil(m - k)$ бит. Но, к сожалению, такое представление элементов массива D не ведет к эффективному алгоритму вычислений. По этой причине в алгоритме BYN для элементов $D[j]$ используется представление вида $0^{k+1-D[j]}1^{D[j]}$. Тогда значение $D[j] = 0$ представимо как 0^{k+1} , $D[j] = 2$ — как $0^{k-1}1^2$, $D[j] = k + 1$ — как 1^{k+1} и т.д. Если ввести дополнительно нулевой разделяющий бит слева от каждого значения $D[j]$, тогда весь

массив $D[1..m - k]$ можно представить в виде битовой строки

$$D = 0 \ 0^{k+1-D[1]} 1^{D[1]} 0 \ 0^{k+1-D[2]} 1^{D[2]} \dots 0 \ 0^{k+1-D[m-k]} 1^{D[m-k]}$$

длиной $(k + 2)(m - k)$ бит. Примерно такой же объем памяти необходим для хранения массива A .

Имея такое представление элементов $D[j]$, легко реализовать в виде бинарных операций вычисления, выполняемые над этими элементами в формуле (11.8), используя для этого такие замены.

- O1.** Вычисление $D[j] + 1$ заменяется на $leftshift(D[j], 1) \vee 0^k 1$.
- O2.** Вычисление $\min\{D[j_1], D[j_2]\}$ заменяется на $D[j_1] \wedge D[j_2]$.
- O3.** Получение доступа к $D[j \pm 1] - leftshift(D[j + 1], k + 2)$ или $rightshift(D[j - 1], k + 2)$.

Для того чтобы реализовать функцию g , сначала надо записать в обратном порядке для каждой буквы $r \in 1..\alpha$ строку бинарного массива t :

$$t'[r] = t[r, m]t[r, m - 1] \dots t[r, 1].$$

Затем опять для каждой буквы r зададим шаблон для массива $D[1..m - k]$ вида

$$T'[r] = 0rightshift(t'[r], 0)0rightshift(t'[r], 1) \dots 0rightshift(t'[r], m - k - 1).$$

Делая обычное предположение о том, что после реализации любого сдвига, независимо от того, вправо был сделан этот сдвиг или влево, в позиции, “освобожденные” сдвигом заносятся нули, наконец можем представить все $m - k$ вычислений по формуле (11.8) (когда j меняется от 1 до $m - k$) в виде одной формулы бинарных операций, выполняемых над машинным словом, в котором записан массив D :

$$\begin{aligned} D &= (leftshift(D, k + 3) \vee (0^{k+1} 1)^{m-k-1} 0 \ 1^{k+1}) \wedge \\ &\quad \wedge (leftshift(D, 1) \vee (0^{k+1} 1)^{m-k}) \wedge \\ &\quad \wedge rightshift((D' + (0^{k+1} 1)^{m-k}) \wedge D', 1) \wedge \\ &\quad \wedge D^{(0)}[1..m - k], \end{aligned} \tag{11.9}$$

где $D^{(0)}[1..m - k] = (0 \ 1^{k+1})^{m-k}$ и $D' = rightshift(D, k + 2) \vee T[x[i]]$.

Первые три строки в формуле (11.9) соответствуют трем выражениям в правой части формулы (11.8), однако здесь эти выражения записаны с помощью бинарных операций O1–O3. Массив $D^{(0)}$ в четвертой строке формулы введен для того, чтобы восстановить разделяющие нули между элементами массива D .

Текущий номер i буквы строки x будет выводиться как результат выполнения алгоритма только тогда, когда после вычисления по формуле (11.9) будет выполняться равенство

$$D \wedge 0^{(k+2)(m-k-1)} 0 \ 10^k = 0,$$

поскольку в этом случае $D^{(i)}[m - k] \leq k$. Затем, чтобы не было “ложных срабатываний” на следующих этапах алгоритма, значение $D[m - k]$ очищается с помощью операции

$$D \leftarrow D \vee 0^{(k+2)(m-k-1)} 0 10^{k+1}.$$

Вычисления по формуле (11.9) слишком сложны для быстрого анализа, поэтому оставим для упражнения 11.2.10 получение точной оценки сложности алгоритма BYN. В работе [26] в качестве меры скорости выполнения алгоритма предложена форма цикла с перескоками (раздел 8.1): текстовая строка x сканируется для обнаружения совпадения с одной из первых $k + 1$ букв паттерна и только после обнаружения такого совпадения начинает работу автомат BYNFA.

До сих пор мы описывали автомат BYNFA в предположении, что длина компьютерного слова достаточна для хранения всего массива D , т.е. предполагали, что $w \geq (k + 2)(m - k)$. Если это предположение выполняется, тогда на каждом из n шагов алгоритма BYN требуется только константное время для выполнения бинарных операций над одним машинным словом и, следовательно, алгоритм завершится за время порядка $\Theta(n)$.

Декомпозиция BYNFA

В общем случае, когда $w < (k + 2)(m - k)$, задача распределения по машинным словам всей информации, необходимой для вычисления в автомате BYNFA, вследствие двухмерной структуры этого автомата становится более сложной, чем аналогичная задача для других подобных алгоритмов, основанных на бинарных отображениях (раздел 10.4). Для автомата BYNFA различают три ситуации, в каждой из которых используется свой подход к декомпозиции автомата.

1. Ситуация, когда $w \geq k + 2$.

В этой ситуации в одно машинное слово можно поместить значения элементов $\lceil w/(k + 2) \rceil$ диагоналей из общего числа $m - k$ диагоналей, формируя таким образом “подавтоматы”. Всего для автомата BYNFA потребуется $\lceil (m - k) / \lceil w/(k + 2) \rceil \rceil$ машинных слов. Поскольку каждый элемент $D[j]$ зависит от элементов $D[j - 1]$ и $D[j + 1]$, лежащих на соседних диагоналях, в алгоритме необходимо предусмотреть взаимодействие каждого подавтомата со своими соседями. Однако и в этом случае, как описано в работе [26] и показано в упражнении 11.2.11, каждое машинное слово обрабатывается за константное время. Поэтому общее время выполнения алгоритма имеет порядок $O(k(m - k)n/w)$.

2. Ситуация, когда $w < k + 2$, но $w \geq 2(m - k)$.

В этом случае для хранения значений элементов одной диагонали не хватает одного машинного слова. В такой ситуации можно попробовать сделать декомпозицию BYNFA на основе вертикальных компонентов, состоящих из $k + 1$ элементов. Однако, поскольку эти элементы все равно будут содержать

части $m - k$ диагональных значений, необходимо предусмотреть $m - k$ разделяющих битов. Следовательно, для хранения самих значений в машинном слове останется только $w - (m - k)$ бит. (Именно поэтому мы поставили условие, что $w \geq 2(m - k)$.) Итак, при вертикальной декомпозиции автомата для хранения одного значения, соответствующего одной диагонали, в машинном слове остается

$$\lfloor (w - (m - k)) / (m - k) \rfloor = \lfloor w / (m - k) \rfloor - 1$$

бит. Всего для полного представления автомата BYNFA необходимо $\lceil (k + 1) / (\lfloor w / (m - k) \rfloor - 1) \rceil$ машинных слов. Здесь, как и в предыдущей ситуации, в алгоритме также необходимо предусмотреть взаимодействие между отдельными компонентами, однако и в этом случае каждое машинное слово можно обработать за константное время. Поэтому, как и ранее, общее время выполнения алгоритма имеет порядок $O(k(m - k)n/w)$.

3. Ситуация, когда $w < k + 2$ и $w < 2(m - k)$.

Подробно данная ситуация рассмотрена в [26], где показано, что и в этом случае общее время выполнения алгоритма имеет порядок $O(k(m - k)n/w)$, а для хранения автомата BYNFA необходима память объемом $\Theta(k(m - k)n/w)$.

Напомним (см. обсуждение 1.3.1), что всегда можно положить $w \in \Omega(\log n)$. Поэтому справедлива следующая теорема.

Теорема 11.2.3. Для заданного одиночного паттерна $p = p[1..m]$ и неотрицательного целого числа $k < m$ алгоритм BYN на основе расстояния преобразования вычисляет все k -приближенные совпадения с текстовой строкой $x = x[1..n]$ за время порядка $O(k(m - k)n / \log n)$ с использованием памяти объемом $\Theta(k(m - k)n/w)$, где w — длина машинного слова. ■

В работе [26] предложено несколько стратегий для уменьшения времени выполнения алгоритма BYN и проведен детальный анализ временных и пространственных требований алгоритма для того, чтобы показать завышенность оценок в самом худшем случае, данных в теореме 11.2.3. В частности, доказано, что модифицированный алгоритм BYN в среднем выполняется за время порядка $O(n)$, если величина k/n является медианным значением некоторого распределения. Таким образом, мы видим, что и в самом худшем случае и в среднем временные границы алгоритма BYN (в случае приближенного сравнения с одним паттерном) превышают аналогичные границы для модифицированного алгоритма Майерса (подраздел 10.3.2).

Алгоритм BYN для множественного паттерна

Основная идея алгоритма BYN для множественного паттерна заключается в том, что сначала с помощью алгоритма BYN для одного паттерна вычисляется

сверхмножество всех k -приближенных совпадений заданной строки x с набором паттернов $P = \{p_1, p_2, \dots, p_r\}$, затем это сверхмножество эффективно фильтруется для получения искомого множества решений. При описании алгоритма BYN для множественного паттерна сделаем упрощающее предположение, что все паттерны имеют одинаковую длину m/r , где m — общая длина всех паттернов. Это предположение не ограничивает общности, поскольку для того, чтобы длина всех паттернов была одинаковой, всегда в конец любого паттерна можно добавить символы, которые не будут совпадать с буквами алфавита.

Для вычисления сверхмножества эффективным оказался способ совмещения (наложения) паттернов, когда, например, на первой итерации автомат $\text{BYNFA}(P, k)$ распознает совпадение не просто с первой буквой одного паттерна, а совпадение с любой из букв $p_1[1], p_2[1], \dots, p_r[1]$; или на j -итерации автомат распознает совпадение с любой из букв $p_1[j], p_2[j], \dots, p_r[j]$.

Для примера вместо автомата $\text{BYNFA}(\text{this}, 2)$, показанного на рис. 11.6, рассмотрим его расширение — автомат $\text{BYNFA}(\{\text{this}, \text{worm}\}, 2)$. Этот автомат показан на рис. 11.7. Если на вход данного автомата подается строка $x = \text{what}$, то будут обнаружены 2-совпадения при $i = 2$, хотя $d_E(\text{what}, \text{this}) = d_E(\text{what}, \text{worm}) = 3$. С другой стороны, автомат корректно найдет совпадение при $i = 4$ в слове $x = \text{wham}$ (расстояние преобразования между словом wham и паттерном worm равно 2), а также совпадения при $i = 3$ и $i = 4$ в слове $x = \text{whim}$ (расстояние преобразования между этим словом и обоими паттернами this и worm равно 2).

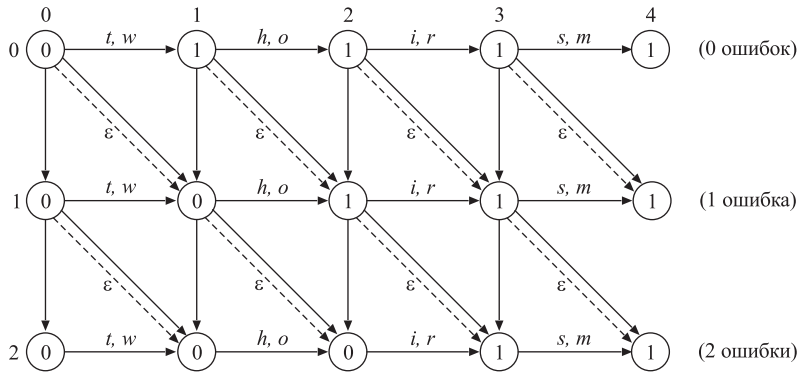


Рис. 11.7. Автомат $\text{BYNFA}(\{\text{this}, \text{worm}\}, 2)$

Такой подход для произвольной строки x позволяет определить в этой строке все подстроки, которые имеют k -совпадения хотя бы с одним частным паттерном из P . Кроме того, такой подход не сложно реализовать алгоритмически: для того чтобы на каждой итерации распознать, имеет ли место совпадение с несколькими

буквами, достаточно переопределить бинарный массив $t[1..a, 1..m]$ таким образом, чтобы равенство $t[h, j] = 0$ для *любой* буквы h выполнялось только тогда, когда эта буква присутствует в позиции j хотя бы в одном частном паттерне. Эти вычисления можно выполнить на предварительном этапе алгоритма так, как показано в подразделе “Сравнение множеств букв” раздела 10.4. Каждый новый пересчет массива t выполняется за константное время. Поэтому, согласно теореме 7.4.4, массив t можно вычислить за время порядка $\Theta(\alpha \lceil m/w \rceil + m)$, где, как и ранее, m — сумма длин всех r паттернов из P . После предварительного вычисления массива t сверхмножество решений задачи приближенного сравнения с множественным паттерном, как следует из теоремы 11.2.3, можно вычислить за время порядка $O(k(m/r - k)n/\log n)$.

Пусть $S = S(P, k, x)$ обозначает вычисленное алгоритмом BYN супермножество для множественного паттерна P и заданной строки x , k — аппроксимационная константа. Элементами множества S являются позиции в строке x , для которых автоматом BYNFA(P, k) обнаружены совпадения.

Для того чтобы отфильтровать элементы множества S , не дающие k -приближенное совпадение хотя бы с одним частным паттерном из P , необходимо подстроки строки x , определяемые каждым значением i_s ($s = 1, 2, \dots, q$), сравнить с каждым паттерном p_t ($t = 1, 2, \dots, r$). Если нет такого паттерна p_t , что вычисленное расстояние преобразования между такой подстрокой и этим паттерном равно k , то элемент i_s удаляется из множества S . Каждое такое сравнение на основе расстояния преобразования можно выполнить с помощью алгоритма Укконена–Майерса (раздел 9.5) за время порядка $O(km/r)$. Тогда на обработку q элементов множества S потребуется время порядка $O(kmq)$, а поскольку $q \leq n$, то время, необходимое для фильтрации этого множества, имеет порядок $O(kmn)$.

Эффективность такого подхода сравнима с эффективностью модифицированного алгоритма Ву–Менбера (подраздел 11.2.3), поскольку зависит как от величины q , так и от числа q' элементов множества S , удаляемых в процессе фильтрации этого множества. Если значение q ограничено малым целым числом, тогда время, необходимое для выполнения алгоритма Укконена–Майерса, уменьшается до величины порядка $O(km)$. Если же значение q' мало, то основное время выполнения алгоритма Укконена–Майерса будет потрачено на проверку элементов множества S , а не на их отбраковку. В работе [27] предложено несколько стратегий, уменьшающих величину q' .

Упражнения 11.2

1. В алгоритме ACFA функция *match* при некоторых значениях i может вызываться дважды. Перепишите этот алгоритм таким образом, чтобы исключить повторные вызовы функции *match*.

2. Как показано на рис. 11.4, в автомате ACFA(P), построенном для множественного паттерна, нумерация многих состояний может не совпадать с их естественной нумерацией, определяемой для одного конкретного паттерна. Объясните, как надо правильно назначать номера состояниям в процессе построения автомата ACFA(P).
3. Дано множество паттернов $P = \{p_1, p_2, \dots, p_r\}$. Разработайте эффективный алгоритм построения автомата ACFA(P) при выполнении *одного* из следующих условий:
 - алфавит индексирован;
 - алфавит упорядочен, но не индексирован.

В каждом случае оцените временную сложность вашего алгоритма и исследуйте эффективность построенного автомата ACFA.

4. Для каждого случая, определенного в предыдущем упражнении, разработайте подходящую функцию *match* и охарактеризуйте ее временные параметры.
5. Предположим, что в множестве $P = \{p_1, p_2, \dots, p_r\}$ паттерн p_1 является гранью паттерна p_2 . Тогда, как видно на рис. 11.4, поглощающему состоянию, определяемому паттерном p_1 , не будет соответствовать конечная вершина автомата ACFA(P). В этом случае, если алгоритм ACFA фиксирует в позиции i строки x совпадение с паттерном p_2 , он не может зафиксировать в этой же позиции совпадение с паттерном p_1 . Обсудите возможность и необходимость изменить алгоритм ACFA, чтобы он в подобных ситуациях фиксировал совпадение и с паттерном p_1 .
6. Напомним, что в алгоритме КМП (раздел 7.1) на предварительном этапе вычисляется вспомогательный массив граней β' . Однако более эффективен массив β'' , где j -й элемент равен наибольшей грани подстроки $p[1..j - 1]$, за которой *не следует* буква $p[j]$.
7. Покажите, что с помощью автомата ACFA(p) можно вычислить массив граней β , и запишите алгоритм (с линейным временем выполнения) для вычисления массива, подобного массиву β'' .
8. Напишите процедуру *hctam* для алгоритма CWFA.
9. (Исследовательский проект!) Для знакомого вам варианта алгоритма БМ (раздел 7.2) разработайте конечный автомат, подходящий для сравнения с множественным паттерном, и соответствующую процедуру сдвига. На этой основе создайте алгоритм сравнения с множественным паттерном. Оцените асимптотическую сложность этого алгоритма и сравните ее с аналогичными оценками алгоритмов ACFA и CWFA.

10. Покажите, что массив A , определенный в алгоритме BYNFA, содержит в точности $(k + 1)(2m - k)/2$ элементов, соответствующих первоначально неактивным состояниям.
11. Предположим, что длина машинного слова достаточна для хранения всего массива $D = D[1..m - k]$. Для этого случая запишите реализацию алгоритма BYN, выполняющего сравнение с одним паттерном.
12. Предположим, что длина машинного слова w удовлетворяет неравенствам $w < (k + 2)(m - k)$, $w \geq k + 2$. Для этого случая измените алгоритм BYN таким образом, чтобы время его выполнения имело порядок $O(k(m - k)n)$.
13. Покажите, как надо изменить алгоритм BYN, чтобы он выводил не только позицию i строки x , где имело место h -приближенное совпадение ($h \in 0..k$), но и само значение h . Возрастет ли сложность алгоритма после внесения в него таких изменений?
14. В тексте главы очень кратко говорилось о циклах с перескоками применительно к конечным автоматам, когда автомат начинает работу только с позиции i строки x , такой, что буква $x[i]$ совпадает хотя бы с одной буквой подстроки $p[1..k + 1]$. Используя идею циклов с перескоками, докажите, что k -приближенное совпадение строки x с паттерном p будет иметь место только тогда, когда по крайней мере одна из букв подстроки $p[1..k + 1]$ будет представлена в любой подстроке строки x .
15. Запишите алгоритм для вычисления бинарного массива $B[1..\alpha]$, в котором r -й элемент показывал бы, входит ли буква r ($r \in 1..\alpha$) в подстроку $p[1..k + 1]$ паттерна p . Какова сложность такого алгоритма?
16. Покажите, что при выполнении условия $m \leq 2(\sqrt{w} - 1)$ автомат BYNFA($p[1..m], k$) можно сохранить в одном машинном слове длиной w независимо от значения k . Из приведенного условия вычислите максимально возможное значение m для $w = 32$ и $w = 64$.

ЧАСТЬ IV

Вычисление характеристических паттернов

Мой отец до сих пор каждый день читает словарь. Он говорит, что наша жизнь зависит от нашего умения пользоваться словами.

— Артур Скарджилл (1938).
The Sunday Times, 10 января 1982 г.

Наше путешествие в мир алгоритмов анализа строковых последовательностей завершается изучением методов вычисления характеристических паттернов — таких паттернов, которые описывают некоторые заданные свойства строк, а не указывают, какие символы в них содержатся. Например, мы распознаем строку *aa* и строку *bcbc* как кратные строки, поскольку они обладают общим свойством; и нас не интересует, какие именно символы они содержат.

Характеристические паттерны, как правило, связаны со свойствами периодичности строковых последовательностей. Но, конечно, в качестве “характеристических” можно также рассматривать и другие свойства, не связанные с периодичностью. В главе 12 сначала (в разделе 12.1) будут показаны “классические” алгоритмы вычисления кратных строк в произвольной строковой последовательности, а затем (в разделе 12.2) будут описаны недавние алгоритмы, которые предлагают более компактное представление кратных строк в виде серий (определение 2.3.3).

В главе 13 рассматриваются более слабые формы периодичности. Первой из них является “квазипериодичность”: периодическая строка содержит последовательные вхождения образующей подстроки *u*, но в квазипериодической строке подстроки *u* могут также накладываться друг на друга. В соответствии с определением 2.3.3, мы говорим, что квазипериодическая строка “покрывается” подстрокой *u*. В разделе 13.1 описан алгоритм, который вычисляет все оболочки каждого префикса строки за линейное время. В разделах 13.2 и 13.3 обобщается идея кратных строк до понятия “раппорта”, в котором не только не требуется последовательности вхождений повторяющейся подстроки, но и они могут быть отделены друг от друга. Мы рассмотрим алгоритмы построения раппортов, как точных, так и приближенных, и затем, наконец, используем обобщенное понятие “приближенной периодичности”.

ГЛАВА 12

Периодичность

Благодаря словам мы смогли подняться над животными;
и благодаря же словам мы часто опускаемся до уровня де-
монов.

— Олдос Хаксли (1894–1963).
Адонис и Алфавит

В этой главе описаны алгоритмы для вычисления кратных строк в произвольной строковой последовательности $x = x[1..n]$. Прежде чем приступить к этому материалу, читателю необходимо повторить начало раздела 2.3, где представлена основная терминология и определены кодировки кратных строк в виде троек чисел (i, p^*, r^*) и серий (i, p^*, r^*, t) .

В конце 60-х годов было показано [157], что в любой строке x имеется порядка $O(n \log n)$ квадратов при условии, что образующие квадратов сами не являются кратными строками. В начале 80-х были разработаны три алгоритма [69, 18, 176], которые вычисляют все кратные строки в строковой последовательности x за время $O(n \log n)$. Два из них представлены в двух первых разделах этой главы. Но сначала покажем, что верхняя граница времени выполнения таких алгоритмов наиболее вероятна для строк, определенных на неупорядоченном алфавите, — результат, установленный в [176].

Теорема 12.0.4. Пусть $x = x[1..n]$ — строковая последовательность, определенная на неупорядоченном алфавите. Любой алгоритм, предназначенный для нахождения кратных строк в строке x и основанный на посимвольном сравнении, в самом худшем случае требует времени выполнения порядка $\Omega(n \log n)$.

Доказательство. Напомним (см. раздел 4.1), что символы из обычного алфавита могут сравниваться только для определения их равенства или неравенства. Предположим, что x — это строка, в которой каждый из n символов отличается от другого. Приняв для простоты, что n является четным числом, покажем, что для установления факта, что x не содержит квадратов и, следовательно, ни одной кратной строки, необходимо $\Omega(n \log n)$ посимвольных сравнений.

Пусть для некоторых целых i и j , таких, что $1 \leq i < j \leq n$, при сравнении букв $x[i]$ и $x[j]$ обнаружено, что $x[i] \neq x[j]$. В строке x единственными квадратами, которые могут исключаться по этому сравнению, будут те, для которых необходимо выполнение равенства $x[i] = x[j]$. Если обозначить $k = j - i$, то не более k таких квадратов, каждый длиной $2k$, исключаются:

$$x[i - k + 1..j], x[i - k + 2..j + 1], \dots, x[i..j + k - 1]. \quad (12.1)$$

Отметим, что в действительности может быть исключено менее k квадратов, поскольку, например, список исключенных квадратов, получаемый на основе неравенства $x[i + 1] \neq x[j + 1]$,

$$x[i - k + 2..j + 1], x[i - k + 3..j + 2], \dots, x[i + 1..j + k]$$

содержит только один квадрат, не включенный ранее в список (12.1).

Для любого целого числа $k \in 1..n/2$ существует (как отмечалось в разделе 1.2) всего $n - 2k + 1$ различных возможных подстрок длиной $2k$, следовательно, возможно $n - 2k + 1$ различных квадратов длиной $2k$. Таким образом, для исключения всех возможных квадратов длиной $2k$ необходимо сделать не менее $(n - 2k + 1)/k$ символьных сравнений.

Заметим также, что после исключения всех возможных квадратов длиной $2k$ мы не получим информации о наличии или отсутствии каких-либо квадратов длиной $2k'$, если $k' \neq k$, поскольку для определения того, является ли какая-либо подстрока длиной $2k'$ квадратом, необходимо сравнить пары символов, отделенные друг от друга на k' позиций. Таким образом, исключение квадратов различной длины $2k$, $k = 1, 2, \dots, n/2$, должно рассматриваться как независимые действия. Следовательно, чтобы быть уверенным в том, что в строке x нет квадратов, необходимо сделать не менее

$$\sum_{k=1}^{n/2} (n - 2k + 1)/k \geq (n + 1)H_{n/2} - n$$

символьных сравнений, где H_i — i -е гармоническое число¹. Поскольку $H_i > \ln i + \gamma$ [135], где $\gamma \approx 0,577$ — константа Эйлера, то отсюда находим, что необходимое

¹Гармоническое число H_i — это просто обозначение суммы $1 + 1/2 + \dots + 1/i$. — *Примеч. ред.*

количество символьных сравнений превышает величину

$$(n + 1)(\ln n - \ln 2 + \gamma) - n > (n + 1)(\ln n - 1,13) \in \Omega(n \log n). \quad \blacksquare$$

Этот результат является одним из основных при вычислении кратных строк, однако он порождает новые принципиальные вопросы. Например, такие.

- Справедлива ли нижняя граница $n \log n$ в тех случаях, когда алфавит не общего вида, а, например, упорядочен? В работе [70] Крочемор (Crochemore) описывает алгоритм, который на основе деревьев суффиксов (раздел 5.2) за время порядка $O(n \log \alpha)$ определяет, действительно ли данная строка $x[1..n]$, определенная на упорядоченном алфавите размером α , не содержит квадратов. В статье [170] Мейн и Лоренц (Main and Lorentz) предлагают алгоритм с подобным временем выполнения, но в котором не используются деревья суффиксов. Если $\alpha \in \Theta(n)$, тогда оценка временной сложности этих алгоритмов не является улучшением результата теоремы 12.0.1. С другой стороны, если α — это известная константа, то и $\log \alpha$ является константой, и, таким образом, эти алгоритмы можно рассматривать как алгоритмы, которые выполняются за время порядка $O(n)$. См. также подраздел 5.2.5.
- Можно ли для строковой последовательности $x[1..n]$, определенной на упорядоченном алфавите, вычислить ее кратные строки за время $O(n)$ или, по крайней мере, за время, меньшее чем $\Theta(n \log n)$? До некоторого времени казалось, что вряд ли можно на этот вопрос ответить утвердительно, поскольку, принимая во внимание результат Крочемора (см. [79] и раздел 3.4), строки Фибоначчи длиной f_n содержат порядка $\Theta(f_n \log f_n)$ кратных строк, закодированных в тройки (i, p^*, r^*) , которые по существу являются нормальными формами (раздел 1.2) любых максимальных повторяемых подстрок. Поэтому казалось, что для вычисления всех кратных подстрок строки длиной n необходимо время порядка $\Theta(n \log n)$. Но относительно недавно было введено понятие серии (раздел 2.3), кодирующей строки в виде кортежа из 4 элементов (i, p^*, r^*, t) [168], которое было использовано различными авторами [108, 123, 138]. Затем в [138, 139] было показано, что количество серий в любой строке линейно по отношению к длине строки. Таким образом, возникает возможность, по крайней мере для строк, определенных на известном упорядоченном алфавите, за линейное время вычислить кратные строки, закодированные в виде серий. В [139] для вычисления серий в строке x представлен алгоритм, в котором используются деревья суффиксов (следовательно, предполагается упорядоченность алфавита). Если дополнительно предположить, что алфавит не очень велик либо что он фиксирован или индексирован (раздел 4.1), то потенциальным фактором $\log n$, скрытым в конструкции дерева суффиксов, можно пренебречь. Тогда получим алгоритм со временем выполнения порядка $\Theta(n)$. Этот алгоритм мы рассмотрим

в разделе 12.2. Вопрос о существовании алгоритма вычисления серий для упорядоченных алфавитов со временем выполнения порядка $\Theta(n)$, который не использовал бы деревьев суффиксов, на сегодняшний день остается открытым.

12.1 Все кратные строки

В этом разделе представлены два алгоритма, значительно отличающиеся друг от друга, но которые с сопоставимой сложностью достигают одинаковой цели. Это алгоритм Крочемора [79] со временем выполнения $O(n \log n)$ и алгоритм Мейна и Лоренца [176] со временем выполнения $\Theta(n \log n)$, которые вычисляют все кратные строки в произвольной строке $x = x[1..n]$.

12.1.1 Алгоритм Крочемора

Прежде чем перейти к объяснению алгоритма Крочемора, вначале опишем простой алгоритм, который на каждом этапе достигает такого же результата, что и алгоритм Крочемора, а именно: в строке x он вычисляет все повторяющиеся подстроки длиной ℓ , где $\ell = 1, 2, \dots, n - 1$.

Рассмотрим строку Фибоначчи

$$\begin{array}{ccccccccccccccc}
 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\
 f_6 = & a & b & a & a & b & a & b & a & a & b & a & a & b
 \end{array}$$

Предположим, что на этапе (уровне) $\ell = 1$ в строке f_6 вычислены последовательности позиций, в которых встречаются одинаковые символы:

$$\begin{array}{ccc}
 \ell = 1 & \langle 1, 3, 4, 6, 8, 9, 11, 12 \rangle & \langle 2, 5, 7, 10, 13 \rangle \\
 & a & b
 \end{array}$$

Напомним (раздел 4.1), что для строк, определенных на неупорядоченном алфавите, эти вычисления можно выполнить за время порядка $\Theta(n^2)$, а для определенных на упорядоченных и индексированных алфавитах — за время, соответственно, $O(n \log n)$ и $\Theta(n)$.

Далее для $\ell = 2$ мы хотим найти все повторяющиеся подстроки длиной 2. Поскольку повторяющиеся подстроки длиной $\ell \geq 2$ будут иметь общий префикс длиной $\ell - 1$, то, очевидно, вычисления уровня ℓ должны привести к последовательностям, которые будут *подпоследовательностями* последовательностей, вычисленных на уровне $\ell - 1$. Другими словами, на уровне $\ell \geq 2$ разбиение множества позиций строки на отдельные непересекающиеся множества (последовательности) — это **декомпозиция** разбиения на уровне $\ell - 1$. Пример такой декомпозиции показан в табл. 12.1.

Таблица 12.1. Последовательная декомпозиция строки $f_6 = abaababaabaab$

$\ell = 2$	$\langle 1, 4, 6, 9, 12 \rangle$	$\langle 3, 8, 11 \rangle$	$\langle 2, 5, 7, 10 \rangle$	$\langle 13 \rangle$	
	ab	aa	ba	$b\$$	
$\ell = 3$	$\langle 1, 4, 6, 9 \rangle$	$\langle 12 \rangle$	$\langle 3, 8, 11 \rangle$	$\langle 2, 7, 10 \rangle$	$\langle 5 \rangle$
	aba	$aa\$$		baa	bab
$\ell = 4$	$\langle 1, 6, 9 \rangle$	$\langle 4 \rangle$	$\langle 3, 8 \rangle$	$\langle 11 \rangle$	$\langle 2, 7, 10 \rangle$
	$abaa$	$abab$	$aaba$	$aab\$$	$baab$
$\ell = 5$	$\langle 1, 6, 9 \rangle$	$\langle 3 \rangle$	$\langle 8 \rangle$	$\langle 2, 7 \rangle$	$\langle 10 \rangle$
	$abaab$	$aabab$	$aabaa$	$baaba$	$baab\$$
$\ell = 6$	$\langle 1, 6 \rangle$	$\langle 9 \rangle$	$\langle 2 \rangle$	$\langle 7 \rangle$	
	$abaaba$	$abaab\$$	$baabab$	$baabaa$	
$\ell = 7$	$\langle 1 \rangle$	$\langle 6 \rangle$			
	$abaabab$	$abaabaa$			

Примечание. В этой таблице одноэлементные множества после их первого появления были опущены, поскольку их дальнейшая декомпозиция невозможна.

Нетрудно реализовать процесс декомпозиции так, чтобы в самом худшем случае для его реализации потребовалось время порядка $\Theta(n^2)$. Это процесс *непосредственной* декомпозиции каждой последовательности, полученной на уровне ℓ , на подпоследовательности для уровня $\ell + 1$. Например, строка $x = (ab)^{n/2}$ для каждого уровня ℓ потребует $n - \ell + 1$ сравнений, следовательно, всего необходимо $\Theta(n^2)$ сравнений. Однако в [117] описаны две идеи, которые уменьшают временную сложность этого процесса.

1. Декомпозицию каждой последовательности можно получить косвенным путем, а не путем прямых вычислений. Идея такого подхода состоит в следующем.

На каждом уровне ℓ выполняется непосредственная декомпозиция любой последовательности $c_j^{(\ell)}$, основанной на конкретных символах, находящихся в первой позиции справа от подстроки, соответствующей $c_j^{(\ell)}$. Более точно, если $c_j^{(\ell)} = \langle p_1, p_2, \dots, p_r \rangle$, то необходимо проверить совпадение букв $x[p_1 + \ell], x[p_2 + \ell], \dots, x[p_r + \ell]$, и если какие-либо пары букв равны, например $x[p_{h_1} + \ell] = x[p_{h_2} + \ell]$, тогда p_{h_1} и p_{h_2} помещаются в одну и ту же последовательность на уровне $\ell + 1$.

Таким образом, например, в приведенном выше примере строки f_6 последовательность $\langle 1, 4, 6, 9 \rangle$ на уровне $\ell = 3$ разбивается на уровне $\ell = 4$ на последовательности $\langle 1, 6, 9 \rangle$ и $\langle 4 \rangle$, поскольку

$$f_6[1 + 3] = f_6[6 + 3] = f_6[9 + 3] \neq f_6[4 + 3].$$

Но декомпозицию можно также выполнить, основываясь не на разбиваемой последовательности, а на последовательностях, *относительно которых* будут разбиваться другие последовательности. Снова рассмотрим уровень $\ell = 3$ и последовательность

$$c_1^{(3)} = \langle p_1, p_2, p_3, p_4 \rangle = \langle 1, 4, 6, 9 \rangle,$$

относящуюся к подстроке aba . Для каждой позиции $p_h > 1$ известно, что подстрока $f_6[p_h - 1..p_h + \ell - 1]$ (длиной $\ell + 1$) принадлежит к некоторой последовательности $c_{j'}^{(\ell+1)}$ на уровне $\ell + 1$. Поскольку последовательность $c_j^{(\ell)}$ соответствует уникальной подстроке строки f_6 , каждая такая последовательность $c_{j'}^{(\ell+1)}$ должна формироваться из тех же позиций $p_{h_1}, p_{h_2}, \dots, p_{h_k}$ последовательности $c_j^{(\ell)}$, которые определяют класс эквивалентности

$$f_6[p_{h_1} - 1] = f_6[p_{h_2} - 1] = \dots = f_6[p_{h_k} - 1].$$

В нашем примере $c_j^{(\ell)} = c_1^{(3)} = \langle 1, 4, 6, 9 \rangle$, и, поскольку $f_6[4 - 1] = f_6[9 - 1] = a$, последовательность $\langle 3, 8, 11 \rangle$ уровня 3 на уровне 4 разбивается на подпоследовательности $\langle 3, 8 \rangle$ и $\langle 11 \rangle$. Далее, так как $f_6[5 - 1] = b$, последовательность $\langle 5 \rangle$ уровня 3 на уровне 4 переходит просто в саму себя. Таким образом, декомпозиция последовательностей $\langle 3, 8, 11 \rangle$ и $\langle 5 \rangle$ выполняется *по отношению* к последовательности $\langle 1, 4, 6, 9 \rangle$.

Подобным образом декомпозиция относительно последовательности $\langle 2, 7, 10 \rangle$ на уровне 3 приводит к равенствам

$$f_6[2 - 1] = f_6[7 - 1] = f_6[10 - 1] = a,$$

и, следовательно, на уровне 4 формируется последовательность $\langle 1, 6, 9 \rangle$, которая является декомпозицией последовательности $\langle 1, 4, 6, 9 \rangle$ уровня 3. Таким способом декомпозицию на уровне $\ell + 1$ можно выполнить косвенным путем, рассматривая каждую последовательность уровня ℓ с позиции, находящейся на одну позицию левее от начальной позиции этой последовательности.

2. Косвенная декомпозиция не обязательно должна выполняться относительно *каждой* последовательности уровня ℓ : в каждом наборе последовательностей, порожденных одной последовательностью уровня $\ell - 1$, всегда можно исключить использование одной из них для декомпозиции последовательностей на уровне ℓ .

Снова рассмотрим декомпозицию последовательности $\langle 1, 4, 6, 9 \rangle$ уровня 3 на подпоследовательности $\langle 1, 6, 9 \rangle$ и $\langle 4 \rangle$ уровня 4. Последовательность

$\langle 1, 4, 6, 9 \rangle$ соответствует строке с собственным суффиксом ba , представленной на уровне 2 последовательностью $\langle 2, 5, 7, 10 \rangle$: таким образом, любая декомпозиция последовательности $\langle 1, 4, 6, 9 \rangle$ на уровне 4 получается из соответствующей декомпозиции последовательности $\langle 2, 5, 7, 10 \rangle$ уровня 3. Фактически последовательность $\langle 1, 6, 9 \rangle$ получена из подпоследовательности $\langle 2, 7, 10 \rangle$ (суффикса baa), в то время как $\langle 4 \rangle$ получена из подпоследовательности $\langle 5 \rangle$ (суффикс bab). Этот пример показывает, что для декомпозиции последовательности $\langle 1, 4, 6, 9 \rangle$ на уровне 4 достаточно использовать *все, за исключением одной*, декомпозиции последовательности $\langle 2, 5, 7, 10 \rangle$ уровня 3 — конечным разбиением будет то, что останется от последовательности $\langle 1, 4, 6, 9 \rangle$ после того, как будут вычислены другие разбиения. В этом примере есть два выбора: можно получить или $\langle 1, 6, 9 \rangle$ из $\langle 2, 7, 10 \rangle$, оставляя $\langle 4 \rangle$ как остаток, или $\langle 4 \rangle$ из $\langle 5 \rangle$, оставляя $\langle 1, 6, 9 \rangle$ в виде остатка. Конечно, второй подход предпочтительнее, поскольку требует только одного вычисления.

Эти примеры подводят к следующему определению.

Определение 12.1.1. В декомпозиции последовательности $c_j^{(\ell)}$ на подпоследовательности $(c_1^{(\ell+1)}, c_2^{(\ell+1)}, \dots, c_q^{(\ell+1)})$ ($q \geq 1$) назовем одну подпоследовательность с самым большим количеством элементов **большой**, а остальные $q - 1$ подпоследовательностей — **малыми**. Для уровня $\ell = 1$ каждая последовательность **малая**. ■

Отметим, что в случае $q = 1$ не будет ни одной малой подпоследовательности, в частности, всегда, когда $c_j^{(\ell)}$ состоит из одного элемента.

Теперь основную идею алгоритма Крочемора кратко можно изложить в следующем виде: выполнять декомпозицию последовательностей на каждом уровне ℓ *только* относительно последовательностей, которые на этом уровне являются малыми. Как показывает следующий результат, эта стратегия также определяет эффективность процесса декомпозиции.

Лемма 12.1.2. Предположим, что декомпозиция последовательностей, соответствующих произвольной строке $x = x[1..n]$, выполняется для уровней $\ell = 1, 2, \dots, \ell^*$, где ℓ^* ($\ell^* \in 1..n - 1$) — наименьший уровень, на котором каждая последовательность содержит единственную позицию. Тогда каждая позиция i строки x входит в малые последовательности $O(\log n)$ раз.

Доказательство. Заметим, что если последовательность $c_j^{(\ell)}$ разбивается на подпоследовательности $(c_1^{(\ell+1)}, c_2^{(\ell+1)}, \dots, c_q^{(\ell+1)})$, то каждая малая последовательность $c_{j'}^{(\ell+1)}$ должна удовлетворять условию $|c_{j'}^{(\ell+1)}| \leq |c_j^{(\ell)}|/2$. Другими словами,

каждая малая подпоследовательность при $\ell \geq 2$ не превышает половины размера своей исходной последовательности. Поскольку для $\ell - 1$ начальная малая последовательность может содержать не более n позиций, то из этого следует, что ни одна из позиций не может входить в больше чем $\lceil \log_2(n + 1) \rceil$ малых последовательностей. ■

Поскольку строка x содержит n позиций, то из леммы 12.1.2 следует, что всего в малых последовательностях на всех уровнях содержится $O(n \log n)$ позиций. Таким образом, если время обработки последовательностей на каждом уровне ℓ пропорционально количеству элементов в малых последовательностях этого уровня, то полный процесс декомпозиции будет выполнен за время порядка $O(n \log n)$. Представленный ниже алгоритм 12.1.1 показывает схему этого процесса, который выглядит в достаточной степени простым. Это и стало причиной известности алгоритма Крочемора, однако основная сложность этого алгоритма (и других, ему подобных) заключается в структурах данных, необходимых для получения времени выполнения алгоритма порядка $O(n \log n)$.

Алгоритм 12.1.1 (Алгоритм Крочемора)

▷ *Вычисление всех кратных строк в строке $x = x[1..n]$*
 $\ell \leftarrow 1$
 вычисление последовательностей на уровне 1,
 все они помечаются как *малые*
while на уровне ℓ есть малая последовательность **do**
 output кратные строки с периодом ℓ (если они есть)
 декомпозиция последовательности уровня ℓ на подпоследовательности
 уровня $\ell + 1$, используя позиции только
 в малых последовательностях
 $\ell \leftarrow \ell + 1$
 вычисление малых последовательностей на уровне ℓ

Прояснив основную идею алгоритма Крочемора, классифицируем используемые им структуры данных в соответствии с их основными функциями.

1. Запись текущей последовательности для каждой позиции в строке x .

Массив $SEQ[\ell..n]$: $SEQ[i]$ содержит индекс текущей последовательности, которой принадлежит i -я позиция в строке x .

Массив $SEQLIST[1..2n]$: $SEQLIST[j]$ — указатели на двусвязный список позиций, принадлежащих последовательности с индексом j и расположенных в порядке их возрастания.

Массив $SEQSIZE[1..2n]$: $SEQSIZE[j]$ равно количеству позиций в последовательности с индексом j , т.е. количеству позиций в списке, на который указывает элемент $SEQLIST[j]$.

Стек $INDEXSTACK$: стек неиспользованных индексов последовательностей, инициализированный для размещения индексов $1, 2, \dots, 2n$.

Всякий раз, когда необходимо сформировать новую последовательность или на уровне 1, или в качестве части декомпозиции существующей последовательности, из стека $INDEXSTACK$ извлекается индекс последовательности. Если последовательность становится пустой в результате ее декомпозиции на подпоследовательности, индекс последовательности возвращается в стек $INDEXSTACK$. Всякий раз, когда к последовательности с индексом j добавляется позиция i , выполняются операторы присваивания

$$SEQ[i] \leftarrow j; SEQSIZE[j] \leftarrow SEQSIZE[j] + 1,$$

а номер i добавляется в конец списка $SEQLIST[j]$. Если i -я позиция удаляется из последовательности с индексом j в результате его вхождения в подпоследовательность с индексом j , выполняется присваивание $SEQSIZE[i] \leftarrow SEQSIZE[j] - 1$, а i -я позиция удаляется из списка $SEQLIST[i]$.

Как показано в упражнении 12.1.4, в любой момент используется не более $2n$ индексов последовательности. Следовательно, массива размером $[1..2n]$ будет достаточно для выполнения всех операций. Заметим также, что поскольку i -е позиции всегда сохраняются в списке $SEQLIST$ в порядке их возрастания и поэтому всегда удаляются из этого списка в таком же порядке, то операции и сохранения и удаления можно выполнить за константное время. Отсюда следует, что любую операцию, выполняемую в процессе декомпозиции для записи изменений принадлежности позиций к разным последовательностям, можно реализовать за константное время.

2. Управление малыми последовательностями.

Очередь $SMALL$: очередь индексов j малых последовательностей, организованная таким образом, что (для $\ell > 1$) все имеющиеся в ней малые последовательности являются подпоследовательностями одной и той же последовательности предыдущего уровня.

- Очередь QUEUE: очередь пар (i, j) , где i — позиция в малой последовательности с индексом j , которая должна использоваться для декомпозиции текущего уровня; для каждого j позиции i располагаются в возрастающем порядке.
- Массив LASTSMALL[1..2n]: LASTSMALL[j'] = $j > 0$ тогда и только тогда, когда j — это индекс малой последовательности, использованной последней по времени для декомпозиции последовательности с индексом j' .

Для $\ell = 1$ очередь SMALL содержит все индексы, которые первоначально вычисляются в соответствии с каждым символом строки x ; для $\ell > 1$ очередь SMALL создается повторно путем просмотра подпоследовательностей тех последовательностей, для которых была выполнена декомпозиция на уровне ℓ , и добавлением к очереди SMALL всех подпоследовательностей, кроме единственной “большой” подпоследовательности (см. описание массивов SPLIT и SUBSEQ в следующем пункте). Для каждого элемента j в очереди SPLIT эти действия можно выполнить за один просмотр списка, на который указывает элемент SUBSEQ[j]. Таким образом, очередь SMALL модифицируется за время, пропорциональное количеству имеющихся в ней элементов.

Очередь QUEUE создается в начале обработки для каждого уровня $\ell > 1$ путем удаления элементов из очереди SMALL. Для каждой последовательности j , обнаруженной в очереди SMALL, извлекаются найденные в списке SEQLIST[j] позиции $i > 1$, а элементы (i, j) добавляются к очереди QUEUE. Основная обработка каждого из уровней $\ell > 1$ выполняется после того, как из очереди QUEUE один за другим удаляются элементы (i, j) . Для каждой удаленной позиции i последовательность $j' = \text{SEQ}[i-1]$ является той последовательностью, которая будет подвергнута декомпозиции. Следовательно, если LASTSMALL[j'] $\neq j$, то последовательность j' была последней по времени подвергнута декомпозиции относительно некоторого другого класса, поэтому необходим новый индекс последовательности. Соответственно с этим устанавливаем LASTSMALL[j'] $\leftarrow j$, извлекаем из стека INDEXSTACK следующий номер последовательности и добавляем его к списку, на который указывает SUBSEQ[j'] (см. следующий пункт).

Таким образом, для выполнения операций, связанных с каждым элементом очереди QUEUE, необходимо одинаковое константное время.

3. Организация подпоследовательностей.

Массив $SPLITFLAG[1..2n]$: $SPLITFLAG[j] = 1$ тогда и только тогда, когда последовательность с индексом j будет подвергнута декомпозиции на текущем уровне (с использованием последовательностей, которые хранятся в очереди $SMALL$).

Список $SPLIT$: список последовательностей j , которые будут подвергнуты декомпозиции на текущем уровне (т.е. для которых $SPLITFLAG[j] = 1$).

Массив $SUBSEQ[1..2n]$: при декомпозиции последовательности j на текущем уровне $SUBSEQ[j]$ указывает на список индексов подпоследовательностей последовательности j , первым элементом в этом списке будет сам индекс j .

Все эти структуры данных инициализируются для каждого уровня $\ell > 1$ одновременно с формированием из очереди $SMALL$ очереди $QUEUE$. В частности, когда (i, j) добавляется к очереди $QUEUE$, известно, что i будет использоваться для декомпозиции последовательности j' , в которой имеется позиция $i - 1$ ($j' = SEQ[i - 1]$). Поэтому, если $SPLITFLAG[j'] = 0$, то j' в это время хранится в $SPLIT$ — списке последовательностей, которые подлежат декомпозиции на уровне ℓ , в то время как выполняются следующие присваивания:

$$SPLITFLAG[j'] \leftarrow 1; SUBSEQ[j'] \leftarrow \langle j' \rangle; LASTSMALL[j'] \leftarrow 0.$$

Из этого следует, что для помещения в очередь $QUEUE$ любого элемента требуется одинаковое константное время, и, как было показано во втором пункте, на удаление любого элемента из очереди $QUEUE$ также необходимо одинаковое константное время. Поскольку элементы в очереди $QUEUE$ являются позициями в малых последовательностях, то, очевидно, подпоследовательности и малые последовательности можно организовать заранее пропорционально количеству малых последовательностей, на что, в соответствии с леммой 12.1.2, необходимо время порядка $O(n \log n)$.

По завершении обработки каждого уровня ℓ ($\ell > 1$) последовательности j удаляются один за другим из списка $SPLIT$: если $SEQSIZE[j] = 0$, то эта последовательность уже пуста и поэтому ее индекс больше не требуется. Поэтому j помещается в стек $INDEXSTACK$ и удаляется из списка, на который указывает $SUBSEQ[j]$.

4. Вычисление кратных строк.

Массив $\text{GAP}[\ell..n]$: элемент $\text{GAP}[i]$ равен положительной разности между i -й позицией в строке x и следующей большей позицией в той же самой последовательности на текущем уровне; если нет большей позиции, то $\text{GAP}[i] = \infty$.

Массив $\text{GAPLIST}[\ell..n]$: $\text{GAPLIST}[g]$ указывает на двусвязный список i -х позиций, для которых $\text{GAP}[i] = g$.

Для $\ell = 1$ массивы GAP и GAPLIST инициализируются непосредственно путем обработки каждого списка, на который указывает $\text{SEQLIST}[j]$, где $j = 1, 2, \dots, \alpha'$, вычисляя $\text{GAP}[i]$ для каждой i -й позиции в списке $\text{SEQLIST}[j]$, затем добавляя i -ю позицию к списку, на который указывает $\text{GAPLIST}[\text{GAP}[i]]$.

При $\ell > 1$ массивы GAP и GAPLIST обновляются при удалении элемента (i, j) из очереди QUEUE , что приводит к декомпозиции последовательности $j' = \text{SEQ}[i - 1]$, при которой позиция $i - 1$ в $\text{SEQLIST}[j']$ будет перемещена в конец $\text{SEQLIST}[j']$. Поскольку списки, на которые указывают и SEQLIST и GAPLIST , двусвязные, повторное вычисление элементов массива GAP , необходимое из-за удаления позиции $i - 1$ из одного списка и ее вставки в другой список, можно выполнить непосредственно.

Мы видим, что, после создания последовательностей для уровня ℓ , каждый элемент в $\text{GAPLIST}[\ell]$ описывает две идентичные подстроки строки x длиной ℓ , первые буквы которых находятся друг от друга на расстоянии ℓ , другими словами, имеем квадрат этой подстроки. В этом случае массивы GAP и GAPLIST можно использовать так, как показано в следующем фрагменте листинга, для вывода нормальной формы максимальных кратных строк для каждого уровня ℓ за время, пропорциональное количеству квадратов.

```

while  $\text{GAPLIST}[\ell]$  not пусто do
     $i_0 \leftarrow$  следующая позиция в  $\text{GAPLIST}[\ell]$ 
     $i \leftarrow i_0; r \leftarrow 1$ 
    repeat
        удаление  $i$  из  $\text{GAPLIST}[\ell]$ 
         $i \leftarrow i + \ell; r \leftarrow r + 1$ 
    until  $\text{GAP}[i] \neq \ell$ 
    output  $(i_0, \ell, r)$ 

```

Предположим, что в строке x встречается подстрока u^2 , где $|u| = \ell$, а u сама является кратной строкой v^k для некоторого целого числа $k \geq 2$. Тогда $u^2 = v^{2k}$ содержит $k + 1$ вхождений u с пропуском $\ell/k < \ell$, поэтому подстрока u никогда не будет выбрана в качестве образующей кратных строк. Поскольку, как было отмечено ранее, в строке x имеется не более $O(n \log n)$ квадратов, образующие которых сами не являются кратными строками [157], тогда все кратные строки можно найти за время $O(n \log n)$.

Этим завершается наше описание алгоритма Крочемора. Из-за необходимости реализации описанных пп. 1–4, этот алгоритм очень объемён и технически сложен. Мы не представляем здесь его во всех деталях; заинтересованного читателя отправляем к работе [69]. Тем не менее мы попытались показать читателю, каким образом сложные структуры данных могут использоваться для декомпозиции каждого уровня и обеспечивать время вычисления кратных строк порядка $O(n \log n)$. Вспомнив, что декомпозицию уровня 1 на упорядоченном алфавите можно вычислить за время $O(n \log n)$, получаем следующий важный результат.

Теорема 12.1.3. Алгоритм 12.1.3 вычисляет все кратные строки в произвольной строке $x = x[1..n]$, определенной на упорядоченном алфавите, за время порядка $O(n \log n)$ с использованием памяти объёмом $\Theta(n)$.

Доказательство. Истинность приведенной оценки верхней границы времени выполнения алгоритма следует из истинности трех отдельных утверждений:

- алфавит является упорядоченным;
- лемма 12.1.2;
- количество обработанных квадратов в любой строке длиной n имеет порядок $O(n \log n)$ [157]. ■

Важно отметить, что для вычисления кратных строк с помощью алгоритма Крочемора необходимо просмотреть не более $\lfloor n/2 \rfloor$ уровней, поскольку в любой строке x максимальный период ℓ ее кратных строк удовлетворяет условию $\ell \leq \lfloor n/2 \rfloor$. Но оценка верхней временной границы $O(n \log n)$ будет справедлива даже тогда, когда декомпозиция уровней продолжается до тех пор, пока в SEQLIST не останутся только одноэлементные списки. В этом случае придется рассмотреть $\ell > \lfloor n/2 \rfloor$ уровней, например, как для строки $x = abaababaababa$, содержащей длинные перекрывающиеся подстроки. Из этого видно, что алгоритм Крочемора можно легко расширить для вычисления всех раппортов строки x за время порядка $O(n \log n + Q)$, где Q — размер выводимых данных. Мы возвратимся к этой задаче в разделе 13.2, где покажем, как все повторяемые подстроки в строке x можно представить в виде данных объёмом $Q \in O(n)$.

Нетрудно заметить, что алгоритм Крочемора в расширенном виде по существу является алгоритмом для построения дерева суффиксов. Это легко увидеть, если

представить декомпозицию строки f_6 из предыдущего примера в виде дерева (рис. 12.1), где каждое ребро помечено буквой, на которой основывается текущая декомпозиция, а конечный узел создается тогда, когда эта буква оказывается единственным элементом декомпозиции.

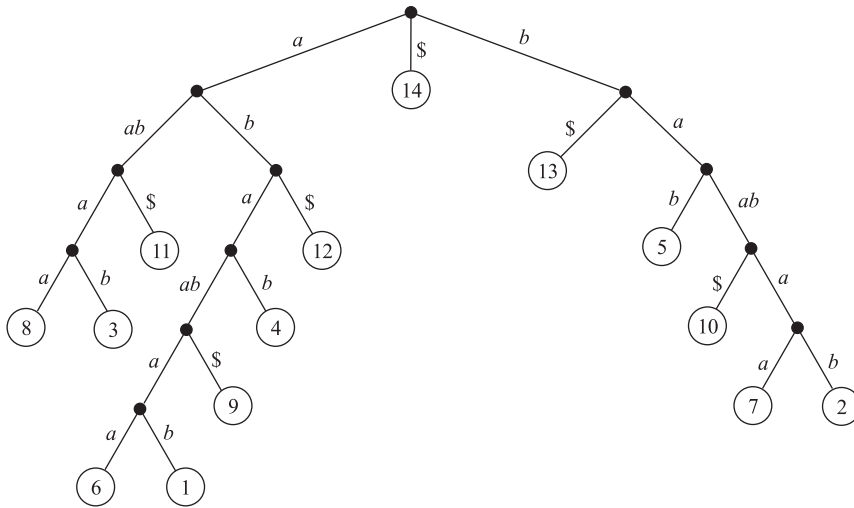


Рис. 12.1. Дерево декомпозиции Крочемора для строки $f_6 = abaababaabaab\$$

Сравнивая дерево на рис. 12.1 с деревом суффиксов, изображенным на рис. 5.2, мы видим, что для преобразования одного в другое достаточно простого переименования меток некоторых ребер. Кроме того, необходимо удостовериться, что в дереве Крочемора все ребра, которые ведут к конечным узлам, представляют суффиксы. Таким образом, расширенная версия алгоритма Крочемора позволяет вычислять все общие префиксы суффиксов $x[l..n]$ — эта информация содержится в дереве суффиксов T_x .

К сожалению, алгоритм Крочемора также обладает и другим свойством, присущим всем алгоритмам построения деревьев суффиксов: он требует для своего выполнения большого объема памяти. Действительно, в силу широкого использования массивов, очередей и двусвязных списков, известные реализации алгоритма Крочемора потребовали не менее $80n$ байт памяти, хотя это количество было недавно урезано наполовину с помощью нового подхода [94]. Напомним, что в подразделе 5.2.5 было сделано предположение о том, что любая реализация дерева суффиксов потребует не менее $10n$ байт памяти, — это нижняя граница объема памяти, необходимой любому алгоритму построения дерева суффиксов.

12.1.2 Алгоритм Мейна и Лоренца

Подход, предпринятый в алгоритме Мейна–Лоренца (Main–Lorentz, сокращенно — алгоритм МЛ), полностью отличается от подхода, используемого в алгоритме Крочемора. Интересно, что алгоритм МЛ работает с неупорядоченным алфавитом, но все равно выполняется за время порядка $\Theta(n \log n)$. Этот алгоритм не требует лексикографического упорядочения букв и подстрок строки x , тогда как эти условия являются критическими на первом шаге алгоритма Крочемора. Алгоритм МЛ — это алгоритм типа “разделяй и властвуй” [5], который вычисляет все квадраты подстрок, следовательно, все кратные строки в строковой последовательности $x = x[\ell..n]$, распределяя их на три различных класса.

- Класс C1 — класс квадратов w^2 , которые начинаются и заканчиваются в подстроке $x[1..[n/2]]$, т.е. $w^2 = x[i..i + 2|w| - 1]$, где $i + 2|w| - 1 \leq [n/2]$.
- Класс C2 — класс квадратов w^2 , которые начинаются и заканчиваются в подстроке $x[[n/2] + 1..n]$, т.е. $w^2 = x[i..i + 2|w| - 1]$, где $i > [n/2]$.
- Класс C3 — класс квадратов w^2 , которые начинаются в подстроке $x[1..[n/2]]$ и заканчиваются в подстроке $x[[n/2] + 1..n]$, т.е. здесь $w^2 = x[i..i + 2|w| - 1]$, где $1 \leq i \leq [n/2]$ и $[n/2] < i + 2|w| - 1 \leq n$.

Обозначим $u = [1..[n/2]]$ и $v = x[[n/2] + 1..n]$. Основная идея алгоритма МЛ заключается в вычислении квадратов подстрок в строке $x = uv$ на основе эффективного вычисления новых квадратов класса C3. Действительно, рекурсивное разбиение строки x на половинные подстроки сводит вычисление *всех* квадратов до квадратов класса C3. Эта идея приводит к следующему прямому рекурсивному алгоритму.

Алгоритм 12.1.2 (Алгоритм Мейна–Лоренца)

▷ Вычисление всех квадратов подстрок в строке $x = x[1..n]$

procedure ML(x)

if $|x| > 1$ **then**

▷ Вычисление квадратов класса C3

C3($x[1..[n/2]]$, $x[[n/2] + 1..n]$)

▷ Применение рекурсии

ML($x[1..[n/2]]$)

ML($x[[n/2] + 1..n]$)

Далее мы увидим, что для произвольных строк u и v процедуру C3(u, v) можно реализовать за время $\Theta(|uv|)$. Если обозначить через t_n максимальное время, необходимое для выполнения процедуры ML на любой строке длиной n , то очевидно, что для $n > 1$ существует константа k_2 , которая удовлетворяет условию

$$t_n \leq k_2 n + 2t_{[n/2]}, \quad (12.2)$$

тогда как для $0 \leq n \leq 1$ существует константа k_1 , такая, что

$$t_1 \leq k_1. \quad (12.3)$$

Как показано в упражнении 12.1.5, разрешая неравенства (12.2) и (12.3) относительно t_n , получим, что $t_n \in O(n \log n)$. Но поскольку время, необходимое для выполнения процедуры СЗ, в точности пропорционально величине $|uv|$, то, следовательно, существуют константы k'_1 и k'_2 , которые удовлетворяют условиям

$$t'_1 \geq k'_1, t'_n \geq k'_2 n + 2t'_{\lceil n/2 \rceil}, n > 1, \quad (12.4)$$

где t'_n — минимальное время, необходимое для выполнения алгоритма МЛ на любой строке длиной n . Здесь опять, разрешая рекуррентное соотношение (12.4), можно показать, что $t'_n \in \Omega(n \log n)$. Таким образом, имеем следующую теорему.

Теорема 12.1.4. Если процедура СЗ выполняется за линейное время, то алгоритм 12.1.2 вычисляет все квадраты подстрок в произвольной строке $x = x[1..n]$ за время $\Theta(n \log n)$. ■

Для описания процедуры СЗ(u, v) вначале отметим, что новые квадраты w^2 в строке uv могут быть двух видов:

- **правый квадрат**, когда второе вхождение w в строку uv полностью находится в подстроке v ;
- **левый квадрат**, когда в подстроке u имеется непустой префикс второго вхождения w .

Например, подстроки $u = aba$ и $v = ba$ порождают правый квадрат $w^2 = (ba)^2$ и левый квадрат $w^2 = (ab)^2$. Здесь мы опишем процесс вычисления правых квадратов и оставим для упражнения 12.1.10 симметричный процесс вычисления левых квадратов.

Процедура СЗ реализуется как комбинация процессов вычисления правых и левых квадратов. Пусть $u = u[\ell..n_1]$ и $v = v[\ell..n_2]$. Чтобы в uv найти все правые квадраты, необходимо вычислить два массива, которые являются некоторым подобием массива граней (раздел 1.3). Для определения этих массивов напомним (подраздел 5.2.1), что $\text{lcp}(x_1, x_2)$ обозначает длину наибольшего общего префикса строк x_1 и x_2 . Используем подобное обозначение $\text{lcs}(x_1, x_2)$ для длины наибольшего общего суффикса строк x_1 и x_2 . (Не путать lcs с LCS — наибольшей общей подпоследовательностью из главы 9!) Новые массивы можно кратко определить следующим образом.

- Массив $\text{LP}[2..n_2 + 1]$: для каждой позиции $i = 2, 3, \dots, n_2$ строки v $\text{LP}[i] = \text{lcp}(v[i..n_2], v)$, при этом для $i = n_2 + 1$ принимаем, что $\text{LP}[i] = 0$.
- Массив $\text{LS}[1..n_2]$: для каждой позиции $i = 1, 2, \dots, n_2$ строки v $\text{LS}[i] = \text{lcs}(v[1..i], u)$.

Например, для строки $v = abaababa$ массив $LP[2..9] = 01303010$. Если $u = abaab$, тогда $LS[1..8] = 02005020$. Приблизительно можно сказать, что массив LP рекуррентно определяет максимальные длины префиксов v внутри самой v , а массив LS определяет вхождения суффиксов u максимальной длины в строку v .

Эти массивы, которые похожи на массивы граней, можно вычислить за время $\Theta(n_2)$. Позже мы опишем процесс их вычисления, но сначала покажем, каким образом эти массивы можно использовать для вычисления правых квадратов строк u и v . Для этого необходима следующая лемма.

Лемма 12.1.5. Пусть $u[1..n_1]$ и $v[1..n_2]$ — произвольные непустые строки и пусть $p \leq n_2$ и $i \in p..2p-1$ — положительные целые числа. В uv имеется квадрат длиной $2p$, заканчивающийся на i -й позиции строки v , тогда и только тогда, когда

$$2p - LS[p] \leq i \leq p + LP[p + 1].$$

Доказательство. В строке uv для любого квадрата w^2 длиной $2p$, заканчивающегося на i -й позиции, первое вхождение w должно начинаться с буквы $u[n_1 - 2p + i + 1]$, а второе вхождение — с буквы $v[i - p + 1]$. Поэтому такой квадрат имеет “право на существование” тогда и только тогда, когда

$$u[n_1 - 2p + i + 1..n_1]v[1..i - p] = v[i - p + 1..i];$$

т.е. только тогда, когда

- $u[n_1 - 2p + i + 1..n_1] = v[i - p + 1..p]$ и
- $v[1..i - p] = v[p + 1..i]$.

Отметим, что если $i = p$, вследствие чего квадрат формируется из суффикса w строки u и префикса w строки v , второе условие сводится к утверждению, что $\varepsilon = \varepsilon$, и поэтому становится ненужным.

Первое условие, которое говорит о том, что строка $v[i - p + 1..p]$ длиной $2p - i$ является суффиксом строки u , выполняется только тогда, когда $LS[p] \geq 2p - i$, т.е. только в случае выполнения неравенства $2p - LS[p] \leq i$.

Аналогично, как показано в упражнении 12.1.7, второе условие эквивалентно неравенству $i \leq p + LP[p + 1]$, что и завершает доказательство. ■

Правый квадрат возможен только для тех значений p и i , которые удовлетворяют условиям леммы 12.1.5. Кроме того, для каждого допустимого значения $p \in 1..n_2$ соответствующие допустимые значения величины i могут находиться только в интервале $2p - LS[p]..p + LP[p + 1]$. Поэтому при условии, что массивы LP и LS вычислены заранее, фактические значения i определяются за константное время путем вывода границ этих интервалов всякий раз, когда $LP[p + 1] + LS[p] \geq p$. Следовательно, суммируя по всем допустимым значениям p , правые квадраты можно вычислить за время $\Theta(n_2)$.

Оставим для упражнения 12.1.8 запись простого алгоритма, который использует массивы LP и LS для выполнения этих вычислений.

Для вычисления левых квадратов нужны аналогичные массивы.

- Массив $LP'[2..n_1]$: для каждой позиции $i = 2, 3, \dots, n_1$ строки \mathbf{u} $LP'[i] = \text{lcp}(\mathbf{u}[i..n_1], \mathbf{v})$.
- Массив $LS'[1..n_1 - 1]$: для каждой позиции $i = 1, 2, \dots, n_1 - 1$ строки \mathbf{u} $LS'[i] = \text{lcs}(\mathbf{u}[1..i], \mathbf{u})$.

Как показано в упражнении 12.1.9, для массивов LP' и LS' имеет место аналог леммы 12.1.5. Поэтому время вычисления всех левых квадратов в строке \mathbf{uv} имеет порядок $\Theta(n_1)$.

Как было отмечено выше, процедура СЗ вычисляет правые и левые квадраты класса СЗ за время $\Theta(n_1 + n_2) = \Theta(|\mathbf{uv}|)$. Таким образом, если массивы $[LP[2..n_2 + 1]$ и $LS[1..n_2]$ можно вычислить за время $\Theta(n_2)$, а аналогичные массивы $LP'[2..n_1]$ и $LS'[1..n_1 - 1]$ — за время $\Theta(n_1)$, тогда, в соответствии с теоремой 12.0.1, для выполнения алгоритма МЛ необходимо время порядка $\Theta(n \log n)$.

Далее покажем, каким образом можно эффективно вычислить необходимые для алгоритма массивы.

Вычисление массивов LP, LS, LP', LS'

Вспомним, что для произвольной строки $\mathbf{v}[1..n_2]$ и любого целого числа $i \in 1..n_2 + 1$ элемент $LP[i]$ равен длине наибольшего префикса подстроки $\mathbf{v}[i..n_2]$, который совпадает с префиксом самой строки \mathbf{v} . Будем вести вычисление массива LP слева направо для возрастающих значений i , считая, что при вычислении элемента $LP[i]$ все элементы $LP[k]$, где $2 \leq k < i$, уже вычислены. Предположим, что для некоторого значения $i \geq 3$ существует такое целое число $k \in 2..i - 1$, что $k + LP[k] > i$. Как показано на рис. 12.2, это означает, что существует подстрока $\mathbf{w} = \mathbf{v}[i..k + LP[k] - 1]$ длиной $\ell = k + LP[k] - i$, совпадающая с суффиксом $\mathbf{v}[i - k + 1..LP[k]]$ префикса $\mathbf{v}[1..LP[k]]$.

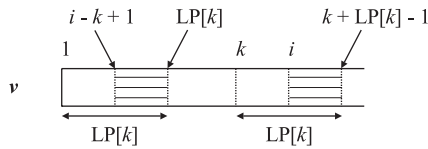


Рис. 12.2. Вычисление $LP[i]$ при $k + LP[k] > i$

Теперь рассмотрим элемент $LP[i - k + 1]$, равный длине наибольшей подстроки, совпадающей с префиксом строки \mathbf{v} . Если $LP[i - k + 1] < |\mathbf{w}| = \ell$, тогда должно выполняться равенство

$$LP[i] = LP[i - k + \ell]. \tag{12.5}$$

С другой стороны, если имеет место неравенство $LP[i - k + 1] \geq \ell$, тогда

$$LP[i] \geq \ell, \quad (12.6)$$

и поэтому w — префикс наибольшей (при данном i) подстроки, которая совпадает с префиксом строки v . Следовательно, если существует k , такое, что $k + LP[k] > i$, то элемент $LP[i]$ вычисляется точно согласно выражению (12.5) или, в соответствии с (12.6), для него установлена нижняя граница.

Конечно, если нет такого целого числа $k \in 2..i - 1$, что $k + LP[k] > i$, тогда единственной нижней границей для $LP[i]$ будет нуль, т.е. в этом случае

$$LP[i] \geq 0. \quad (12.7)$$

Три случая (12.5)–(12.7) в алгоритме вычисления массива LP необходимо рассмотреть отдельно. Как показано в алгоритме 12.1.3, после определения нижней границы для $LP[i]$ (случаи (12.6) и (12.7) обрабатываются одинаково) путем по-символьного сравнения вычисляется точное значение элемента $LP[i]$. Заметим, что алгоритм может работать и со значениями k , которые меньше i , лишь бы это k доставляло максимум выражению $k + LP[k]$. Для удобства вычисления конечных позиций в массиве LP в конец строки v добавлен символ $\$2$, который не совпадает ни с одной буквой алфавита.

Алгоритм 12.1.3 (Вычисление массива LP)

```

▷ Вычисление массива  $LP[2..n_2 + 1]$ , соответствующего строке  $v[1..n_2]\$2$ 
▷ Вначале вычисляется  $LP[2]$  (случай (12.7)) и инициализуется  $k$ 
 $LP[2] \leftarrow 0$ 
while  $v[LP[2] + \ell] = v[LP[2] + 2]$  do  $LP[2] \leftarrow LP[2] + 1$ 
 $k \leftarrow 2$ 
for  $i \leftarrow 3$  to  $n_2 + 1$  do
     $\ell \leftarrow k + LP[k] - i$     ▷ возможно отрицательное значение!
    if  $LP[i - k + 1] < \ell$  then
        ▷ Случай (12.5)
         $LP[i] \leftarrow LP[i - k + 1]$ 
    else
        ▷ Случаи (12.6) и (12.7) — в точности так же, как для  $i = 2!$ 
         $LP[i] \leftarrow \max\{0, \ell\}$ 
        while  $v[LP[i] + \ell] = v[LP[i] + i]$  do  $LP[i] \leftarrow LP[i] + 1$ 
         $k \leftarrow i$ 

```

Проведенное выше обсуждение доказывает правильность алгоритма 12.1.3. Он реализует все эти три случая (12.5)–(12.7) и делает это даже для $v[n_2 + 1] = \$2$.

Для того чтобы оценить временную сложность алгоритма, заметим, что в цикле **for** все операторы, за исключением операторов цикла **while**, выполняются за

константное время. Таким образом, время выполнения алгоритма равно $\Theta(n_2)$ плюс время выполнения циклов **while**, которое пропорционально числу приращений значений $LP[i]$. Но каждое приращение $LP[i]$ увеличивает на единицу значение $k + LP[k]$, которое в последующих итерациях будет в строке v минимальной позицией буквы, подлежащей сравнению. Поскольку начальное значение k равно 2, то, очевидно, оператор $LP[i] \leftarrow LP[i] + 1$ выполняется не более $n_2 - 1$ раз в *обоих* циклах **while**. Отсюда следует, что алгоритм вычисления массива LP выполняется за время порядка $\Theta(n_2)$.

Не трудно применить подход, реализованный в алгоритме 12.1.3, для вычисления массива $LP'[2..n_1]$, т.е. для вычисления длины наибольшего префикса подстроки $u[i..n_1]$ (а не префикса подстроки $v[i..n_2]$), который совпадает с префиксом строки v . Здесь также для каждого $i \geq 3$ определяется целое число $k \in 2..i - 1$, которое обеспечивает максимальное значение выражению $k + LP'[k]$. Как показано на рис. 12.3, это означает, что для $k + LP'[k] > i$ существует подстрока $w = u[i..k + LP'[k] - 1]$ длиной $\ell' = k + LP'[k] - i$, совпадающая с подстрокой $v[i - k + 1..LP'[k]]$.

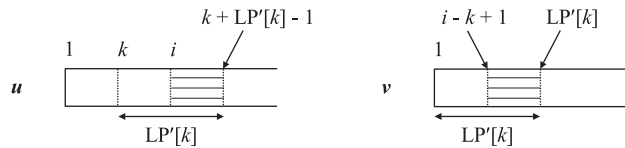


Рис. 12.3. Вычисление $LP'[i]$ при $k + LP'[k] > i$

Здесь также имеют место три случая, подобные (12.5)–(12.7). Предположим, что элемент $LP[2..n_2 + 1]$ уже вычислен и остается неизменным, если $LP[i - k + 1] < \ell'$, тогда

$$LP'[i] = LP[i - k + 1]. \tag{12.8}$$

В противном случае

$$LP'[i] \geq \ell', \tag{12.9}$$

и если $k + LP'[k] \leq i$, тогда

$$LP'[i] \geq 0. \tag{12.10}$$

Реализация трех случаев (12.8)–(12.10) приводит к алгоритму, почти идентичному алгоритму 12.1.3 и представленному в следующем листинге. С помощью таких же аргументов, что были использованы для обоснования оценки временной сложности алгоритма 12.1.3, здесь нетрудно доказать, что алгоритм вычисления массива LP' выполняется за время порядка $\Theta(n_1)$.

Алгоритм 12.1.4 (Вычисление массива LP')

\triangleright Вычисление массива $LP'[2..n_1 + 1]$, соответствующего строкам $u[1..n_1]S_1$
 \triangleright и $v[\ell..n_2]S_2$ при условии, что массив $LP[2..n_2 + 1]$ уже вычислен
 \triangleright Сначала вычисляется массив $LP'[2]$
 \triangleright и инициализируется k
 $LP'[2] \leftarrow 0$
while $v[LP'[2] + 1] = u[LP'[2] + 2]$ **do** $LP'[2] \leftarrow LP'[2] + 1$
 $k \leftarrow 2$
for $i \leftarrow 3$ **to** n_1 **do**
 $\ell' \leftarrow k + LP'[k] - i \quad \triangleright$ возможно отрицательное значение!
 \triangleright Отметим, что использование LP предпочтительнее LP'
if $LP[i - k + \ell] < \ell'$ **then**
 $LP'[i] \leftarrow LP[i - k + 1]$
else
 $LP'[i] \leftarrow \max\{0, \ell'\}$
while $v[LP'[i] + 1] = u[LP'[i] + i]$ **do** $LP'[i] \leftarrow LP'[i] + 1$
 $k \leftarrow i$

Разработку алгоритмов для вычисления массивов LS и LS' оставим для упражнения 12.1.14.

Заключение

Если для вычисления массивов LP , LP' , LS и LS' используются алгоритмы с линейным временем выполнения, тогда можно утверждать, что процедура $C3(u, v)$ выполняется за время $\Theta(|uv|)$. Следовательно, согласно теореме 12.1.4, алгоритм 12.1.2 выполняется за время $\Theta(n \log n)$. Однако вспомним, что алгоритм 12.1.1 выполняется за время $O(n \log n)$. Но с другой стороны, поскольку алгоритм 12.1.2 оперирует только простыми матричными структурами, он требует очень небольшого объема памяти — определенно меньше $40n$ байт, необходимых для алгоритма 12.1.1 даже в его самом экономичном варианте [94].

Интересно сравнить результаты, выводимые алгоритмами Крочемора и Мейна-Лоренца. В алгоритме Крочемора каждая кратная строка выводится в нормальной форме, используя кодирование (i, p^*, r^*) . Однако напомним, что хотя алгоритм Крочемора обрабатывает каждый квадрат в отдельности, он объединяет смежные квадраты до тех пор, пока не будет получен максимальный показатель степени, который может превышать значение 2. С другой стороны, алгоритм МЛ выводит только квадраты, но, принимая во внимание лемму 12.1.5 и ее аналоги для массивов LP' и LS' , он выводит их в блоках, которые близки к понятию “серий”,

введенному в разделе 2.3. Например, рассмотрим строку

$$\begin{array}{ccccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\
 x = & a & b & a & a & b & a & b & a & a & b & a & b & a
 \end{array}$$

Для кратных строк с периодом $p = 5$, алгоритм Крочемора выдаст результаты:

$$(1, 5, 2), (2, 5, 2), (3, 5, 2), (4, 5, 2),$$

соответствующие квадрату $(abaab)^2$ и его трем перестановкам, полученным в результате циклического сдвига, $(baaba)^2$, $(aabab)^2$ и $(ababa)^2$. Но алгоритму МЛ в этом случае необходимо вывести только тройку

$$(p, \text{первая}, \text{последняя}) = (5, 10, 13),$$

указывающую на то, что квадраты с периодом 5 заканчиваются на позициях 10–13.

Из этого видно, что алгоритм МЛ превосходит понятие серии, которое было введено в [168] и применено там для уменьшения асимптотической сложности проблемы кратных строк. Мы исследуем эти вопросы в следующем разделе.

Упражнения 12.1

1. Докажите следующую “декомпозицию” леммы 12.1.2: ни одна из позиций строки x не может входить в более чем $\lceil \log_2(n - \alpha' + 2) \rceil$ малые последовательности, где α' — количество различных букв в строке x .
2. Вектор Париха (Parikh), или частотный вектор строки $x[1..n]$, которая содержит α' различных букв, — это целочисленный вектор $\phi_x = \phi_x[1..\alpha']$, где $\phi_x[h]$ равно числу вхождений h -й буквы в строку x [193]. Выразите сложность алгоритма Крочемора через n и вектора Париха. Используйте это выражение как основу для обсуждения поведения алгоритма Крочемора в тех случаях, когда α' велико.
3. В этом разделе часто использовался термин “последовательность”, но без его формального определения. Имея в виду, что о последовательности можно думать как своего рода о строке на основе индексированного алфавита (раздел 4.1), предоставьте отсутствующие определения для последовательности, подпоследовательности и, возможно, для других подобных терминов.
4. Покажите, что для выполнения декомпозиции всех последовательностей уровня ℓ на подпоследовательности уровня $\ell + 1$ необходимо не более $2n$ индексов последовательностей. Можно ли более точно оценить эту верхнюю границу?

5. Решите неравенства (12.2) и (12.3) и покажите, что $t_n \in O(n \log n)$. Затем на основе неравенств (12.4) докажите, что $t_n \in \Omega(n \log n)$.

Совет. Если для некоторого неотрицательного числа r выполняются неравенства $2^r < n < 2^{r+1}$, тогда сформулированные выше утверждения будут выполняться для n в том случае, если они выполняются для 2^{r+1} .

6. Предположим, что неравенство (12.2) обобщено в виде неравенства

$$t_n \leq k_2 f_n + 2t_{\lceil n/2 \rceil}, f_n \in \Omega(n^2),$$

справедливого для любого $n \geq 1$. Одновременно выполняется неравенство (12.3). Покажите, что в этом случае $t_n \in O(f_n)$.

Замечание. Мы используем этот результат в разделе 13.3 для того, чтобы оценить сложность алгоритма Шмидта для аппроксимирующих раппортов.

7. Покажите, что второе условие леммы 12.1.5 эквивалентно неравенству $i \leq p + \text{LP}[p + 1]$.
8. Предположим, что массивы LP и LS вычислены заранее для строк $u[1..n_1]$ и $v[1..n_2]$. Запишите алгоритм вычисления всех правых квадратов в строке uv с временем выполнения $\Theta(n_2)$.
9. Докажите следующий аналог леммы 12.1.5.

Пусть $u[1..n_1]$ и $v[1..n_2]$ — произвольные строки, а p ($p < n_1$) и i — положительные целые числа, такие, что $n_1 - 2p < i - 1 < n_1 - p$. Тогда в строке uv имеется квадрат длиной $2p$, начинающийся с i -й позиции строки u , только в том случае, если выполняются неравенства

$$(n_1 - p) - \text{LS}'[n_1 - p] \leq i - 1 \leq (n_1 - 2p) + \text{LP}'[(n_1 - p) + 1].$$

10. Предположим, что для строк $u[1..n_1]$ и $v[1..n_2]$ предварительно вычислены массивы LP' и LS' . Запишите алгоритм вычисления всех левых квадратов в строке uv с временем выполнения $\Theta(n_1)$.
11. Используя алгоритмы вычисления массивов и правых и левых квадратов с линейным временем выполнения, запишите процедуру СЗ, которая вычисляет все квадраты класса СЗ в строке uv , где $u = u[1..n_1]$, $v = v[1..n_2]$.
12. Дайте формальное доказательство равенства (12.5).
13. Докажите, что алгоритм 12.1.3 действительно может работать со значениями $k \in 2..i - 1$, при которых достигается максимум выражения $k + \text{LP}[k]$.
14. Аналогично алгоритмам вычисления массивов LP и LP' , запишите алгоритмы для вычисления массивов LS и LS' .

12.2 Серии

В этом разделе мы рассмотрим два подхода к проблеме кратных строк, которые уменьшают время их обработки и вывода результатов путем вычисления серий (i, p^*, r^*, t) , а не кратных строк (i, p^*, r^*) . В обоих подходах используется s -факторизация, вычисление которой с использованием деревьев суффиксов было описано в разделе 6.3. Кроме того, оба алгоритма используют массивы LS, LS', LP, LP', представленные в подразделе 12.1.2 при описании алгоритма МЛ. Второй алгоритм, алгоритм Колпакова–Кучерова (сокращенно, алгоритм КК), является по существу расширением первого алгоритма, алгоритма Мейна, при этом алгоритм КК вычисляет все серии, а алгоритм Мейна — только *крайние левые* серии, т.е. только крайние левые вхождения каждой отдельной серии в строку x .

Поскольку оба алгоритма зависят от использования деревьев суффиксов, то для их выполнения в случае произвольного упорядоченного алфавита, неизвестного заранее, потребуется время порядка $O(n \log n)$. Но для небольшого упорядоченного алфавита, известного заранее, обоим алгоритмам, скорее всего, потребуется одинаковое время порядка $\Theta(n)$. Как отмечается в работе [139], основная сложность этих алгоритмов заключается в “памяти, занятой структурой данных, необходимой для вычисления s -факторизации”.

Для того чтобы алгоритмы Мейна и КК выполнялись за линейное (относительно длины строки) время, размер вывода должен быть линейным (относительно количества серий). Как мы увидим, для алгоритма Мейна это следует из того факта, что сам процесс вычисления, который определяет крайние левые серии, будет линейным. Это также подтверждается более поздним результатом [88], цитируемым ниже, где доказано, что число различных (следовательно, и крайних левых) квадратов линейно. Гораздо более трудной проблемой является доказательство того, что общее количество серий, крайних левых или других, является линейным (относительно длины строки).

Теорема 12.2.1. Пусть $\rho(n)$ обозначает максимальное количество серий, которые могут содержаться в произвольной строке длиной n , определенной на произвольном алфавите. Существуют положительные константы k_1 и k_2 , независимые от n , такие, что для каждого целого $n \geq 1$ справедливо неравенство

$$\rho(n) \leq k_1 n - k_2 \sqrt{n} \log_2 n. \quad (12.11)$$

Доказательство см. в [139]. ■

Мы не приводим доказательство этого значительного результата, поскольку оно является очень длинным и техническим. Хотя теорема 12.2.1 является существенным продвижением в изучении периодичности в строковых последовательностях, она также является показателем того, сколько еще осталось неизученного.

Доказательство неравенства (12.11) не позволяет сделать какую-либо оценку относительно величины констант k_1 и k_2 — k_1 может равняться 10^9 , а k_2 может быть равно 10^{-9} . В то же время исчерпывающие компьютерные эксперименты до $n = 60$ [138] показали такие результаты:

1. $\rho(n) < n$;
2. $0 \leq \rho(n+1) - \rho(n) \leq 2$;
3. среди строк, определенных на алфавите $\{a, b\}$, существует строка без кубов, в которой содержится $\rho(n)$ серий.

Кажется вполне вероятным, что все эти фундаментальные результаты являются истинными, но, похоже, что никто не имеет какого-либо представления, как их доказать!

В недавно опубликованной статье [88] вместо $\rho(n)$ была рассмотрена величина $\sigma(n)$, равная максимальному количеству различных квадратов, которые могут содержаться в любой строке длиной n , и было показано, что $\sigma(n) \leq 2n - 8$ для всех $n \geq 5$. Как предполагается в [139], между $\rho(n)$ и $\sigma(n)$ может быть некоторая связь, которая может помочь в определении более точных границ величины ρ . В работе [140] доказан один более сильный результат, чем в теореме 12.2.1, а именно: там показано, что сумма показателей всех серий в любой строке длиной n также линейна относительно n .

12.2.1 Крайние левые серии — алгоритм Мейна

Крочемор [70], кажется, был первым, кто осознал, что s -факторизацию [156], использованную в работе [235] для сжатия строк (раздел 6.3), можно также применить для эффективного вычисления кратных строк. В этом подразделе мы описываем алгоритм Мейна [168], который использует s -факторизацию вместе с методами, разработанными в [176] (подраздел 12.1.2) для вычисления крайнего левого вхождения каждой отдельной серии в произвольной строке x . Основой этого алгоритма является следующая теорема.

Теорема 12.2.2. Предположим, что строка $x[\ell..n]$ имеет s -факторизацию $w_1 w_2 \dots w_k$, а r обозначает крайнее левое вхождение серии $x[i..j]$ в строку x , конечная позиция j которой находится в факторе w_h для некоторого $h \leq k$. Тогда

- а) начальная позиция i серии r находится в факторе $w_{h'}$ для некоторого $h' \leq h$;
- б) если $r = r_L u_R$ для некоторого префикса r_R фактора w_h , то

$$|r_L| < 2|w_{h-1}| + |w_h|.$$

Доказательство. Предположим, что утверждение а) не выполняется. Тогда серия r будет подстрокой фактора w_h , и поэтому $|w_h| \geq 2$. В этом случае по

определению s -факторизации должно существовать другое вхождение w_h в строку x слева от фактора w_h . Но тогда существует слева другое вхождение серии r , что противоречит условию теоремы. Таким образом, утверждение a) должно выполняться.

На основании утверждения a) можно считать, что $r = r_L u_R$ для некоторой непустой подстроки r_L и некоторого непустого префикса r_R фактора w_h . Предположим, что утверждение b) не выполняется, поэтому $|r_L| \geq 2|w_{h-1}| + |w_h|$. Тогда по крайней мере половина серии r находится слева от фактора w_{h-1} . Рассмотрим суффикс $v = w_{h-1} u_R$ серии r , как показано на рис. 12.4.

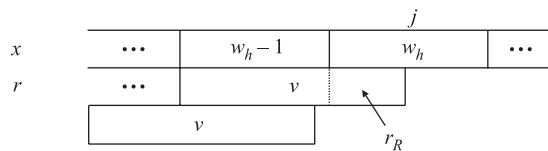


Рис. 12.4. Периодический суффикс v крайней левой серии r

Поскольку v является суффиксом серии r с периодом p^* и показателем $r^* \geq 2$, а $|v|$ не превышает половины длины серии r , то существует другое вхождение v слева на расстоянии p^* , как показано на рис. 12.4. Но поскольку w_{h-1} — собственный префикс v , а v — это повторяющаяся подстрока, то w_{h-1} не может быть повторяющейся подстрокой максимальной длины, необходимой для s -факторизации, что является противоречием с условием теоремы. Отсюда заключаем, что утверждение b) выполняется. ■

Фактически, эта непритязательная теорема предоставляет намного больше информации, чем кажется на первый взгляд. В частности, заметим, что доказательство утверждения b) теоремы требует только того, чтобы серия r , не будучи подстрокой фактора w_h , заканчивалась в w_h , т.е. она не должна быть крайней левой серией в определенном выше смысле. Таким образом, все серии строки x можно разделить на два непересекающихся класса.

- Класс I содержит серии, которые являются подстроками одного фактора w_h , $2 \leq h \leq k$, полученного при s -факторизации строки x .
- Класс II содержит серии, которые заканчиваются в некотором факторе w_h , $h \geq 2$, и начинаются в некотором предыдущем факторе $w_{h'}$, $h' \leq h$.

В этом подразделе покажем, как вычислить все серии класса II, серии класса I рассмотрим в следующем подразделе.

Проявив некоторую интуицию относительно серий класса II, рассмотрим наш универсальный пример строки Фибоначчи f_6 , s -факторизация которой имеет вид

$$w_1 w_2 \dots w_6 = a/b/a/aba/baaba/ab.$$

Тогда легко связать серии класса II с соответствующими факторами следующим образом:

$$(aba)^2(w_1 - w_4), a^2(w_3 - w_4), (ab)^2a(w_3 - w_4), \\ (abaab)^2a(w_1 - w_5), (aba)^2ab(w_4 - w_6).$$

Отметим, что серия класса I aa из w_5 в этот список не включена.

Алгоритм 12.2.1 (Алгоритм Мейна)

- ▷ Вычисление серий класса II в строке $x = x[\ell..n]$
- ▷ Выполнение алгоритма 6.3.1
- вычисление s -факторизации $x = w_1w_2 \dots w_k$
- $j^* \leftarrow 0$;
- for** $h \leftarrow 2$ **to** k **do**
 - ▷ Вычисление префикса r_L^* максимальной длины серии,
 - ▷ заканчивающейся в w_h (теорема 12.2.2)
 - $j^* \leftarrow j^* + |w_{h-1}|$
 - $\ell \leftarrow \min\{j^*, 2|w_{h-1}| + |w_h|\}$
 - $r_L^* \leftarrow x[j^* - \ell + 1..j^*]$
 - ▷ Вычисление серий класса II, заканчивающихся в w_h
 - $calcruns(r_L^*, w_h)$

Заметим, что на основании утверждения б) теоремы 12.2.2 длина префикса r_L^* каждой серии класса II ограничена значением $2|w_{h-1}| + |w_h|$. Как оказывается, это ограничение позволяет применить метод алгоритма МЛ для вычисления всех серий класса II, заканчивающихся в любом факторе w_h , $h \geq 2$, за время, пропорциональное $2|w_{h-1}w_h|$. Поскольку

$$2n \leq \sum_{h=2}^k 2|w_{h-1}w_h| < 4n, \tag{12.12}$$

отсюда следует, что все серии класса II в строке x можно вычислить за время $\Theta(n)$ на основе предварительно вычисленной s -факторизации. Поэтому мы вправе формулировать следующую теорему.

Теорема 12.2.3. При наличии ранее вычисленной s -факторизации $w_1w_2 \dots w_k$ строки $x = x[\ell..n]$ и при условии линейного времени выполнения процедуры $calcruns$ алгоритм 12.2.1 вычисляет все серии класса I в строке x за время порядка $\Theta(n)$. ■

Сделаем важное замечание. Как мы видим, каждый запрос к программе $calcruns(u, v)$ приводит к выводу не более чем $|u| + |v|$ серий класса II. Тогда оценка верхней границы, задаваемая неравенствами (12.12), применима также

к общему количеству серий класса II. Таким образом, общее количество серий класса II не может превышать $4n$. Это более удобная оценка, чем задаваемая теоремой 2.2.1. Но эта оценка, к сожалению, также не полна, поскольку не учитывает серии класса I.

Чтобы завершить описание алгоритма Мейна, осталось объяснить работу процедуры $calcruns(u, v)$, вычисляющей в строке uv все серии, которые не принадлежат полностью ни u , ни v . Пусть $u = u[1..n_1]$ и $v = v[1..n_2]$. Напомним, что в подразделе 12.1.2 мы различали (и отдельно вычисляли) правые и новые квадраты в зависимости от того, находится ли второе вхождение образующей квадрата полностью в пределах v . Здесь вводится подобное различие, разделяя процедуру $calcruns$ на две отдельные процедуры $calcR$ и $calcL$, которые вычисляют “правые максимальные” и “левые максимальные” серии соответственно, определяемые следующим образом.

- **Правая максимальная серия** — эта серия класса II из uv имеет суффикс в v , по меньшей мере, длиной p , где p — период этой серии.
- **Левая максимальная серия** — эта серия класса II из uv имеет непустой суффикс в v длиной $\ell < p$.

В приведенном выше примере строки f_6 все перечисленные серии, за исключением $(aba)^2ab$, являются правыми максимальными сериями; $(aba)^2ab$ — это левая максимальная серия, поскольку она имеет период $p = 3$.

Аналогия с подразделом 12.1.2 продолжается. Определим два массива, соответствующие правым максимальным сериям, следующим образом.

- Массив $LP[2..n_2 + 1]$: для каждой позиции $i = 2, 3, \dots, n_2$ в v $LP[i] = \text{lcp}(v[i..n_2], v)$, тогда как для $i = n_2 + 1$ $LP[i] = 0$.
- Массив $LS[\ell..n_2]$: для каждой позиции $i = 1, 2, \dots, n_2$ в v $LS[i] = \text{lcs}(uv[\ell..i], u)$.

Определение массива LP идентично тому, которое было дано в подразделе 12.1.2; определение массива LS немного изменено — вместо суффикса $v[1..i]$ в нем рассматривается суффикс $uv[1..i]$.

Подобным образом определяются массивы LP' и LS' , соответствующие левым максимальным сериям.

- Массив $LP'[2..n_1]$: для каждой позиции $i = 2, 3, \dots, n_1$ в u $LP'[i] = \text{lcp}(u[i..n_1]v, v)$.
- Массив $LS'[\ell..n_1 - 1]$: для каждой позиции $i = 1, 2, \dots, n_1 - 1$ в u $LS'[i] = \text{lcs}(u[\ell..i], u)$.

В оставшейся части этого подраздела мы ограничимся описанием только процедуры $calcR$ и ее массивов LP и LS , оставляя процедуру $calcL$ вместе с массивами LP' и LS' для упражнений 12.2.3 и 12.2.4. Кроме того, опустим описание алгоритмов вычисления массивов LP и LS , поскольку они отличаются только

незначительными деталями от аналогичных алгоритмов, описанных в подразделе 12.1.2. Далее будем предполагать, что каждый из этих четырех массивов можно вычислить за время, пропорциональное их размерам.

Чтобы понять работу процедуры *calcR*, необходима следующая лемма, подобная лемме 12.1.5.

Лемма 12.2.4. Пусть $u[\ell..n_1]$ и $v[\ell..n_2]$ — произвольные строки, а $r = r_L r_R$ — подстрока строки uv , такая, что r_L является непустым суффиксом строки u , а r_R — непустым префиксом строки v . Предположим, что $|r_R| \geq p$ для некоторого целого $p \in 1.. \lfloor |r|/2 \rfloor$. Тогда r имеет период p тогда и только тогда, когда выполняются следующие условия:

- а) $|r_L| \leq \text{LS}[p]$;
- б) $|r_R| \leq p + \text{LP}[p + 1]$.

Доказательство. Для того чтобы серия r имела период p , необходимо и достаточно, чтобы серия r имела грань длиной $r - p$, где $r = |r|$. В соответствии с условиями теоремы имеем $r - p \geq r - |r_R| = |r_L|$. Поэтому эта грань будет существовать тогда, когда выполняются условия

- А) $r[1..|r_L|] = r[p + 1..p + |r_L|]$;
- Б) $r[|r_L| + 1..r - p] = r[p + |r_L| + 1..r]$.

Условие А) эквивалентно утверждению, что подстрока $r[p + 1..p + |r_L|]$ является суффиксом u . Но эта подстрока также является суффиксом подстроки $uv[1..p]$. Поэтому условие А) выполняется только тогда, когда длина $|r_L|$ этого суффикса не превышает $\text{LS}[p]$. Таким образом, условие А) эквивалентно условию а).

Условие Б) эквивалентно утверждению, что подстрока $r[p + |r_L| + 1..r]$ является префиксом v . Но эта подстрока является префиксом (возможно, пустым) подстроки $v[p + 1..n_2]$. Поэтому условие Б) выполняется только тогда, когда длина $r - p - |r_L|$ этого префикса не превышает $\text{LP}[p + 1]$, что эквивалентно условию б). ■

Заметим, что если подстрока r , описанная в этой лемме, существует, то она может быть записана в следующем виде:

$$r = (r[1..p])^k r[1..p'] \tag{12.13}$$

для некоторых целых $k \geq 2, p' \in 0..p-1$. Кроме того, лемма говорит, что существует максимальной длины подстрока $r^* = r_L^* r_R^*$ с левой частью r_L^* , удовлетворяющей условию

$$|r_L^*| = \text{LS}[p] > 0, \tag{12.14}$$

и с правой частью r_R^* , удовлетворяющей условию

$$|r_R^*| = p + \text{LP}[p + 1]. \tag{12.15}$$

Из определения 2.3.3 следует, что серия должна быть представима в форме (12.13), но, кроме того, она должна быть непродолжаемой, т.е. она не может быть собственной подстрокой любой другой подстроки с тем же периодом p . Поскольку, в соответствии с теоремой 12.2.2 и нашим выбором uv можем предполагать, что каждая серия класса II не может быть продолжена влево или вправо за пределы uv , то, следовательно, подстрока r^* , определяемая условиями (12.14) и (12.15), имеет максимальную длину и непродолжаема в пределах строки x , поэтому она является серией строки x . Более того, она, по определению, является правой максимальной серией.

Таким образом, имея заранее вычисленные значения $LP[p + 1]$ и $LS[p]$, соответствующие u и v для каждого $p \in 1..n_2$, можно в uv выделить каждую правую максимальную серию r^* , применяя условия (12.14) и (12.15) вместе с условием $|r^*| \geq 2p$. Эти три условия эквивалентны неравенствам

$$LS[p] > 0 \text{ и } LS[p] + LP[p + 1] \geq p.$$

На этой основе строится следующий алгоритм.

Алгоритм 12.2.2 (Процедура *calcR*)

▷ Вычисление в $uv = u[\ell..n_1]v[\ell..n_2]$ всех правых максимальных серий
вычисление $LP[2..n_2 + 1]$, $LS[1..n_2]$ для uv

for $p \leftarrow 1$ **to** n_2 **do**

if $\{LS[p] > 0 \text{ and } LS[p] + LP[p + 1] \geq p\}$ **then**

$i \leftarrow n_1 - LS[p] + 1$ ▷ начало серии в u

$j \leftarrow p + LP[p + 1]$ ▷ конец серии в v

output i, j

На выполнение цикла **for** алгоритма процедуры *calcR* требуется время $\Theta(n_2)$, а для вычисления массивов LP и LS , как показано в упражнении 12.2.2, необходимо время $\Theta(n_1 + n_2)$. Таким образом, время выполнения процедуры *calcR*(u, v) линейно относительно $|uv|$ и требует памяти, объем которой линеен относительно $|v|$. Аналогично время выполнения процедуры *calcL*(u, v) линейно относительно $|uv|$ и требует памяти, объем которой линеен относительно $|u|$. Следовательно, процедура *calcruns*(u, v) выполняется за время порядка $\Theta(|u| + |v|)$ с использованием памяти такого же объема, при этом теорема 12.2.3 сохраняет силу, а алгоритм Мейна определяет все серии класса II в строке $x[\ell..n]$ за время $\Theta(n)$ (при условии, что s -факторизация вычислена заранее).

Необходимо отметить, что алгоритм Мейна может вывести одну и ту же серию класса II более одного раза, если серия имеет более одного периода. Рассмотрим, например, строку $x = (ab)^3(aba)^4a$ с s -факторизацией

$$w_1w_2w_3w_4 = a/b/ababab/abaabaaba.$$

Правая максимальная серия $(aba)^4a$, соответствующая фактору w_4 , будет выведена дважды: один раз для $p = 3$ (образующая aba) и второй раз для $p = 6$ (образующая $(aba)^2$). Из этого видно, что период, использованный для выделения серии, не всегда должен быть минимальным, следовательно, образующая серии может иногда сама быть кратной строкой. Но даже с этим потенциальным дублированием серий, общее количество выводов должно все еще быть меньше чем $4n$.

Отметим также, что в соответствии с необходимым условием алгоритма процедуры *calcR*

$$LS[p] + LP[p + 1] \geq 1$$

в дополнение ко всем сериям мы выводим все минимальной длины непримитивные подстроки с периодом p , удовлетворяющие условиям $|r_L^*| \geq 1$, $|r_R^*| \geq p$. Однако количество выводов, следовательно, и общее количество таких подстрок должно быть меньше $4n$.

Как отмечалось в [168], “дальнейшая модификация может позволить алгоритму находить все [серии класса I] простым способом”. В следующем подразделе мы увидим, что это действительно может быть сделано также за время $\Theta(n)$.

12.2.2 Все серии — алгоритм Колпакова и Кучерова

В этом подразделе кратко описывается модификация [139] алгоритма 12.2.1, которая вычисляет серии класса I в произвольной строке $x[\ell..n]$ за время $\Theta(n)$. Таким образом, первоначальный алгоритм Мейна и эта модификация позволяют вычислить все серии в строке x за линейное время, предполагая, что заранее вычислена s -факторизация строки x и, следовательно, дерево суффиксов T_x строки x (раздел 6.3).

Алгоритм Колпакова–Кучерова (сокращенно, алгоритм КК) представляет каждую серию класса II $r = x[i, j]$, найденную с помощью алгоритма Мейна, как пару чисел (i, j) , где i — начальная позиция, а j — конечная позиция серии в строке x . (В случае, если алгоритм Мейна нашел больше одной серии (i, j) , дубликаты удаляются, как мы это вскоре увидим.) Формируется массив OCCURS $[\ell..n]$, в котором OCCURS $[i]$ — это указатель на связный список позиций $j_{i1}, j_{i2}, \dots, j_{ir_i}$, которые первоначально являются конечными позициями серий класса II, начинающихся с i -й позиции строки x . Будем придерживаться соглашения, что

$$j_{i1} < j_{i2} < \dots < j_{ir_i}. \tag{12.16}$$

Если ни одна из серий не начинается в i -й позиции, то значение OCCURS $[i]$ полагается равным NULL.

Чтобы вычислить массив OCCURS и его связанные списки, сначала необходимо выполнить блочную сортировку пар (i, j) , в результате которой они будут размещаться в порядке возрастания значений j . Затем массив OCCURS модифи-

цируется путем обеспечения доступа к отсортированным парам в порядке их возрастания и дополнения каждого значения i соответствующим значением j в конце списка, на который указывает OCCURS[i], — это является гарантией того, что условие (12.16) будет выполняться. Поскольку количество серий класса Π равно $O(n)$ и поскольку i и j являются целыми числами из интервала $1..n$, то блочную сортировку и формирование OCCURS можно выполнить за время $\Theta(n)$. Разумеется, дубликаты пар (i, j) можно легко устранить в процессе сортировки.

Для того чтобы сделать эффективным использование этой структуры данных, для каждого фактора w_h ($1 \leq h \leq k$) необходимо сохранять две величины: его начальную позицию i_h и начальную позицию $i'_h < i_h$ предыдущей копии w_h в строке x . В случае, когда w_h является крайним левым вхождением некоторого символа, тогда не существует предыдущей копии w_h , и поэтому полагаем $i'_h \leftarrow 0$. Итак, создается два массива $I = I[\ell..k + 1]$ (с $I[k + 1] = n + 1$) и $I' = I'[\ell..k]$ для сохранения этой промежуточной информации, определяемой по первоначально вычисленной s -факторизации, значения массива I' легко вычисляются из дерева суффиксов.

Используя эти структуры данных, можно определить, что копия w_h (если она существует) расположена в позиции $I'[h] < I[h]$. Затем можно переместить все серии подходящей длины из копии непосредственно в фактор w_h . Поскольку для каждой начальной позиции i доступ к сериям из массива OCCURS[i] осуществляется в порядке возрастания их длины, перемещение, соответствующее каждой i , можно реализовать за время, пропорциональное количеству перемещаемых серий. Так как, согласно теореме 12.2.1, общее количество серий равно $O(n)$, поэтому дополнительное время, необходимое для сохранения серий класса I в структуре данных, также равно $O(n)$. Детали этого процесса представлены в алгоритме 12.2.3. Основным результатом данного раздела сформулирован в теореме 12.2.5.

Алгоритм 12.2.3 (Алгоритм КК)

▷ Заданы массивы $I[\ell..k + 1]$ и $I'[\ell..k]$, которые описывают
 ▷ s -факторизацию $w_1 w_2 \dots w_k$ строки $x[1..n]$, и заданы списки массива
 ▷ OCCURS[$1..n$], которые определяют (в порядке возрастания
 ▷ их длины) все серии класса Π , встречающиеся в каждой позиции i .
 ▷ Алгоритм модифицирует списки для включения всех серий строки x

```

for  $h \leftarrow 2$  to  $k$  do
    if  $I'[h] > 0$  then
         $\delta \leftarrow I[h] - I'[h]$     ▷ смещение фактора  $w_h$  относительно его копии
        for  $i \leftarrow I[h]$  to  $I[h + 1] - 1$  do
             $\forall j \in \text{список}(\text{OCCURS}[i - \delta])$ 
            ▷ Серия класса  $I$  должна начинаться и заканчиваться в  $w_h$ ,
            if  $(j + \delta) - i < I[h + 1] - i$  then
                вставка списка (OCCURS[ $i$ ])  $\leftarrow j + \delta$ 
    
```

Теорема 12.2.5. Если s -факторизация $w_1 w_2 \dots w_k$ строки $x[\ell..n]$ вычислена заранее, то массив OCCURS $[\ell..n]$ и связанные с ним списки вместе с массивами $I[\ell..k + \ell]$ и $I'[\ell..k]$ можно вычислить за время порядка $\Theta(n)$. На основе этих структур данных алгоритм 12.2.3, модифицирующий массив OCCURS и связанных с ним списков с целью включения всех серий строки x , выполняется за время порядка $O(n)$. ■

Имеется несколько других алгоритмов, которые вычисляют наборы кратных строк в $x[\ell..n]$ за время $\Theta(n)$ на алфавите размером α . Во всех из них используются деревья суффиксов, и поэтому в оценках времени выполнения присутствует множитель $\log \alpha$. В [141] описывается алгоритм, который вычисляет самый короткий квадрат в каждой позиции в x . На основе упомянутого ранее в [88] результата, что максимальное количество различных квадратов в любой строке равно $O(n)$, в работе [107] предлагается алгоритм вычисления всех различных квадратов в строке x . В [108] описан алгоритм, который вычисляет все отличающиеся друг от друга непродолжаемые (определение 2.3.2) повторяющиеся подстроки в x : эти подстроки — различные образующие не только кратных строк, но и раппортов, которые могут быть расщеплением (непоследовательным) или покрытием (определение 2.3.1). Вычисление всех пар непродолжаемых повторяющихся подстрок, которые удовлетворяют различным ограничениям, описано в [41].

В заключение скажем, что в упражнении 12.2.10 мы предлагаем другой подход к вычислению серий класса I — подход, сложность которого пока не совсем ясна!

Упражнения 12.2

1. Докажите неравенства (12.12).
2. Определение массива LP и алгоритм его вычисления даны в подразделе 12.1.2. Массив LS немного отличается от массива LP. Запишите алгоритм вычисления массива LS для произвольных строк u и v и покажите, что для этого требуется время порядка $\Theta(|u| + |v|)$.
3. Запишите алгоритм вычисления массива LP' для произвольных строк u и v , выполняемый за время $\Theta(|u| + |v|)$.
4. Сформулируйте и докажите аналог леммы 12.2.4 для левых максимальных серий. На этой основе опишите алгоритм процедуры *calcL*.
5. Предположим, что подстрока $x[i..j]$ — это серия с периодом p в строке x . Докажите, что для любого значения $i' \in i + p..j - p + 1$ ни одна подстрока в строке x , начинающаяся в позиции i' , не может быть серией с периодом p . Покажите на примере, что границы для таких значений i' являются неулучшаемыми.

6. Предположим, что подстрока $x[i..j]$ — это серия с периодами p и p' ($p' \neq p$) в строке x . Докажите, что наибольший общий делитель чисел p и p' также является периодом подстроки $x[i..j]$.
Совет. Не забудьте о лемме периодичности!
7. Объясните, каким образом можно вычислить массив I' за линейное время на основе дерева суффиксов T_x во время выполнения s -факторизации строки x .
8. Пусть h — наименьшее целое число, при котором фактор w_h в s -факторизации строки x содержит конечную позицию серии. Является ли необходимым, чтобы каждая серия в факторе w_h была серией класса II?
9. Предположим, что фактор w_h в s -факторизации строки x начинается с позиции i и что он является копией подстроки x , начинающейся с позиции $i' < i$. Далее предположим, что $i - i' \leq |w_h|/2$. Покажите, что подстрока $[i'..i + |w_h| - 1]$ является серией строки x .
10. Заметим, что многие серии класса I, которые располагаются в пределах какого-либо фактора w_h в s -факторизации строки x , сами являются сериями класса II в s -факторизации фактора w_h . Используйте этот факт при разработке алгоритма, который рекурсивно обращается к алгоритму Мейна и вычисляет все серии в строке x . Что можно сказать о сложности этого алгоритма?
11. Пусть $z_0 = (\lambda_0)^2$ и $z_i = z_{i-1}\lambda_i z_{i-1}$ для всех $i = 1, 2, \dots, n$, где символы $\lambda_1, \lambda_2, \dots, \lambda_{n-1}$ попарно различаются. Как предполагает Роман Колпаков, строка z_n представляет худший случай для рекурсивного алгоритма, описанного в предыдущем упражнении. На основе s -факторизации z_n определите временную сложность применяемого к z_n алгоритма.

ГЛАВА 13

Обобщение периодичности

Для мудреца достаточно одного слова,
а множество слов не заполнят сосуд.

— Бенджамин Франклин (1706–1790).
Альманах Бедного Ричарда

В этой главе мы обсуждаем темы, являющиеся прямыми обобщениями тех тем, которые возникают из понятия периодичности. Если говорить кратко, то основными проблемами, связанными с периодичностью, являются следующие две:

- вычисление всех периодов произвольной строковой последовательности $x = x[1..n]$;
- вычисление всех кратных строк в строке x .

Первая из этих проблем решается с помощью алгоритма вычисления массива граней (раздел 1.3). Это первый алгоритм, который изучается в этой книге: он вычисляет за время порядка $\Theta(n)$ массив граней $\beta[1..n]$ строки x , который определяет каждую грань, следовательно, и каждый период не только строки x , но и каждого префикса x . Массив граней, как обсуждалось в разделе 5.1, соответствует структуре дерева, называемого деревом граней. Корневой узел дерева помечен нулем, а другие его узлы помечены позициями $i \in 1..n$ букв строки x . Метки узлов на пути от корня дерева к любому узлу i (за исключением самого узла i) являются точными длинами (в возрастающем порядке) граней подстрок $x[1..i]$.

В разделе 13.1 приведен довольно примечательный факт, что квазипериоды строки x (см. определения раздела 2.3) можно вычислить с использованием структур, формально идентичных тем, которые применяются для вычисления периодов:

- за время порядка $\Theta(n)$ можно вычислить *массив оболочек* $\gamma[1..n]$, который определяет каждую оболочку (т.е. каждый квазипериод) каждого префикса строки x ;
- массив оболочек определяет *дерево оболочек* с узлами, помеченными числами из интервала $0..n$ (как и в случае дерева граней, метки, отличные от нуля, на пути от корня до любого узла i дерева оболочек — это точные длины (в возрастающем порядке) оболочек $x[1..i]$).

Как и о массиве граней, о массиве оболочек, наверное, следует думать как о внутреннем паттерне: каждая строка имеет массив оболочек. Однако из-за его тесной связи с периодичностью здесь мы рассмотрим его как характеристический паттерн.

В разделах 13.2 и 13.3 вместо рассмотрения естественных расширений понятия кратной строки обратимся к вычислению точных и приближенных раппортов (раздел 2.3). В разделе 13.2 показано, что, как и кратные строки, непродолжаемые раппорты вычисляемы за время порядка $O(n \log n)$ и, подобно сериям, требуют памяти порядка $O(n)$. Попутно обнаружим, что алгоритм Крочемора для вычисления всех кратных строк (подраздел 12.1.1) является, в конечном счете, алгоритмом построения дерева суффиксов, что дерево суффиксов компактно представимо в виде массива пар целых чисел и что вычисление непродолжаемых раппортов является естественным началом для вычисления всех непродолжаемых покрываемых подстрок (задача 2.17 — раздел 2.3).

В разделе 13.3 рассмотрим дальнейшее обобщение проблемы раппортов до задачи приближенных раппортов, которые используем для вычисления расстояния и для приближенного сравнения с паттерном (вместо понятий и методов, представленных в главах 9 и 10). Вычисление приближенных раппортов имеет особенно важное значение для молекулярной биологии, например для идентификации длинных последовательностей ДНК, которые повторяются в другом месте генома, но с незначительным изменением.

Наконец, в разделе 13.4 рассматриваются некоторые проблемы, которые возникают при определении приближенных образующих строковой последовательности.

13.1 Все оболочки — алгоритм Ли–Смита

Вспомним из раздела 2.3, что раппорт

$$M_{x,u} = (p; i_1, i_2, \dots, i_r) \tag{13.1}$$

является *оболочкой* подстроки $x[i_1..i_r + p - 1]$ тогда и только тогда, когда $r \geq 2$, $p = |u|$ и

$$u = x[i_1..i_1 + p - 1] = x[i_2..i_2 + p - 1] = \dots = x[i_r..i_r + p - 1],$$

а каждая *лакуна* (gap) $g_j = i_{j+1} - i_j$ ($j = 1, 2, \dots, r - 1$) лежит в интервале $1..p$. В этом разделе рассмотрим такую проблему: существует ли для произвольной строки $x = x[1..n]$ такая подстрока u , что раппорт (13.1) является оболочкой для случая, когда $i_1 = 1$ и $i_r + p - 1 = n$. Если нет необходимости в определении местонахождения r , в котором появляется подстрока u , то можно сказать, что подстрока u является оболочкой строки x или что u *покрывает* строку x . Например, для строки $x = abaababaaba$ подстроки $u = abaaba$ и $u' = aba$ являются оболочками строки x . Заметим, что поскольку в (13.1) значение r должно быть не меньше 2, то сама строка x никогда не будет своей оболочкой.

Перечислим некоторые элементарные факты о строке x и ее оболочке u , оставляя их доказательство для упражнения 13.1.1.

- F1.** $0 < |u| < |x|$.
- F2.** Подстрока u является гранью строки x .
- F3.** Пусть u' — такая строка, что $|u'| < |u|$. Тогда u' будет оболочкой u только в том случае, если u' будет оболочкой x .
- F4.** Подстрока u имеет грань длиной $p - g_j$, где g_j ($j = 1, 2, \dots, r - 1$) — произвольная лакуна, определенная в (13.1).

В 1990-х годах было опубликовано несколько алгоритмов вычисления оболочек строки x со временем выполнения порядка $\Theta(n)$: в [19] вычислялась только самая короткая оболочка x , затем в [40] — самая короткая оболочка каждого префикса, позже в [182, 181] — каждая оболочка строки x , но без оболочек ее собственных префиксов, наконец, в [160] вычислялась каждая оболочка каждого префикса x . Здесь опишем последний из этих алгоритмов, алгоритм Ли–Смита (Li–Smith, сокращенно — алгоритм ЛС), выполнение которого стало возможным благодаря введению понятия *массива оболочек* $\gamma = \gamma[1..n]$ строки x , который определяется следующим образом: для каждого $i \in 1..n$ элемент массива $\gamma[i]$ равен длине наибольшей оболочки подстроки $x[1..i]$; $\gamma[i] = 0$, если такой оболочки не существует. Приведем для сравнения пример массивов граней и оболочек.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
$x =$	a	b	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b	a	b	a	b	a
$\beta =$	0	0	1	1	2	3	2	3	4	5	6	4	5	6	7	8	9	10	11	7	8	2	3
$\gamma =$	0	0	0	0	0	3	0	3	0	5	6	0	5	6	0	8	9	10	11	0	8	0	3

Непосредственным следствием фактов F1 и F3 является то, что длины оболочек подстроки $x[1..i]$ являются элементами убывающей последовательности

$$\langle \gamma[i], \gamma^2[i], \dots, \gamma^{k-1}[i] \rangle, \tag{13.2}$$

Свойства массива оболочек

Наблюдательный читатель в процессе решения упражнения 2.2.6 уже убедился в том, что для произвольной строки $x[1..n]$ справедливо неравенство

$$d_H(x, R_j(x)) \neq 1, \tag{13.3}$$

где d_H — расстояние Хемминга (раздел 2.2), а $R_j(x)$ — j -й циклический сдвиг строки x , $j \in 0..n-1$ (раздел 1.4). Этот результат мы используем для доказательства леммы, которая, в свою очередь, будет использована для доказательства первых двух интересных свойств массивов граней.

Лемма 13.1.1. Если u является оболочкой строки x и $|u| = p$, тогда каждая подстрока v длиной p в строке x удовлетворяет неравенству

$$d_H(u, v) \neq 1.$$

Доказательство. Рассмотрим два различных вхождения u в $x = x[1..n]$, которые являются смежными или перекрывающимися. Обозначим левое вхождение u как $u_1 = x[i_1..i_1 + p - 1]$, а правое вхождение — как $u_2 = x[i_2..i_2 + p - 1]$, где $1 \leq i_1 \leq n - p + 1$ и $i_1 < i_2 \leq i_1 + p$.

Покажем, что утверждение леммы справедливо для каждой подстроки v строки $U = x[i_1..i_2 + p - 1]$, следовательно, и для самой строки x .

Сначала рассмотрим случай $i_2 = i_1 + p$. Здесь $U = u^2$, поэтому каждая подстрока v длиной p в U является циклическим сдвигом u . Тогда утверждение леммы является прямым следствием неравенства (13.3).

Теперь предположим, что $i_2 < i_1 + p$. В этом случае u_1 и u_2 перекрываются некоторой непустой гранью u' строки u . Тогда, как показано на рис. 13.2, можно записать $u = u'u''$ для некоторого непустого суффикса u'' строки u .

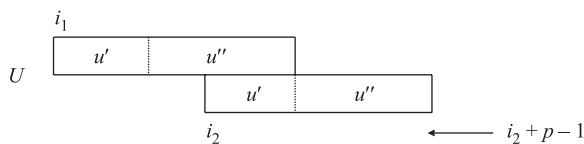


Рис. 13.2. Перекрывающиеся вхождения u в U

Пусть $u'' = w^k$ для некоторого наибольшего целого числа $k \geq 1$ (если u'' не является кратной строкой, то $w = u''$ и $k = 1$). Тогда U имеет период $|w|$, и поэтому каждая подстрока v'' длиной $|u''|$ в U также имеет период $|w|$. Таким образом, подстрока v'' должна удовлетворять одному из следующих двух условий.

- v'' является циклическим сдвигом u'' , но не равным u'' , и поэтому $d_H(u'', v'') \geq 2$ в соответствии с (13.3). Следовательно, для каждого вхождения $v = v'v''$ длиной p в U выполняется неравенство $d_H(u, v) \geq 2$.

- $v'' = u''$. В этом случае, поскольку v'' является подстрокой U с периодом $|w|$ и поскольку строка $u = u'u'' = u'v''$ входит в U , по меньшей мере, один раз, то каждому вхождению v'' в U предшествует вхождение u' , такое, что $d_H(u, u'v'') = 0$.

Таким образом, показано, что в U не существует ни одной подстроки v , такой, что $d_H(u, v) = 1$, что и требовалось доказать. ■

Этот результат впервые был доказан (очень изящно) в статье [210] путем использования варианта леммы о периодичности [32]; более слабая версия была сформулирована (с тяжеловесным определением) в [160]. С их помощью можно доказать следующее фундаментальное свойство массива оболочек $\gamma[1..n]$.

Теорема 13.1.2. Если $\gamma[i] \neq 0$, тогда или $\gamma[i + 1] > \gamma[i]$, или $\gamma[i + 1] = 0$ ($\forall i \in 1..n - 1$).

Доказательство. Пусть $p = \gamma[i]$. Тогда $p > 0$ и $x[1..p]$ является наибольшей оболочкой подстроки $x[1..i]$. Предположим, что теорема неверна, и поэтому существует положительное целое число $p' \leq p$, такое, что $x[1..p']$ будет наибольшей оболочкой подстроки $x[1..i + 1] = x[1..i]\lambda$. Отсюда следует, что $x[1..p']$ является оболочкой подстроки $x[1..p]\lambda$ (рис. 13.3).

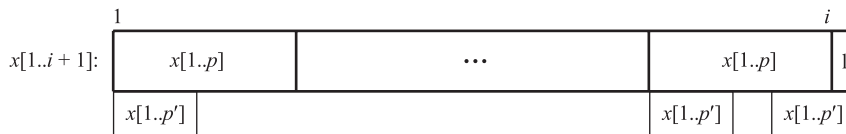


Рис. 13.3. Подстрока $x[1..p']$ является оболочкой строки $x[1..p]\lambda$

Если за каждым вхождением строки $x[1..p]$ в $x[1..i + 1]$ будет следовать буква λ , тогда $x[1..p + 1]$ будет оболочкой $x[1..i + 1]$, что входит в противоречие с предположением о том, что $x[1..p']$ является наибольшей оболочкой. Таким образом, за некоторым вхождением $x[1..p]$ в $x[1..i]$ должна следовать буква $\lambda \neq \lambda'$. Это говорит о том, что

- имеется, по меньшей мере, две различные буквы в строковой последовательности $x[1..i + 1]$, следовательно, и в $x[1..p']$, так что $p' \geq 2$;
- поскольку $x[1..p'] = x[1..p' - 1]\lambda$ является оболочкой, тогда в строке $x[1..i + 1]$ существует подстрока $u = x[1..p' - 1]$, такая, что $d_H(x[1..p'], u) = 1$. Это противоречит лемме 13.1.1.

Таким образом, предположение $0 < p' \leq p$ является неверным, что доказывает теорему. ■

Следующий результат связывает значения массива оболочек $\gamma[1..n]$ со значениями массива граней $\beta[1..n]$: значения $\gamma[i + 1]$ могут быть ненулевыми только тогда, когда $\beta[i + 1] = \beta[i] + 1$.

Теорема 13.1.3. Если $\beta[i + 1] \leq \beta[i]$, тогда $\gamma[i + 1] = 0$ ($\forall i \in 1..n - 1$).

Доказательство. Из факта F2 следует, что если $\beta[i + 1] = 0$, тогда $\gamma[i + 1] = 0$. Следовательно, без потери общности можем предположить, что $\beta[i + 1] > 0$. Поэтому, в соответствии с гипотезой, $\beta[i] > 0$.

Теперь предположим, что теорема неверна, поэтому $\gamma[i + 1] = p' > 0$. Тогда $x[1..p']$ является оболочкой подстроки $x[1..i + 1]$. Таким образом, положив $p = \beta[i]$, имеем

$$0 < p' = \gamma[i + 1] \leq \beta[i + 1] \leq p.$$

Обозначив $x[i + 1] = \lambda$, видим, что оболочка $x[1..p'] = x[i - p' + 2..i]\lambda$ является собственным суффиксом строковой последовательности

$$x[i - p + 1..i + 1] = x[i - p + 1..i]\lambda = x[1..p] = x[1..p + 1]. \quad (13.4)$$

Но поскольку $\beta[i + 1] \neq \beta[i] + 1$, то невозможно равенство $x[p + 1] = \lambda$, что противоречит (13.4). Поэтому утверждение $\gamma[i + 1] > 0$ не имеет места. ■

Из этой теоремы естественным образом вытекает понятие *ступени* (staircase) — максимальной последовательности $S_{i,h}$, состоящей из h номеров $i, i + 1, \dots, i + h - 1$ элементов массива граней $\beta[1..n]$, таких, что $\beta[j + 1] = \beta[j] + 1$ для каждого $j \in i..i + h - 2$. Отметим, что “максимальность” ступени означает, что $\beta[i] \leq \beta[i - 1]$ для $i - 1 \geq 1$ и $\beta[i + h - 1] \geq \beta[i + h]$ для $i + h \leq n$. Например, в массиве граней $\beta = 0011232345645$ имеется 5 ступеней: $S_{1,1}, S_{2,2}, S_{4,3}, S_{7,5}, S_{12,2}$.

Теорема 13.1.3 говорит, что $\gamma[i] = 0$ в начальной позиции i каждой ступени массива β . Позднее мы увидим, насколько полезным может быть понятие ступени для вычисления массива оболочек.

Еще одно полезное понятие — это понятие “живой” вершины. Если строка x является префиксом строки y , то будем говорить, что y является *правым расширением* строки x . Если префикс u строки x может быть оболочкой некоторого правого расширения, тогда об u говорят, что он *живой* (live) по отношению к x ; иначе об u говорят, что он *мертвый* (dead) по отношению к x . Например, если $x = abaab$, то $u = aba$ живой по отношению к x , поскольку aba покрывает правое расширение xa ; с другой стороны, $u = a$ и $u = ab$ не могут покрывать правое расширение x , и поэтому они по отношению к x мертвые. Заметим, что x всегда жива по отношению к самой себе.

Поскольку префикс определяется своей длиной и поскольку мы всегда говорим об одной данной строке x , расширим понятия “мертвый–живой” до отдельных позиций в строке x : будем говорить, что позиция $j \in 1..i$ жива по отношению

к позиции i тогда и только тогда, когда $x[1..j]$ может быть оболочкой некоторого правого расширения строки $x[1..i]$. Итак, поскольку позиции в x соответствуют узлам дерева оболочек, мы далее расширим эту терминологию до живых и мертвых узлов в дереве T_γ . Заметим, что если позиция j мертва по отношению к i , она также мертва и по отношению к $i + 1$. Таким образом, при просмотре строки x слева направо, т.е. при возрастании величины i , количество мертвых позиций будет монотонно неубывающей последовательностью. В строках, как и в реальной жизни, если умер, то это навсегда.

Лемма 13.1.4. Если $j \in i - \beta[i]$, то j жива по отношению к i .

Доказательство. В соответствии с условием леммы, подстрока $x = x[1..i]$ имеет минимальный период $p^* = i - \beta[i]$ и поэтому может быть записана в нормальной форме как $x = x[1..p^*]^{r^*} x[1..p']$ для целых $r^* = \lfloor i/p^* \rfloor$ и $p' = i \bmod p^*$. Для каждого $j \in p^*..i$ префикс

$$x[1..j] = x[1..p^*]^{\lfloor i/p^* \rfloor} x[1..j \bmod p^*]$$

является оболочкой правого расширения

$$xx[p' + 1..p^*]x[1..j \bmod p^*] = x[1..p^*]^{r^*+1} x[1..j \bmod p^*].$$

Поэтому позиция j жива по отношению к i . ■

В качестве примера к этой лемме рассмотрим строку $x[1..10] = abcaabcaab$ длиной $i = 10$ и периода $p^* = 4$. Тогда для каждого $j \in 4..10$ подстроки $x[1..j]$ являются оболочками правого расширения

$$xx[3..4]x[1..j \bmod 4] = (abca)^3 x[1..j \bmod 4].$$

Однако заметим, что достаточное условие леммы 13.1.4 не является необходимым, например $x[1..3] = aba$ является оболочкой строки $x[1..8] = ababaaba$ длиной $i = 8$ и периода $p^* = 5$, хотя $j = 3$ не принадлежит интервалу $5..8$. Следующая лемма дает полную характеристику живых позиций.

Лемма 13.1.5. По отношению к каждому $i \in 1..n$ позиция j жива тогда и только тогда, когда выполняется одно из следующих двух условий:

- а) $j \in i - \beta[i]..i$;
- б) $x[1..j]$ является оболочкой подстроки $x[1..j']$ для некоторого $j' \in i - \beta[i]..i$.

Доказательство. Чтобы доказать, что условия а) и б) необходимы, сначала предположим, что позиция j жива по отношению к позиции i . Принимая во внимание лемму 13.1.4, необходимо показать только то, что условие б) выполняется для $j < p^* = i - \beta[i]$. Поэтому предположим, что условие б) не выполняется: $x[1..j]$

не является оболочкой $x[1..j']$ ни при каком $j' \in p^*..i$. Но тогда, поскольку $j < p^*$, $x[1..j]$ не может покрывать ни одно из расширений x — приходим к противоречию. Итак, если условие *a*) не выполняется, тогда должно выполняться условие *b*) для того, чтобы позиция j могла быть живой по отношению к позиции i .

Достаточность условия *a*) уже установлена леммой 13.1.4. Чтобы доказать достаточность условия *b*) предположим, что для некоторого целого числа $j \leq p^* - 1$ подстрока $x[1..j]$ является оболочкой некоторой подстроки $x[1..j']$, $j' \in p^*..i$.

Сначала положим, что $j' < i - p^*$. Тогда, поскольку $x[1..i]$ имеет период p^* , если $x[1..j]$ покрывает $x[1..j']$, то она также должна покрывать $x[p^* + 1..p^* + j']$ и, следовательно, $x[1..p^* + j']$. Фактически, для каждой позиции $j' < i - p^*$ $x[1..j]$ покрывает $x[1..p^* + j']$, и поэтому без потери общности можем допустить, что $j' \geq i - p^*$.

Теперь можем показать, что правое расширение $y = xx[i - p^* + 1..j']$ длиной $p^* + j' \geq i$ имеет грань $x[1..j']$ длиной $j' \geq (p^* + j')/2$ и поэтому покрывается оболочкой $x[1..j']$. Тогда, по определению, j жива по отношению к i , что и требовалось доказать.

Оставим для упражнения 13.1.4 доказательство равенства

$$y = x[1..j']x[j' + 1..i]x[i - p^* + 1..j'] = x[1..p^*]x[1..j']. \quad (13.5)$$

■

Как будет видно из следующего подраздела, для вычисления массива оболочек очень важно охарактеризовать мертвые узлы по отношению к началу ступени.

Лемма 13.1.6. Пусть i — позиция, с которой начинается ступень в массиве граней β . Тогда позиция $j < i - \beta[i]$ будет мертвой по отношению к i только тогда, когда j не будет иметь узлов-сыновей в T_γ , которые живы по отношению к i .

Доказательство. Предположим, что позиция $j < i - \beta[i]$ является мертвой по отношению к i . Тогда $x[1..j]$ не может быть оболочкой ни одного правого расширения $x[1..i]$. Если j имеет узла-сына j' в T_γ , который жив по отношению к i , тогда $x[1..j']$ будет потенциальной оболочкой $x[1..k]$ для некоторого $k \geq i$. Но поскольку j' является сыном j , то тогда $x[1..j]$ покрывает $x[1..j']$ и поэтому, в соответствии с фактом F3, покрывает также $x[1..k]$, что является противоречием. Таким образом, необходимость условий теоремы доказана.

Чтобы доказать их достаточность, предположим, что $j < i - \beta[i]$ не имеет сынов, живых по отношению к i , но сам узел j жив по отношению к узлу i . Тогда, согласно лемме 13.1.5, $x[1..j]$ должна покрывать некоторую подстроку $x[1..k]$, $i - \beta[i] \leq k \leq i$. Поэтому j имеет сына k в T_γ , который, снова согласно лемме 13.1.5, должен быть живым по отношению к i , что является противоречием. Мы приходим к заключению, что j мертв по отношению к i .

■

Следствия леммы: если узел j жив по отношению к узлу i , тогда жив и его узел-родитель $\gamma[i]$ в T_γ .

В заключение докажем простое свойство функции $i - \beta[i]$, которое необходимо для вычисления массива оболочек.

Лемма 13.1.7. Функция $i - \beta[i]$ инвариантна для каждого i в одной и той же ступени и монотонно неубывающая по i ; в частности, для любой позиции $i > 1$, которая начинает новую ступень в β , выполняется неравенство

$$(i - 1) - \beta[i - 1] < i - \beta[i].$$

Доказательство является следствием того факта, что неравенство $\beta[i] \leq \beta[i - 1]$ справедливо тогда и только тогда, когда i начинает новую ступень. ■

Исходя из лемм 13.1.5 и 13.1.7, приходим к заключению, что совокупность мертвых узлов для всех значений i одной и той же ступени остается неизменной; таким образом, ее повторное вычисление (с использованием леммы 13.1.6) необходимо только тогда, когда начинается новая ступень.

Вычисление массива оболочек

Алгоритм ЛС, обрабатывая слева направо для каждого $i = 1, 2, \dots, n$ как строку x , так и массив граней β , выполняет две основные задачи:

- вычисляет $\gamma[i]$ и добавляет i в качестве нового сына узла $\gamma[i]$ в дерево оболочек T_γ — вспомним, что по определению $\gamma[i]$ является предком i в T_γ ;
- для каждой позиции i , которая является началом ступени в β , вычисляет узлы в T_γ , которые мертвы по отношению к i .

В дополнение к массивам $\beta[1..n]$ и $\gamma[1..n]$, алгоритм также использует следующие массивы.

- Массив $dead[0..n - 1]$: согласно лемме 13.1.6, $dead[i] = \text{TRUE}$ тогда и только тогда, когда для текущего значения $ij < i - \beta[i]$ и j не имеет сыновей в T_γ , которые живы по отношению к i (как можно видеть, корневой узел 0 в T_γ всегда живой);
- Массив $livechildren[0..n]$: $livechildren[i] = k$ тогда и только тогда, когда узел i имеет в дереве оболочек T_γ ровно k живых узлов-сыновей;
- Массив $largestlive[0..n]$: $largestlive[i] = j$, где j — наибольший предок i в дереве оболочек T_γ (разумеется, узел i всегда живой по отношению к самому себе).

Этот алгоритм начинается с помещения узла 0 в дерево оболочек T_γ и установки каждого элемента в массиве $dead$ в состояние FALSE, каждого элемента в $livechildren$ в состояние 0, а $largestlive[i]$ в состояние i для каждого $i \in 1..n$. Затем для текущей i -й позиции в строке x ($i = 1, 2, \dots, n$) в алгоритме ЛС выполняются следующие три шага.

- Шаг 1 обновляет массив *largestlive*.
 - Шаг 2 присоединяет *i* к его собственному узлу-родителю в T_γ .
 - Шаг 3 гарантирует обновление массива *dead* в начале каждой новой ступени.
- Эти шаги более подробно описываются в алгоритме 13.1.1.

Алгоритм 13.1.1 (Алгоритм ЛС)

```

▷ Вычисление массива оболочек  $\gamma$  для строки  $x[1..n]$ 
for  $i \leftarrow 1$  to  $n$  do
  ▷ Шаг 1: если  $\beta[i]$  мертвый в  $T_\gamma$ , то  $\beta[i]$  должен иметь
  ▷ того же наибольшего предка, что и его родитель
  if  $dead[\beta[i]] = \text{TRUE}$  then
     $largestlive[\beta[i]] \leftarrow largestlive[\gamma[\beta[i]]]$ 
  ▷ Шаг 2: Вычисление  $\gamma[i]$ 
  ▷ если  $\beta[i]$  живой, положить  $\gamma[i] \leftarrow \beta[i]$ ; иначе, поскольку каждая
  ▷ оболочка  $x[1..i]$  должна покрывать  $x[1..\beta[i]]$ , положить  $\gamma[i]$  равным
  ▷ самому большому живому предку  $\beta[i]$  (возможно, 0)
   $\gamma[i] \leftarrow largestlive[\beta[i]]$ 
   $livechildren[\gamma[i]] \leftarrow livechildren[\gamma[i]] + 1$ 
  ▷ Шаг 3: В  $T_\gamma$  определяются все узлы, которые становятся
  ▷ мертвыми в результате начала новой ступени
  if  $i > 1$  and  $\beta[i] \leq \beta[i - 1]$  then
     $c_1 \leftarrow i - \beta[i]; c_2 \leftarrow (i - 1) - \beta[i - 1]$       ▷ Лемма 13.1.7
    for  $j \leftarrow c_1$  downto  $c_2$  do                          ▷ Лемма 13.1.6
      if not  $dead[j]$  and  $livechildren[j] = 0$  then
         $dead[j] \leftarrow \text{TRUE}$ 
         $j' \leftarrow \gamma[j]$ 
         $livechildren[j'] \leftarrow livechildren[j'] - 1$ 
        ▷ Установить мертвыми все предки,
        ▷ у которых нет живых сынов
        while  $livechildren[j'] = 0$  do
           $dead[j'] \leftarrow \text{TRUE}$ 
           $j' \leftarrow \gamma[j']$ 
           $livechildren[j'] \leftarrow livechildren[j'] - 1$ 

```

Справедлива следующая теорема.

Теорема 13.1.8. Алгоритм 13.1.1 корректно вычисляет массив оболочек $\gamma[1..n]$ произвольной строки $x[1..n]$.

Доказательство. Рассмотрим каждый шаг алгоритма в отдельности, начиная с шага 3. Пусть $1 = i_1 < i_2 < \dots < i_k \leq n$ обозначают позиции (номера)

элементов массива β , с которых начинаются ступени. На шаге 3 для каждой позиции $j \in i_{h-1} \dots i_h - 1$ ($h = 2, 3, \dots, k - 1$) проверяется, будет ли эта позиция мертвой по отношению к i_{h+1} . В соответствии с леммой 13.1.6, узел-родитель $\gamma[j]$ любого элемента j , установленного мертвым, должен быть проверен: если $\gamma[j]$ ранее был живым, то теперь он может быть мертвым тогда и только тогда, когда у него нет живых сыновей. Таким образом, на шаге 3 должна выполняться рекурсивная проверка родителя *любой* позиции, установленной мертвой, до тех пор пока не будет найден предок, имеющий, по меньшей мере, одного живого сына (поэтому этот предок останется живым). Отметим, что поскольку позиция i_{h+1} по отношению к самой себе всегда живая, то, согласно лемме 13.1.6, должен существовать путь, ведущий от i_{h+1} до корня и содержащий только живые узлы. Таким образом, корень всегда живой, а выполнение шага 3 будет всегда заканчиваться на первом живом узле по пути от j до корня.

Отметим, что позиции $j \in i_{h-1} \dots i_h - 1$ просматриваются в обратном порядке $i_h - 1, i_h - 2, \dots, i_{h-1}$ из-за того, что, возможно, j или также и $\gamma[j] \in i_{h-1} \dots i_h - 1$ установлены мертвыми. Именно по этой самой причине, чтобы избежать избыточной обработки того же пути к корню, необходимо проверить в цикле **for** шага 3 условие **not dead[j]**: узел j может быть установлен мертвым из-за того, что он является родителем большего узла из интервала $i_{h-1} \dots i_h - 1$, который уже установлен мертвым.

Поскольку $livechildren[\gamma[j]]$ всегда уменьшается на единицу для каждого узла j , который устанавливается мертвым, приходим к заключению, что шаг 3 корректно справляется с задачей установки мертвыми узлов, находящихся в начале каждой новой ступени.

Шаг 2 заключается в непосредственном обновлении T_γ , базирующемся на текущей позиции i . Его корректность полностью зависит от корректного обновления $largestlive[\beta[i]]$ в шаге 1.

Оператор присваивания в шаге 1 будет выполнен только в том случае, если $\beta[i]$ был установлен мертвым во время выполнения одного из предыдущих циклов **for** в шаге 3. Как было отмечено в доказательстве шага 3, если некоторый узел j должен быть установлен мертвым относительно i_{h+1} , он должен быть живым по отношению к i_h , следовательно, согласно лемме 13.1.5, он живой относительно $i_{h+1} - 1$. Поэтому $x[1..j]$ может быть *потенциальной* оболочкой $x[1..i_{h+1}]$; таким образом, она может покрывать некоторое правое расширение $x[1..\beta[i_{h+1}]]$, которое является суффиксом подстроки $x[1..i_{h+1}]$. Следовательно, если какая-либо позиция j устанавливается мертвой, тогда j будет превышать $\beta[i_{h+1}]$, где i_{h+1} — начало ступени, которое используется на шаге 3.

Поскольку, согласно лемме 1.3.1, значения в массиве β могут увеличиваться не более чем на единицу от одной позиции к следующей, тогда каждая такая мертвая позиция j должна позже стать значением в массиве β , т.е. для того чтобы был выполнен шаг 1, необходимо положить $j = \beta[i']$ для некоторых $i' > i_{h+1}$.

В частности, значения j должны обрабатываться в шаге 1 в порядке возрастания их величин: т.е. в порядке их убывания в дереве оболочек T_γ . Это означает, что шаг 1 передаст правильные значения наибольшего живого предка от родителя к сыну. ■

Теперь рассмотрим время, необходимое для выполнения алгоритма ЛС. Каждая операция на шаге 3, кроме цикла **for**, выполняется за константное время. В таком случае для выполнения алгоритма ЛС необходимо время порядка $\Theta(n)$ плюс общее время, использованное циклом **for**. Для оценки этого общего времени заметим, что время, необходимое для *каждого однократного* выполнения цикла **for**, пропорционально выражению

$$\max\{(c_1 - c_2), \text{ количество узлов, установленных мертвыми}\}.$$

Поскольку сумма разностей $c_1 - c_2$ по всем ступеням не превышает $n - 1$ и каждый из n узлов хотя бы один раз может быть установлен мертвым, то из этого следует, что общее время для выполнения всех циклов **for** будет иметь порядок $O(2n)$. Следовательно, справедлива теорема.

Теорема 13.1.9. Алгоритм 13.1.1 требует для своего выполнения время порядка $\Theta(n)$ и объем памяти порядка $\Theta(n)$. ■

Алгоритм ЛС — это оптимальный *онлайновый* алгоритм в слабом смысле (см. определение в разделе 4.1): он вычисляет γ таким образом, что делает доступными все оболочки каждого префикса строки x , затрачивая на это асимптотически наименьшее возможное время и объем памяти. Отметим, что поскольку вычисление β — также онлайновый процесс, то алгоритм ЛС можно легко модифицировать для одновременного онлайнового вычисления обоих массивов: массива граней и массива оболочек.

Читателю было бы полезно посмотреть, каким образом алгоритм можно применить к примеру, данному на рис. 13.1. В частности, отметим, что в результате того факта, что $\beta[22] = 2$, установлены мертвыми следующие узлы: 5, 6, 9, 10, 11, 13–19.

Упражнения 13.1

1. Докажите факты F1–F4.
2. На основе фактов F1–F3 покажите, что $u = x[1..p]$ является оболочкой $x[1..i]$ тогда и только тогда, когда для некоторого целого $h \in 1..k - 1$ $p = \gamma^h[i]$, где k — наименьшее положительное целое, такое, что $\gamma^k[i] = 0$.
3. Строго говоря, массив оболочек наиболее близко соответствует не массиву граней, а, скорее, *массиву периодов* $\pi = \pi[1..n]$, в котором $\pi[i]$ равен

наибольшему периоду строки $x[1..i]$, являющейся кратной строкой, иначе $\pi[i]$ равен нулю. Без вычисления массива граней определите алгоритм, который вычисляет массив π за линейное время. Затем покажите, что этот массив периодов можно представить в виде *дерева периодов* со свойствами, аналогичными свойствам дерева граней и дерева оболочек.

4. Используя факты, что x имеет период p и что $j' \geq \max\{p, i - p\}$, докажите равенство (13.5).
5. Измените алгоритм ЛС таким образом, чтобы обеспечить онлайн-овое вычисление обоих массивов β и γ .

13.2 Все раппорты — алгоритм Франека–Смита–Танга

В этом разделе мы рассматриваем проблему точных раппортов в ее самой общей форме: говоря языком раздела 2.3, мы вычислим все NE (непродолжаемые) полные раппорты в произвольной строковой последовательности $x[1..n]$ за время порядка $O(n \log n)$ и представим эти раппорты в виде массива или дерева, для хранения которых требуется память порядка $\Theta(n)$. Сначала опишем алгоритм Франека–Смита–Танга (Franek–Smyth–Tang [93], сокращенно — алгоритм ФСТ), использующий деревья суффиксов (раздел 5.2), затем покажем, каким образом он может быть реализован с помощью массивов суффиксов (подраздел 5.3.2) при значительной экономии пространства памяти. Алгоритм вычисления всех раппортов основан на простом факте.

Лемма 13.2.1. Обозначим через \hat{x} строку $x[n]x[n-1] \dots x[1]$, обратную к заданной строке x . Тогда раппорт $M_{x,u}$ будет раппортом типа LE (продолжаемым влево) тогда и только тогда, когда раппорт $M_{\hat{x},\hat{u}}$ будет раппортом типа RE (продолжаемым вправо).

Доказательство предложено дать в упражнении 13.2.1. ■

Этот результат предлагает прямой подход к вычислению всех NE-раппортов в строке x .

- **Шаг 1.** Вычисление всех NRE-раппортов (не RE-раппортов) в обеих строках x и \hat{x} .
- **Шаг 2.** Сравнение (каким-либо образом) NRE-раппортов строки x с аналогичными раппортами обратной строки \hat{x} — согласно лемме 13.2.1, совпадающие раппорты будут NE-раппортами.

Этот, простой с виду подход обладает труднопреодолимыми недостатками, связанными с тем, что раппортов типов NRE и NLE может быть очень много,

и с тем, что не совсем понятно, каким образом их можно эффективно сравнивать друг с другом. В частности, этот подход может быть очень дорогостоящим и по времени вычисления, и по объему памяти. И тем не менее, использование, в сущности, этого подхода представляется привлекательным по двум причинам.

- Как отмечено в подразделе 12.1.1, NRE-раппорты естественно представимы в виде дерева суффиксов (и, как мы увидим позже, также в виде массива суффиксов). Таким образом, шаг 1 может быть реализован путем построения деревьев (массивов) суффиксов для строки x и для обратной строки \hat{x} .

Для примера рассмотрим дерево суффиксов, построенное для нашего обычного примера строки Фибоначчи f_6 , показанного на рис. 13.4 в упрощенном виде, немного отличающемся от рис. 12.1. Листья дерева (квадратики) помечены начальными позициями $i \in 1..13$ суффиксов $f_6[i..13]$, в то время как каждый внутренний узел (кружок) помечен числом, равным длине наибольшего общего префикса (lcp) совокупности суффиксов, идентифицированных расположенными ниже листьями. Таким образом, значения lcp определяют длины полных NRE-раппортов в строке. Поэтому такое дерево суффиксов назовем *NRE-деревом*, внутренние (круглые) и конечные (квадратные) узлы которого будем называть соответственно *lcp-узлами* и *позиционными узлами*.

В нашем примере подстроки $u = x[9..13] = x[6..10] = x[1..5]$ определяются поддеревом с корнем в lcp-узле 5 и составляют полный NRE-раппорт $M_{x,u}^* = (5; 9, 6, 1)$.

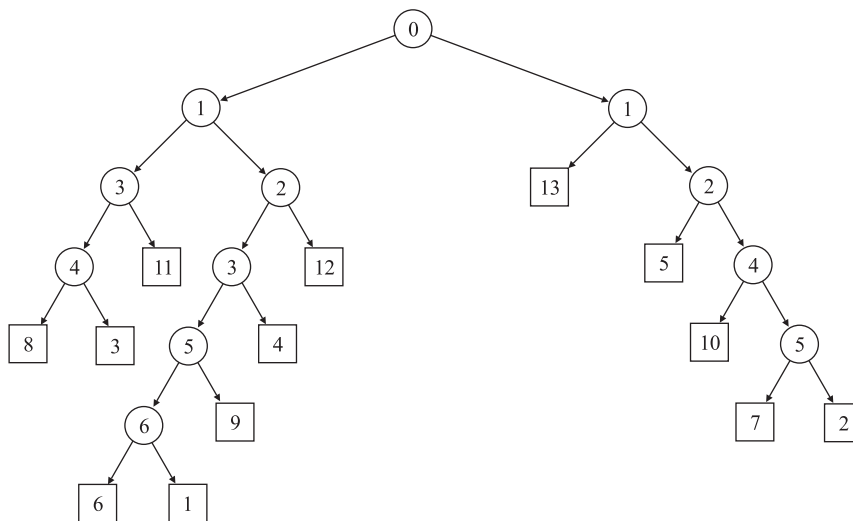


Рис. 13.4. NRE-дерево для строки f_6

- Как показано ниже, шаг 2 можно реализовать с помощью линейной по времени процедуры, которая определяет совпадения между NRE-деревьями (массивами) строк x и \hat{x} для создания NE-дерева (NE-массива), которое описывает все NE-раппорты в строке x .

Основываясь на этих замечаниях, в подразделе 13.2.1 покажем, как, имея NRE-деревья для x и \hat{x} , можно вычислить NE-дерево, затем в подразделе 13.2.2 покажем, как можно вычислить NE-массив на основе известных массивов суффиксов строк x и \hat{x} .

13.2.1 Вычисление NE-дерева

Предполагаем, что деревья суффиксов T_x и $T_{\hat{x}}$ для строк x и \hat{x} уже вычислены, определив таким образом полные NRE-раппорты этих строк. Из леммы 13.2.1 следует, что NRE-раппорты строки \hat{x} определяют NLE-раппорты строки x . Для краткости образующую u раппорта $M_{x,u}$ типа NRE (соответственно, NLE и NE) назовем повторяющейся NRE-подстрокой (соответственно, NLE- и NE-подстрокой), при этом просим читателя иметь в виду, что понятия NRE, NLE, NE связаны именно с раппортом, который является совокупностью повторяющихся подстрок.

Предположим, что некоторый lcr-узел \hat{p} ($\hat{p} > 0$) в дереве $T_{\hat{x}}$ имеет в качестве узла-сына позиционный узел \hat{i} . Тогда в \hat{x} существует повторяющаяся NRE-подстрока $\hat{u} = \hat{x}[\hat{i}.. \hat{i} + \hat{p} - 1]$, обратная повторяющейся NLE-подстроке $u = x[n - (\hat{i} + \hat{p} - 1) + 1.. n - \hat{i} + 1]$ строки x . Таким образом, присваивание

$$i \leftarrow n - (\hat{i} + \hat{p} - 2) \tag{13.6}$$

определяет одну начальную позицию i в x повторяющейся NLE-подстроки u длиной \hat{p} . Заметим, что фактически i является также начальной позицией одной повторяющейся NLE-подстроки из набора таких подстрок, имеющих длину $j \in 1.. \hat{p}$. Таким образом, если в T_x мы находим для некоторых i наибольшее значение \hat{p} , такое, что i будет одной из совокупности подстрок длиной \hat{p} , которые являются и NRE-, и NLE-подстроками, тогда каждый родитель \hat{p} в T_x также определяет совокупности подстрок, которые являются и NRE-, и NLE-подстроками.

Назовем подстроку u строки x *максимальной повторяющейся NE-подстрокой* строки x , если

- u встречается в строке x по меньшей мере дважды;
- u не является собственной подстрокой любой повторяющейся подстроки строки x .

Следующий результат говорит о том, что максимальная повторяющаяся NE-подстрока должна определяться как в T_x , так и в $T_{\hat{x}}$.

Лемма 13.2.2. Если $u = x[i..i + p - 1]$ является максимальной повторяющейся NE-подстрокой строки x , тогда позиционный узел i будет узлом-сыном lcr-узла в T_x , а узел $n - (i - p - 2)$ — сыном p в $T_{\hat{x}}$.

Доказательство предложено дать в упражнении 13.2.4. ■

Вспомним сделанное выше замечание, что каждый lcr-узел в T_x , который является предком повторяющейся NE-подстроки, сам должен быть корнем поддерева T_x , позиционные узлы которого являются повторяющимися подстроками полного NE-раппорта. Поэтому последняя лемма обеспечивает нас простой стратегией для определения всех повторяющихся NE-подстрок: требуется найти только максимальные подстроки (с наибольшим значением p), затем определить местонахождение всех их предков в дереве T_x . Алгоритм 13.2.1 реализует эту схему.

Алгоритм 13.2.1 (Вычисление NE-дерева)

- ▷ Заданы деревья суффиксов T_x и $T_{\hat{x}}$, вычисляется NE-дерево для $x[1..n]$
- 1: Обход дерева T_x для создания таблицы POINTER, где для каждой позиции i в строке x POINTER[i] является указателем на соответствующий узел в T_x
- 2: **for** каждый lcr-узел в T_x **do**
 NE[p] ← FALSE
- 3: **for** каждая пара родитель–сын (\hat{i}, \hat{p}) в $T_{\hat{x}}$ **do**
 ▷ Здесь используется POINTER[i]:
 if $i = n - (\hat{i} + \hat{p} - 2)$ является сыном lcr-узла в T_x **then**
 while not NE[\hat{p}] **do**
 NE[\hat{p}] ← TRUE
 if $\hat{p} \neq 0$ **then**
 $\hat{p} \leftarrow$ родитель \hat{p} в T_x
- 4: Обход T_x для удаления каждого поддерева с корнем в lcr-узле p , для которого NE[p] = FALSE

Алгоритм разбит на четыре шага. Шаг 1 — это обход дерева T_x , во время которого создается таблица, предоставляющая доступ к каждому позиционному узлу i в дереве T_x за константное время. На шаге 2 булева переменная, соответствующая каждому lcr-узлу p в T_x , устанавливается в исходное состояние FALSE, указывая на то, что ни один из конечных узлов в поддереве с корнем в p не был на данный момент определен как NE-узел. На шаге 3 обрабатывается каждая пара родитель–сын (\hat{i}, \hat{p}) дерева $T_{\hat{x}}$ с целью определить, существует ли в T_x эквивалентная пара родитель–сын $(\hat{p}, n - (\hat{i} + \hat{p} - 2))$: если это так, тогда lcr-узел \hat{p} и все его предки в T_x должны быть типа NE. Соответственно, пока не будет найден предок, который уже NE типа, \hat{p} и его предки в T_x идентифицируются как NE типа. На заключительном шаге выполняется обход дерева T_x с целью удаления всех

поддеревьев с корнем в любом узле p , для которого $NE[p] = FALSE$: оставшееся дерево T_x^{NE} является NE-деревом строки x .

На шагах 1, 2 и 4 алгоритма 13.2.1 совершаются обходы NRE-дерева, при этом для обработки каждого узла необходимо время порядка $O(1)$, следовательно, для обработки всех узлов необходимо время порядка $O(n)$. На шаге 3 также совершается обход дерева, но с более сложной обработкой узлов: поскольку оператор **if** использует массив POINTER, он при каждом выполнении также требует только $O(1)$ времени, следовательно, полное время выполнения этого оператора будет также порядка $O(n)$. В цикле **while** шага 3 переменная NE для не более n lcr-узлов устанавливается в состояние TRUE, а также проверяется текущее состояние не более n lcr-узлов, для которых переменная NE уже установлена в TRUE. Таким образом, полное время, затрачиваемое на выполнение цикла **while**, также имеет порядок $O(n)$. Следовательно, справедлива теорема.

Теорема 13.2.3. Если известны деревья суффиксов T_x и $T_{\hat{x}}$ для строки $x[1..n]$, алгоритм 13.2.1 правильно вычисляет NE-дерево за время порядка $O(n)$ с использованием памяти порядка $\Theta(n)$. ■

На рис. 13.5 изображено дерево суффиксов $T_{\hat{f}_6}$, измененное таким образом, что позиционные узлы \hat{i} заменены на узлы $i = n - (\hat{i} + \hat{p} - 2)$. Отметим, что только позиции 1, 6 и 9 приводят к NLE-раппортам. Таким образом, NE-дерево сводится до поддерева дерева T_{f_6} , которое соответствует этим позициям — это пути от корня дерева T_{f_6} до lcr-узла 6, который соответствует подстроке $abaaba$. Как видно из рис. 13.4, все полные раппорты строки f_6 являются вхождениями подстрок $a, ab, aba, abaab$ и $abaaba$. В общем случае вхождения NE-раппортов, всех из них или выбранных в соответствии с какими-либо критериями, можно локализовать и вывести путем обхода NE-дерева за время, пропорциональное их количеству. Такие задачи обсуждаются в работе [41]. Другой, совершенно отличный подход

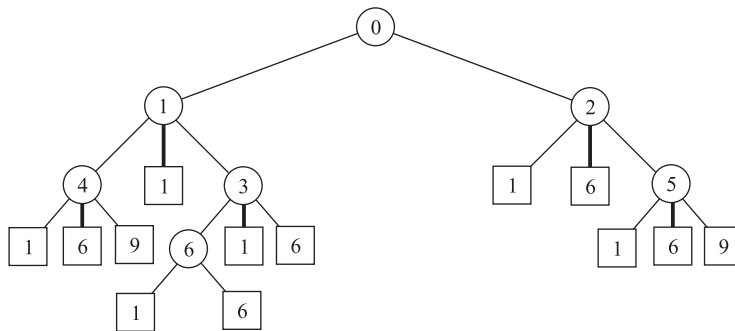


Рис. 13.5. Дерево $T_{\hat{f}_6}$ для \hat{f}_6 с замененными узлами

к вычислению непродолжаемых раппортов (хотя тоже с использованием деревьев суффиксов) описан в [108].

13.2.2 Вычисление NE-массива

В подразделе 5.3.2 мы определили массив суффиксов $\sigma_x = \sigma[1..n]$ для заданной строки $x[1..n]$ как перестановку чисел $1..n$ с таким свойством, что $x[\sigma_x[i]..n] < x[\sigma_x[j]..n]$ для всех $1 \leq i < j \leq n$. Таким образом, σ_x определяет начальные позиции суффиксов строки x в возрастающем лексикографическом порядке. Например, если

$$f_6 = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ a & b & a & a & b & a & b & a & a & b & a & a & b \end{matrix} \quad (13.7)$$

тогда

$$\sigma_{f_6} = 11 \ 8 \ 3 \ 12 \ 9 \ 6 \ 1 \ 4 \ 13 \ 10 \ 7 \ 2 \ 5. \quad (13.8)$$

В подразделе 5.3.2 также показано, как путем присваивания lcp-значений диапазонам позиций в σ_x можно определить вспомогательный массив π_x . Здесь мы определим более простой массив $\lambda_x = \lambda[1..n]$, где $\lambda[1] = 0$, а для каждого $i \in 2..n$

$$\lambda[i] = \text{lcp}(x[\sigma[i-1]..n], x[\sigma[i]..n]).$$

Будем использовать сокращенную форму записи:

$$\lambda[i] = \text{lcp}(\sigma[i-1], \sigma[i]). \quad (13.9)$$

Например, в нашем примере

$$\lambda_{f_6} = 0341256301452, \quad (13.10)$$

а это свидетельствует о том, что $\text{lcp}(11, 8) = 3$, $\text{lcp}(8, 3) = 4$ и т.д. Так же, как и в случае массива π_x , массив λ_x можно вычислить в качестве промежуточного результата при вычислении массива суффиксов. Напомним (подраздел 5.3.2), что массивы суффиксов можно вычислить независимо от дерева суффиксов [172].

Мы представим здесь схему одного алгоритма, в котором не используются деревья суффиксов: все вычисление выполняется на основе массивов суффиксов (как σ_x , так и λ_x), и его конечным результатом является NE-массив, а не NE-дерево. Поскольку для хранения массива суффиксов строки $x[1..n]$ требуется только $2n$ компьютерных слов, то этот алгоритм приводит к существенной экономии памяти.

Пусть $p_j = \lambda[j]$ для каждого $j \in 1..n$. Тогда из выражения (13.9) следует, что для любого $j \in 2..n$

if $p_j > p_{j-1}$ **then**
 p_j — потомок p_{j-1} в дереве суффиксов
else if $p_j < p_{j-1}$ **then**
 p_j — предок p_{j-1} в дереве суффиксов
else
 p_j и p_{j-1} определяют один и тот же узел в дереве суффиксов

Поскольку массив суффиксов содержит такую же информацию о длинах наибольших общих префиксов, которая предоставляется деревом суффиксов, то эти отношения подразумевают, что массив суффиксов должен быть разложен на подмассивы, каждый из которых определяет узлы, лежащие на одном пути от корня до конечного узла в дереве суффиксов. Например, значения Icr в массиве (13.10) могут быть разделены на три подмассива

$$0, 3, 4, \underline{1}/\underline{1}, 2, 5, 6, 3, \underline{0}/\underline{0}, 1, 4, 5, 2,$$

соответствующие трем основным путям в дереве суффиксов, показанном на рис. 13.4. Здесь повторяющиеся значения $\underline{0}$ и $\underline{1}$ определяют корни поддеревьев, в которых начинаются новые пути. Поэтому назовем такие узлы *узлами ветвления*, поскольку они отмечают точки, в которых новый путь отклоняется от предыдущего пути.

Эти понятия можно ввести более точно следующим образом.

Определение 13.2.4. *I-серией* $I_{j,h}$ в λ_x называется последовательность из $h \geq 2$ Icr -значений $p_j, p_{j+1}, \dots, p_{j+h-1}$ (которые не все равны), таких, что

- а) или $j = 1$ или $p_{j-1} > p_j$;
- б) $p_j \leq p_{j+1} \leq \dots \leq p_{j+h-1}$;
- в) или $j + h - 1 = n$ или $p_{j+h-1} > p_{j+h}$. ■

Таким образом, *I-серия* является неубывающей последовательностью Icr -значений в λ_x . Аналогично *D-серией* $D_{j,h}$ назовем невозрастающую последовательность Icr -значений. Теперь, используя эти определения, можно охарактеризовать пути в дереве суффиксов.

Лемма 13.2.5. В λ_x последовательности $I_{j,h}D_{j+h-1,h'}$ и $I_{j,h}$, где $j + h - 1 = n$, определяют узлы, которые в дереве суффиксов лежат на одном пути.

Доказательство. Поскольку $p_1 = 0$, то из определения 13.2.4 следует, что или в λ_x нет серий ($p_j = 0$ для всех $j \in 1..n$), или последовательность серий начинается с *I-серии* $I_{j,h}$. При этом определение 13.2.4 гарантирует, что за каждой *I-серией* следует *D-серия*, иначе достигается конец массива; аналогично за каждой *D-серией* следует *I-серия* либо достигается конец массива.

Теперь рассмотрим I-серии $I_{j,h}$. Поскольку $p_{j+(t-1)} \leq p_{j+t}$ для каждого $t \in 1..h-1$, то из этого следует, что $I_{j,h}$ определяет последовательность неубывающих p_j , которые лежат на одном пути в дереве суффиксов. Если $i+h=n$, тогда $I_{j,h}$ — конечный путь, определенный в λ_x . Если же нет, тогда за $I_{j,h}$ следует D-серия $D_{j+h-1,h'}$, где $p_{(j+h-1)+(t+1)} \leq p_{(j+h-1)+t}$ для каждого $t \in 1..h'-1$. Эти Icr-значения определяют последовательность предков p_{j+h-1} , лежащих на одном пути от корня, уже определенного с помощью $I_{j,h}$. Последующая I-серия, если она существует, естественно, определяет потомки узла ветвления $p_{j+h+h'-2}$, который лежит на пути, отличном от пути, определенного с помощью $I_{j,h}$. ■

Очевидно, что в λ_x ID-серии, по существу, разбивают внутренние (Icr) узлы массива суффиксов на отдельные пути. При этом смежные пары серий определяют в дереве два пути, которые имеют только один общий узел (узел ветвления). Заметим также, что для каждой ID-серии минимальный узел на соответствующем пути должен быть или первым узлом I-серии, или последним узлом D-серии.

Поскольку узел ветвления является как последним узлом D-серии, так и первым узлом I-серии, то значения пар смежных узлов в λ_x определяют способ, каким связаны два пути в дереве суффиксов. Если ввести в λ_x первое Icr-значение ($p_1 = 0$) в качестве узла ветвления и обозначить последовательные узлы ветвления соответственно как b_1 и b_2 , то нетрудно заметить, что есть только три возможных отношения между соответствующими путями, которые показаны на рис. 13.6. Отметим, что в каждом из этих трех случаев будущие пути могут быть добавлены на текущем пути только в качестве потомков b_2 или потомков предков b_2 . Для пути с корнем в b_2 , который был определен предыдущим путем, содержащим b_1 , никакого дальнейшего изменения не может быть сделано.

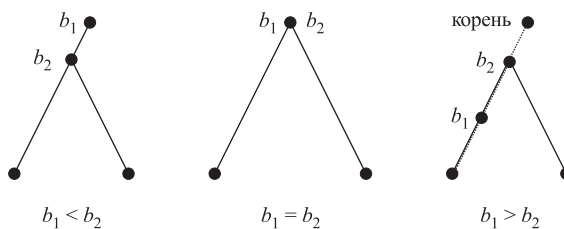


Рис. 13.6. Узлы ветвления b_1 и b_2 смежных путей

Наконец, заметим, что регулярность структуры массива суффиксов дает возможность определить в соответствующем дереве суффиксов T_x родителя каждого позиционного узла $\sigma[j]$, $j = 1, 2, \dots, n$. Сделаем два важных утверждения, доказательство которых оставим для упражнения 13.2.12.

O1. Родитель $\sigma[j]$ в T_x равен $\max\{\lambda[j], \lambda[j+1]\}$, где принимаем, что $\lambda[n+1] = 0$.

O2. (По Эндрю Фрэнсису (Andrew Francis)) Если $\lambda[j]$ — узел ветвления и $j < n$, то $\sigma[j]$ имеет родителя $\lambda[j+1]$ в T_x ; иначе родителем $\sigma[j]$ будет $\lambda[j]$.

Эти утверждения предоставляют альтернативные способы определения узлов ветвления в λ_x .

Теперь опишем метод вычисления NE-раппортов, который основан только на массиве суффиксов. Как и ранее, для строки x обозначим массив позиций и lcr-массив соответственно как σ_x и λ_x ; а для обратной строки \hat{x} обозначим соответствующие массивы как $\hat{\sigma}_x$ и $\hat{\lambda}_x$. Введем также массив LOC, который для каждого позиционного узла \hat{i} в $T_{\hat{x}}$ покажет его местоположение в массиве $\hat{\sigma}_x$; это даст возможность найти любую позицию в $\hat{\sigma}_x$ за фиксированное время. Бинарный вектор NE[1..n] определяет для каждого $j \in 1..n$ состояние lcr-узла $\lambda[j]$ (TRUE или FALSE), а также его тип (NLE или нет). Наконец, чтобы определить начало и конец путей ID в λ_x , используем паттерны серий ID (лемма 13.2.5) для вычисления массива BRANCH, который, кроме узлов ветвления, определенных выше, включает также первые и последние позиции в λ_x . В результате получим более сложную схему вычислений, чем аналогичная схема на основе деревьев. Однако алгоритм 13.2.2, построенный по этой схеме, не требует больших временных затрат, чем алгоритм 13.2.1.

Алгоритм 13.2.2 (Вычисление NE-массива)

▷ Известны массивы суффиксов (NRE-массивы) для строк $x = x[1..n]$ и \hat{x} ,

▷ вычисляется NE-массив для строки x

1: **for** $j \leftarrow 1$ **to** n **do**

LOC[$\hat{\sigma}[j]$] $\leftarrow j$; NE[j] \leftarrow FALSE

2: вычисление BRANCH[1..b*], где для $b \in 1..b^*$ BRANCH[b]

является позицией b -го узла ветвления в λ_x

(позиции $j = 1$ и $j = n$ рассматриваются как узлы ветвления)

3: $b \leftarrow 1$; $p_0 \leftarrow 0$

for $j \leftarrow 1$ **to** n **do**

определение j' , такого, что $(p, i) = (\lambda[j'], \sigma[j])$ является парой родитель–сын в дереве суффиксов строки x

if $p > 0$ **then**

if $p < p_0$ **then**

NE[j'] \leftarrow TRUE; $p_0 \leftarrow p$

else

$\hat{i} \leftarrow n - (i + p - 2)$; $\hat{j} \leftarrow$ LOC[\hat{i}]

определение \hat{j}' , такого, что $(\hat{p}, \hat{i}) = (\hat{\lambda}[\hat{j}'], \hat{\sigma}[\hat{j}])$

является парой родитель–сын в дереве

суффиксов строки \hat{x}

```

if  $p = \hat{p}$  then
     $NE[j'] \leftarrow \text{TRUE}; p_0 \leftarrow p$ 
if  $j = \text{BRANCH}[b + 1]$  then
    определение  $j^* \in \text{BRANCH}[b].. \text{BRANCH}[b + 1]$ ,
    такого, что  $p^* = \lambda[j^*]$  равен максимальной  $\text{lcp}$ ,
    для которой  $NE[j^*] = \text{TRUE}$ 
    for каждого  $j^* \in \text{BRANCH}[b].. \text{BRANCH}[b + 1]$  do
        if  $\lambda[j'] \leq p^*$  then
             $NE[j'] \leftarrow \text{TRUE}$ 
     $b \leftarrow b + 1$ 

```

4: обработка пути, определяемого с помощью BRANCH, в обратном порядке, при установке типа NE (как в шаге 3) для предков p^*

Шаг 1 этого алгоритма — простой цикл **for**, который инициализирует массив LOC и бинарный вектор NE. Шаг 2 на основе леммы 13.2.5 вычисляет массив BRANCH и выполняет линейный просмотр массива $\lambda[1..n]$. На шаге 3 выполняется другой цикл **for**, который сначала на основе утверждения O1 определяет все совпадения пар (p, i) в массивах NRE строк x и \hat{x} , затем для каждого пути выполняет проверку, имеет ли каждый предок совпавших пар (p, i) в пределах этого пути тип NE. На выполнение каждого из двух этапов шага 3 необходимо время порядка $\Theta(n)$. На шаге 4 повторяется обработка второго этапа шага 3, но в обратном порядке — это гарантирует, что все предки NE-элементов установлены должным образом и тоже имеют тип NE. В конце алгоритма NE-массив определяется всеми теми позициями j в λ_x , для которых $NE[j] = \text{TRUE}$. Таким образом, имеет место следующая теорема.

Теорема 13.2.6. Для заданных массивов суффиксов (NRE-массивов), соответствующих строкам $x[1..n]$ и \hat{x} , алгоритм 13.2.2 корректно вычисляет NE-массив за время порядка $\Theta(n)$ с использованием памяти объемом порядка $\Theta(n)$. ■

Отметим, что, при необходимости, для вычисления NE-дерева за линейное время можно использовать соответствующий NE-массив.

Поскольку NE-дерево (или соответствующий ему NE-массив) определяет все NE-раппорты в строке x , то его можно использовать для определения последовательных NE-раппортов (повторений) и NE-подстрок, имеющих оболочки (задача 2.17 из раздела 2.3). Напомним, что листья поддерева с корнем в любом lcp -узле p в NE-дереве являются точными NE-раппортами, соответствующими p . Если эти узлы были отсортированы в последовательность $I = \langle i_1, i_2, \dots, i_r \rangle$, где $i_1 < i_2 < \dots < i_r$, тогда за единственный просмотр последовательности I можно легко определить как последовательные раппорты максимальной длины, так и подстроки максимальной длины, имеющие оболочки и соответствующие узлу p .

В первом случае лакуна между смежными входами в I должна быть в точности равна p , во втором случае не должна превышать p . Для каждого узла p этот процесс можно выполнить за время порядка $O(n)$ с использованием памяти такого же порядка.

Однако более сложной является проблема представления в строке x всех NE-подстрок, имеющих оболочки. Существует три алгоритма, имеющих отношение к проблеме вычисления всех таких подстрок типа NRE (с использованием дерева типа NRE, а не NE). Первый алгоритм, представленный в работе [15], выполняет сложные операции на NRE-дереве (дереве суффиксов) и требует времени порядка $O(n \log n)$, в то время как второй алгоритм [42] применяет еще более сложные методы по отношению к NRE-дереву для уменьшения затрат времени до величины порядка $O(n \log n)$. Третий алгоритм [124], непосредственно обрабатывающий последовательности, определенные декомпозицией Крочемора, также выполняется за время порядка $O(n \log n)$.

Научный мир ожидает “простого” алгоритма с временем выполнения порядка $O(n \log n)$, который вычислял бы NE-подстроки, имеющие оболочки, строки x предпочтительно со скромными требованиями по отношению к необходимому объему памяти.

Упражнения 13.2

1. Докажите лемму 13.2.1.
2. Определение дважды связанного списка, использованное в алгоритме ФСТ, неизбежно приводит к определению строки, данному в разделе 1.1. Если дважды связанный список действительно является строкой, тогда каким алфавитом он определяется?
3. Покажите, что ни один узел в NRE-дереве, возможно за исключением корня, не может иметь только единственного узла-сына. Определите условия, при которых корень действительно имеет единственный узел-сын.
4. Докажите лемму 13.2.2.
5. Вычислите NRE-массив строки \widehat{f}_6 , затем используйте его для проверки корректности NRE-дерева на рис. 13.5. Далее используйте NRE-дерева, показанные на рис. 13.4 и 13.5, для вычисления NE-дерева для строки f_6 в соответствии с алгоритмом 13.2.1.
6. Строка $g = abaababaababa$ отличается от строки f_6 только своими последними двумя позициями. Вычислите NRE-массивы строк g и \widehat{g} , затем NRE-дерева, а потом NE-дерево для строки g . Подумайте над поразительными различиями между этими структурами и аналогичными структурами из предыдущего упражнения.

7. Заметим, что для каждой позиции i , вычисленной для узла $n - (\hat{i} + \hat{p} - 2)$ дерева \widehat{T}_x , имеется только один возможный узел-родитель в дереве T_x . Объясните, каким образом использование поразрядной сортировки может улучшить алгоритм 13.2.1.
8. Запишите алгоритм, который на основе NRE-массива длиной n вычислял бы соответствующую строку $x[1..n]$, определенную на индексированном алфавите $1, 2, \dots, \alpha$.
9. Дайте точное определение D-серии в стиле определения 13.2.4.
10. Опишите алгоритм, который вычисляет NE-дерево строки x на основе NE-массива.
11. Опишите алгоритм, который вычисляет как последовательные NE-раппорты, так и подстроки типа NE, имеющие оболочки, соответствующие произвольному lcr-узлу p в NE-дереве.
12. Докажите утверждения O1 и O2.

13.3 k -приближенные раппорты — алгоритм Шмидта

В этом разделе описана схема алгоритма, который в общем случае эффективно вычисляет аппроксимирующие (приближенные) раппорты. Этот алгоритм особенно интересен специалистам в области молекулярной биологии, где в вычислении расстояния между строками используется общая матрица весов (раздел 2.2). Алгоритм, взятый из работы Шмидта (Schmidt) [205], является сложным алгоритмом с большим количеством довольно беспорядочных специальных терминов. В этой книге мы старались избегать описания таких алгоритмов. Однако здесь мы рассмотрим алгоритм Шмидта в *краткой* форме, поскольку он предоставляет действительно эффективное, возможно, даже асимптотически оптимальное, решение особо важной проблемы обработки строковых последовательностей. Кроме того, алгоритм Шмидта расширяет и объединяет несколько рассмотренных ранее в этой книге алгоритмических идей таким способом, который обеспечивает новое понимание их значения и эффективности.

- Алгоритм Шмидта основан на обычном подходе динамического программирования к массиву стоимостей c (раздел 9.1), но теперь с иной интерпретацией, применимой к повторяющимся подстрокам. Часть этой новой интерпретации включает вычисление максимальных “баллов”, а не минимального расстояния; таким образом, массив стоимостей c , в котором расстояние минимизировано, становится *массивом полезности* (benefit array) b , в котором эти баллы максимизированы.

- Понятие графа зависимостей, введенное в разделе 9.5 в качестве основы алгоритма Укконена–Майерса для вычисления расстояния между строками, вновь появляется здесь в виде *сеточного графа* (grid graph), “пути с наивысшим баллом” которого являются основным предметом рассмотрения алгоритма Шмидта.
- Понятие непродолжаемости, введенное в разделе 2.3 и затем примененное в разделе 12.2 к точным кратным строкам, а в разделе 13.2 к точным раппортам, для приближенных раппортов переформулировано в понятие *локальной оптимальности*. (Получается, что первоначальное определение локальной оптимальности [84] фактически предшествует определению непродолжаемости [168].)
- Рекурсивный метод декомпозиции, используемый в алгоритме Мейна–Лоренца (подраздел 12.1.2) для вычисления квадратов на основе их средних точек, расширен для вычисления приближенных последовательных квадратов в сеточных графах.

В дополнение к расширениям этих знакомых понятий, в алгоритме Шмидта представлены также некоторые новые интересные идеи, и мы просто не можем избежать их обсуждения!

По существу, этот алгоритм не вычисляет раппорты в принятом нами смысле; скорее, он вычисляет набор всех “локально оптимальных” (непродолжаемых) приближенных непокрываемых квадратов, которые удовлетворяют данному значению “точности” k . Как мы вскоре увидим, покрываемые квадраты исключаются, потому что их включение может привести к потере важных последовательных квадратов. Вычисление квадратов, а не раппортов с произвольным показателем степени связано с замечанием, сделанным в конце раздела 2.3, о нетранзитивности приближенного совпадения, следовательно, о присущей приближенным раппортам неопределенности. Например, в подстроке $u_1u_2u_3$ может быть, что $d(u_1, u_2) \leq k$, $d(u_2, u_3) \leq k$. Вследствие этого подстроки u_1u_2 и u_2u_3 являются k -приближенными квадратами. Однако если $d(u_1, u_3) > k$, тогда u_1u_3 таковой не будет. Таким образом, понятие k -приближенных кратных строк или раппортов непросто определить удобным способом: в нестрогом смысле строки вида $u_1u_2u_3$ могут быть k -приближенными кратными строками, хотя u_1u_3 не является даже k -приближенным квадратом. В то время в более строгом смысле обозначение $u_1u_2 \dots u_m$ как k -приближенной кратной строки может потребовать C_m^2 проверок возможных k -приближенных квадратов (как последовательных, так и расщепленных). Как отмечалось в разделе 2.3, образующие u_1 и u_2 k -приближенного квадрата u_1u_2 необязательно должны быть одинаковой длины.

Следовательно, для вычисления приближенных кратных строк, кажется, нет никакой надежды на использование какого-либо аналога тройки Крочемора (i, p, r) для кодирования точных кратных строк (подраздел 12.1.1), который уменьшает ко-

личество выводов до величины порядка $O(n \log n)$, и еще меньше надежды для использования аналога четверок (i, p, r, t) для кодирования серий (раздел 12.2), который уменьшает размер вывода до $O(n)$. Поэтому очевидно, что вычисление приближенных кратных строк или раппортов потребует вывода квадратов, следовательно, в наихудшем случае необходимо $\Omega(n^2)$ выводов (как показано в разделе 2.3 для строки $x = a^n$). Это наблюдение подтверждается теоретическим результатом работы [120], в которой проблема расстояний между строками решена с помощью произвольной матрицы баллов за время порядка $\Omega(n^2)$.

Как мы увидим в следующем разделе, по существу, те же трудности возникают при определении и вычислении приближенных образующих строки x .

Восстановленная симметрия

Напомним, что массив стоимостей c (или так называемая матрица динамического программирования) был сначала введен в разделе 9.1, где рассматривалась проблема вычисления расстояния между строками $x_1[1..n_1]$ и $x_2[1..n_2]$. Эти вычисления выполнялись путем нахождения для всех $i \in 0..n_1$ и $j \in 0..n_2$ расстояний между подстроками $x_1[1..i]$ и $x_2[1..j]$; при этом начальные значения нулевых строки и столбца массива c вычислялись симметрично, определяя расстояния от пустой строки ε до $x_2[j]$ и от $x_1[i]$ до ε . На основе начальных значений другие элементы массива c вычисляются методом динамического программирования, как определено в лемме 9.1.1. Таким образом, на каждом шаге рассматривался *префикс*, и поэтому всегда выполнялось равенство

$$c[i, j] = d(x_1[1..i], x_2[1..j]), \quad (13.11)$$

в котором x_1 и x_2 снова играют симметричные роли. Для этой задачи конечный результат всегда равен $c[n_1, n_2]$.

Однако эта “приятная” симметрия была потеряна в главе 10 при использовании массива стоимостей c для нахождения k -приближенных паттернов $p[1..m]$ в заданной строке $x[1..n]$. Здесь для каждой позиции $i \in 1..n$ в x просматривалась непустая подстрока $x[i'..i]$ ($i' \in 1..i$), что обеспечивало наилучшее k -приближенное совпадение с *префиксом* $p[1..j]$ ($j \in 1..m$). Как объяснялось в разделе 10.1, новая задача требует того, чтобы элементы столбцов инициализировались другим способом. Таким образом, массив c становится массивом $c[0..n, 1..m]$, но, как оказывается, для всех $i \in 1..n$ и $j \in 2..m$ все еще можно использовать тот же самый подход динамического программирования, который приводит к асимметричному соотношению

$$c[i, j] = \min_{1 \leq i' \leq i} d(x[i'..i], p[1..j]) \quad (13.12)$$

В этом случае k -приближенные совпадения паттерна p с x определяются с помощью позиций i в столбце m массива c , для которых $c[i, m] \leq k$.

Дальнейшее обобщение задачи происходит тогда, когда мы пытаемся определить для каждого $i \in 1..n_1$ и каждого $j \in 0..n_2$ совпадения с заданным уровнем k между непустыми подстроками $\mathbf{x}_1[i'..i]$ и $\mathbf{x}_2[j'..j]$ заданных строк $\mathbf{x}_1[1..n_1]$ и $\mathbf{x}_2[1..n_2]$. Тогда восстанавливается симметрия между первыми строкой и столбцом массива $\mathbf{c} = \mathbf{c}[1..n_1, 1..n_2]$, а каждый элемент $\mathbf{c}[i, j]$ удовлетворяет соотношению

$$\mathbf{c}[i, j] = \min_{\substack{1 \leq i' \leq i \\ 1 \leq j' \leq j}} d(\mathbf{x}_1[i'..i], \mathbf{x}_2[j'..j]). \quad (13.13)$$

Здесь о совпадениях, представляющих интерес, сообщают значения $\mathbf{c}[i, j] \leq k$.

Как увидим ниже, массив стоимостей \mathbf{c} , элементы которого определяются выражением (13.13), предоставляет основу для вычисления k -приближенных квадратов при $\mathbf{x} = \mathbf{x}_1 = \mathbf{x}_2$, но с массивом полезности \mathbf{b} , заменяющим массив стоимостей \mathbf{c} . В этом случае элементы $\mathbf{b}[i..j] \geq k$, где $i \neq j$, будут определять квадрат, образованный двумя различными подстроками, удовлетворяющими соотношению

$$\mathbf{x}[i'..i] \stackrel{(k)}{=} \mathbf{x}[j'..j].$$

Однако, чтобы иметь возможность эффективного выполнения вычислений, сначала надо показать, что аналог выражения (13.13) можно вычислить с использованием подхода динамического программирования леммы 9.1.1. Этот вопрос рассмотрен в следующем подразделе.

Массив полезности \mathbf{b}

Вначале рассмотрим реализацию вычисления выражения (13.13) на основе динамического программирования (см. лемму 9.1.1):

$$\mathbf{c}[i, j] \leftarrow \min\{\mathbf{c}[i-1, j] + d(\mathbf{x}_1[i], \varepsilon), \mathbf{c}[i, j-1] + d(\varepsilon, \mathbf{x}_2[j]), \mathbf{c}[i-1, j-1] + d(\mathbf{x}_1[i], \mathbf{x}_2[j])\}, \quad (13.14)$$

используемого теперь для всех $i \in 2..n_1$ и $j \in 2..n_2$. Но мы сразу же столкнемся с трудностями, которые станут очевидными, если, как показано в табл. 13.1, применим выражение (13.13) к нашему примеру $\mathbf{x}_1 = \mathit{rests}$, $\mathbf{x}_2 = \mathit{stress}$, где используется расстояние преобразования d_E .

Сначала заметим, что единственными возможными значениями в этом массиве являются “бесполезные” значения 0 и 1: согласно (13.13), первое значение является минимумом, который имеет место всякий раз, когда есть позиции $i' \leq i$ и $j' \leq j$, для которых $\mathbf{x}_1[i'..i] = \mathbf{x}_2[j'..j]$, в то время как значение 1 имеет место тогда, когда $\mathbf{x}_1[i] \neq \mathbf{x}_2[j]$. Таким образом, выражение (13.13) фактически говорит нам только о совпадении или несовпадении позиций $\mathbf{x}_1[i]$ и $\mathbf{x}_2[j]$, не давая какой-либо информации о любом предыдущем диапазоне позиций в этих двух строках,

Таблица 13.1. Массив стоимостей для $x_1 = rests$ и $x_2 = stress$, полученный с помощью d_E и (13.13)

j	1	2	3	4	5	6
i	s	t	r	e	s	s
1	r	1	1	0	1	1
2	e	1	1	1	0	1
3	s	0	1	1	1	0
4	t	1	0	1	1	1
5	s	0	1	1	1	0

в которых, возможно, имело место частичное совпадение. Но вторая трудность даже еще серьезнее: метод динамического программирования, кажется, больше не работает! Например, используя (13.14), мы не можем вычислить в табл. 13.1 $c[2,2] = 1$ или $c[5,5] = 0$.

Стандартный подход обычно обходит эти трудности. В работе [213] использовался массив, названный нами массивом полезности b , в котором представляет интерес максимальная польза, а не минимальные стоимости. При этом расстояния между буквами $d(\lambda, \mu)$ заменены на **баллы** $S(\lambda, \mu)$, которые положительны в желательных обстоятельствах (например, если $\lambda = \mu$) и отрицательны (или, по крайней мере, имеют меньшие положительные значения) в нежелательных обстоятельствах (например, если $\lambda \neq \mu$). Представим эти баллы как элементы **матрицы баллов** S , аналогичной матрице весов W , представленной в разделе 2.2, за исключением того, что теперь элементы с бóльшим значением, как правило, определяют совпадения букв, а не их несовпадения. Конечно, использование такой матрицы требует индексированного алфавита, как было указано в разделе 4.1.

Предположим, что для нашего примера матрица баллов S представима в следующем виде.

$$S = \begin{matrix} & \varepsilon & e & r & s & t \\ \varepsilon & 0 & -1 & -1 & -1 & -1 \\ e & -1 & 2 & -1 & -1 & -1 \\ r & -1 & -1 & 2 & -1 & -1 \\ s & -1 & -1 & -1 & 2 & -1 \\ t & -1 & -1 & -1 & -1 & 2 \end{matrix} \tag{13.15}$$

Таким образом, в этом простом случае штрафом (отрицательной полезностью), равным -1 , определяется любое несовпадение, удаление или вставка, в то время как полезность, равная $+2$, предоставляется за любое совпадение, за исключением, как обычно, пустой строки: $S(\varepsilon, \varepsilon) = 0$. Заметим, что в некоторых случаях,

например в биологических приложениях, некоторые баллы $S(\lambda, \mu)$ положительны, даже когда $\lambda \neq \mu$.

Таким же образом, как расстояние между буквами было обобщено до расстояния между строками, здесь мы используем баллы между буквами для определения баллов между строками. По-прежнему предполагаем, что строка $x_1 = x[1..n_1]$ преобразуется в строку $x_2 = x[1..n_2]$ с помощью операций редактирования отдельных букв: вставки, удаления, подстановки. Также предполагаем, что эти операции необратимы (нельзя выполнить подстановку $a \rightarrow b$, затем повторно заменить $b \rightarrow a$, даже если, выполняя эти операции, можно увеличить общее количество баллов). Таким образом, преобразование строки x_1 в x_2 можно представить в виде последовательности операций редактирования, при которых не происходит возврата к предыдущему состоянию — такая же модель ранее использовалась для расстояний.

Конечно, как и прежде, может быть больше одной последовательности операций редактирования, преобразующих x_1 в x_2 . Для каждой такой последовательности можно вычислить сумму всех баллов операций редактирования, как это представлено в матрице баллов. Тогда для строк x_1 и x_2 определим баллы $S(x_1, x_2)$ как максимальную сумму по всем последовательностям необратимых операций редактирования, которые преобразовывают x_1 в x_2 .

Теперь можно представить аналог выражения (13.13) для вычисления элементов *массива полезности* $\mathbf{b} = \mathbf{b}[1..n_1, 1..n_2]$: для всех $i \in 1..n_1, j \in 1..n_2$

$$\mathbf{b}[i, j] = \min_{\substack{1 \leq i' \leq i \\ 1 \leq j' \leq j}} S(x_1[i'..i], x_2[j'..j]) \quad (13.16)$$

Это выражение определяет элементы массива полезности, которые вычисляемы методом динамического программирования, используя для этого аналог выражения (13.14): для всех $i \in 2..n_1, j \in 2..n_2$

$$\mathbf{b}[i, j] \leftarrow \min\{0, \mathbf{b}[i-1, j] + S(x_1[i], \varepsilon), \mathbf{b}[i, j-1] + S(\varepsilon, x_2[j]), \mathbf{b}[i-1, j-1] + S(x_1[i], x_2[j])\}. \quad (13.17)$$

Исходные значения для $j = 1$ устанавливаются на основе выражений, аналогичных (10.2) и (10.3). Для удобства предположим, что существует единственное значение $\mathbf{b}[0, 1] = 0$, такое, что для всех $i = 1, 2, \dots, n_1$ можно вычислить

$$\mathbf{b}[i, 1] \leftarrow \max\{0, \mathbf{b}[i-1, 1] + S(x_1[i], \varepsilon), S(x_1[i], x_2[1])\}. \quad (13.18)$$

Подобным образом для $i = 1$ предполагаем, что $\mathbf{b}[1, 0] = 0$, и выполняем инициализацию с помощью следующего выражения: для $j = 1, 2, \dots, n_2$

$$\mathbf{b}[1, j] \leftarrow \max\{0, \mathbf{b}[1, j-1] + S(\varepsilon, x_2[j]), S(x_1[1], x_2[j])\}. \quad (13.19)$$

Теперь нас интересуют элементы $b[i, j] \geq k > 0$, другими словами, те элементы, которые обеспечивают полезность, по меньшей мере равную некоторому положительному порогу k .

Отметим, что выражения (13.17)–(13.19) накладывают на значения $b[i, j]$ условие неотрицательности. Это сделано из-за того, что в (13.16) мы стремимся определить такие i' и j' , которые увеличивают до максимума положительное количество баллов; таким образом, в вычислении на основе динамического программирования (13.17) предыдущие значения $b[i - 1, j]$, $b[i, j - 1]$ и $b[i - 1, j - 1]$ не будут представлять интереса, если они отрицательные, а для удобства вычислений отрицательные значения целесообразно преобразовать в нулевые. Поэтому сохранение в массиве b только неотрицательных значений гарантирует, что присвоение (13.17) эквивалентно выражению (13.16), и, таким образом, мы избегаем описанных выше сложностей, которые возникли при использовании формулы (13.14).

В табл. 13.2 показан массив полезности для строк из нашего примера, в котором выделены локальные максимумы $b[4, 2]$, $b[3, 5]$ и $b[5, 6]$. Эти значения определяют точные совпадения st и res , так же как и приближенные совпадения $rests$ с $ress$. Значение этих максимумов станет ясным позже, когда будем обсуждать понятие локальной оптимальности.

Таблица 13.2. Массив полезности для $x_1 = rests$ и $x_2 = stress$, полученный с помощью S и (13.16)

j	1	2	3	4	5	6	
i	s	t	r	e	s	s	
1	r	0	0	2	1	0	0
2	e	0	0	1	4	3	2
3	s	2	1	0	3	6	5
4	t	1	4	3	2	5	5
5	s	2	3	3	2	4	7

Сеточный граф

В разделе 9.5 был введен граф зависимостей — ориентированный граф, вершины которого помечались элементами $[i, j]$ массива стоимостей c ($i \in 1..n_1, j \in 1..n_2$), а дуги исходили из предшествующих элементов $[i - 1, j]$, $[i, j - 1]$ и $[i - 1, j - 1]$, значения которых явились вкладом в значение $c[i, j]$. Эта структура уточнялась формулами (9.8) и показана на рис. 9.2 для $x_1 = rests$ и $x_2 = stress$.

Здесь мы используем, по существу, ту же структуру, называемую теперь сеточным графом и переопределенную для массива полезности b . Для заданных строк x_1 и x_2 и связанной с ними матрицы баллов S **сеточный граф** $G = G(x_1, x_2, S)$

будет ориентированным ациклическим графом (V, E) , вершины V которого помечены элементами $[i, j]$ массива полезности \mathbf{b} , где каждая вершина $[i, j]$ имеет входящие дуги E тогда и только тогда, когда $\mathbf{b}[i, j] > 0$, и где применяются следующие правила, аналогичные (9.8).

$$\begin{aligned}
 ([i - 1, j], [i, j]) \in E, & \text{ если } \mathbf{b}[i, j] - \mathbf{b}[i - 1, j] = S(\mathbf{x}_1[i], \varepsilon), \\
 ([i, j - 1], [i, j]) \in E, & \text{ если } \mathbf{b}[i, j] - \mathbf{b}[i, j - 1] = S(\varepsilon, \mathbf{x}_2[j]), \\
 ([i - 1, j - 1], [i, j]) \in E, & \text{ если } \mathbf{b}[i, j] - \mathbf{b}[i - 1, j - 1] = S(\mathbf{x}_1[i], \mathbf{x}_2[j]).
 \end{aligned}
 \tag{13.20}$$

Таким образом, согласно этим правилам, каждая дуга в сеточном графе входит в путь, соответствующий набору операций редактирования, приводящему к максимальному количеству баллов $\mathbf{b}[i, j]$. Заметим, что, поскольку вершин $[0, j]$, $[i, 0]$ и $[0, 0]$ в этом графе нет, вершина $[1, 1]$ является источником (не имеет входящих дуг). В общем случае справедлива следующая лемма.

Лемма 13.3.1. Вершина $[i, j]$ в сеточном графе является источником тогда и только тогда, когда целые $i' = i$ и $j' = j$ приводят к максимальному значению $\mathbf{b}[i, j]$ в выражении (13.16).

Доказательство предложено дать в упражнении 13.3.4. ■

На рис. 13.7 показан сеточный граф для нашего примера строк $\mathbf{x}_1 = \text{rests}$ и $\mathbf{x}_2 = \text{stress}$. Всего в графе имеется 10 вершин-источников, из которых восемь фактически изолированы, — ситуация, которая невозможна в графах зависимостей.

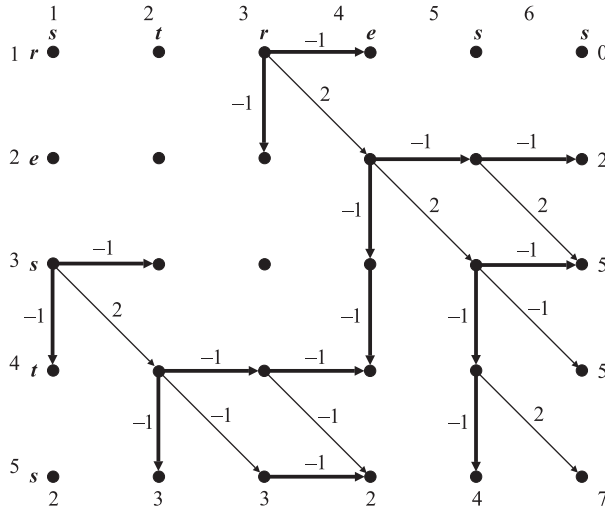


Рис. 13.7. Сеточный граф $G = G(\text{rests}, \text{stress}, S)$

Как видно на рис. 13.7, сеточный граф G создает структуру, которая показывает, каким образом можно вычислить каждый балл $b[i, j]$: после создания графа G проход по любому пути от вершины $[i', j']$ к вершине $[i, j]$ позволяет легко вычислить величину $b[i, j]$ путем суммирования баллов всех дуг, входящих в путь, и прибавления полученной суммы к величине $b[i, j]$. Назовем эту сумму баллов дуг **баллом пути** от $[i', j']$ к $[i, j]$ и обозначаем ее как $S_{ij}^{i'j'}$, т.е. $S_{ij}^{i'j'} = b[i, j] - b[i', j']$.

Таким образом, естественно при каждом выборе $[i, j]$ находить такие вершины-источники $[i', j']$, от которых начинаются пути, ведущие к $[i, j]$. Как показывает следующий результат, для фиксированных i и j эти вершины-источники приводят к максимальным значениям величин $S_{ij}^{i'j'}$.

Лемма 13.3.2. Пусть P обозначает путь в сеточном графе G от вершины-источника $[i', j']$ до некой выбранной вершины $[i, j]$. Тогда величина $S_{ij}^{i'j'}$ является единственным максимальным баллом пути.

Доказательство предложено дать в упражнении 13.3.4. ■

Фактически, на основе формулы (13.16) и этого результата, приходим к вполне удовлетворительному заключению: вершины-источники $[i', j']$ определяют целые числа i' и j' , такие, что подстроки $x_1[i'..i]$ и $x_2[j'..j]$ обеспечивают максимальный балл $S(x_1[i'..i], x_2[j'..j])$. Таким образом, вычисление путей с наивысшим баллом в сеточном графе $G = G(x_1, x_2, S)$ является естественным подходом к решению задачи, определяемой формулой (13.16).

Как видно из примера сеточного графа, показанного на рис. 13.7, вершина $[1, 3]$ является источником пути с наивысшим баллом, ведущего к каждой из вершин $[i, 6]$, $i \in 2..5$. Отсюда, принимая во внимание лемму 13.3.2, заключаем, что суффикс $x_2[3..6] = \text{ress}$ строки x_2 обеспечивает лучшее совпадение с каждой из подстрок $x_1[1..i]$. С другой стороны, заметим, что вершина $[4, 4]$ имеет два источника $[1, 3]$ и $[3, 1]$. Это говорит о том, что наилучшее совпадение суффиксов подстрок $x_1[1..4] = \text{rest}$ и $x_2[1..4] = \text{stre}$ можно получить двумя способами: или совпадением $x_1[1..4]$ с $x_2[3..4] = \text{re}$, или совпадением $x_1[3..4] = \text{st}$ с $x_2[1..4]$.

Принимая во внимание лемму 13.3.2, заманчиво предположить, что максимальный балл пути обязательно достигается как источниками $[i', j']$, так и стоками $[i, j]$. Но это не так — на рис. 13.7 вершина $[3, 5]$ лежит на пути от $[1, 3]$ до стока $[4, 6]$, но тем не менее $b[4, 6] = 5 < 6 = b[3, 5]$.

В общем случае мы не будем интересоваться всеми путями с высоким баллом, превышающим заданный порог k , а только их подмножеством. Как мы только что видели на рис. 13.7, есть пути с высоким баллом, ведущие и к вершинам $[3, 5]$ и $[4, 6]$, но путь, ведущий к вершине $[3, 5]$, кажется более интересным. Точно так же вершина $[4, 2]$, вероятно, будет более интересной, чем $[5, 3]$. В следующем под-

разделе рассмотрим важный вопрос, каким образом в сеточных графах определять самые интересные пути с высоким баллом.

Локальная оптимальность

При обсуждении точных раппортов (раздел 13.2) мы увидели, что понятие непродолжаемости играет определяющую роль в отборе выводимых данных, но при этом значительно уменьшает объем вывода. В данном случае для сопоставления с приближенным паттерном понятие локальной оптимальности играет аналогичную роль. И непродолжаемость, и локальная оптимальность зависят от максимального продолжения некоторых свойств как влево (случай б) в следующем определении), так и вправо (случай в)).

Определение 13.3.3. [84] Пара вершин $[i', j']$, $[i, j]$ в сеточном графе G называется **локально оптимальной**, если выполняются три условия:

- а) $S(\mathbf{x}_1[i'..i], \mathbf{x}_2[j'..j]) > 0$;
- б) для всех $I \in 1..i$ и $J \in 1..j$

$$S(\mathbf{x}_1[i'..i], \mathbf{x}_2[j'..j]) \geq S(\mathbf{x}_1[I..i], \mathbf{x}_2[J..j]) > 0;$$

- в) для всех $I \in i'..n_1$ и $J \in j'..n_2$,

$$S(\mathbf{x}_1[i'..i], \mathbf{x}_2[j'..j]) \geq S(\mathbf{x}_1[i'..I], \mathbf{x}_2[j'..J]).$$

■

Из определения (13.16) массива полезности видно, что условие б) будет удовлетворено путем перечисления всех вершин $[i', j']$, которые соответствуют максимуму $\widehat{b}[i, j]$, вычисленному для каждой вершины $[i, j]$. Условие в) можно удовлетворить путем вычисления массива полезности \widehat{b} для обратных строк \widehat{x}_1 и \widehat{x}_2 и перечисления всех вершин $[i', j']$, которые соответствуют каждому максимальному значению $\widehat{b}[i, j]$. Сравнение этих двух списков вершин с учетом преобразований позиций из \widehat{x}_1 в \mathbf{x}_1 и \widehat{x}_2 в \mathbf{x}_2 , определяет набор пар вершин, которые удовлетворяют обоим условиям определения 13.3.3.

Для сеточных графов эта процедура означает идентификацию вершин-источников в G , соответствующих каждой вершине $[i, j]$, как было описано выше, и выполнение того же для “обратного” сеточного графа \widehat{G} . Требование условия а) определения 13.3.3 позволяет избежать рассмотрения изолированных вершин. И чтобы не затруднять ни читателя, ни себя, мы опускаем дальнейшие подробности этой процедуры.

Однако полезно посмотреть, что получается, если понятие локальной оптимальности применить к нашему примеру. В табл. 13.3 показан массив полезности \widehat{b} для обратных строк $stser$ и $sserts$ с выделенными локальными максимумами в $\widehat{b}[3, 6]$ и $\widehat{b}[5, 4]$, которые соответствуют точному совпадению ts и хорошему

совпадению подстрок $stser$ и $sser$. Обращаясь к табл. 13.2, мы видим, что эти совпадения соответствуют точному совпадению st и хорошему совпадению подстрок $rests$ и $ress$, которые были выделены в массиве \mathbf{b} . Таким образом, эти совпадения могут определяться как локально оптимальные с помощью только что описанной процедуры или с помощью вывода $([1, 3], [5, 6])$ и $([3, 1], [4, 2])$, определяющего пары вершин в графе G , либо, в качестве альтернативы, с помощью вывода $([1, 1], [5, 4])$ и $([2, 5], [3, 6])$, определяющего пары вершин в графе \hat{G} .

Таблица 13.3. Массив полезности $\hat{\mathbf{b}}$ для $\hat{x}_1 = stser$ и $\hat{x}_2 = sserts$

j	1	2	3	4	5	6	
i	s	s	e	r	t	s	
1	s	2	2	1	0	0	2
2	t	1	1	0	0	2	1
3	s	2	3	2	1	1	4
4	e	1	2	5	4	3	3
5	r	0	1	4	7	6	5

Однако отметим отсутствие вывода точного совпадения раппорта ser , выделенного в табл. 13.2, поскольку $\hat{\mathbf{b}}[3, 2]$ не соответствует узлу-источнику в обратном сеточном графе \hat{G} и, следовательно, не удовлетворяет требованию локальной оптимальности. В нашем примере это не является серьезной проблемой: поскольку ser является суффиксом как строки $stser$, так и $sser$, о совпадении которых показывает вывод, то можно сказать, что ser не попал в вывод потому, что он является “продолжаемым влево”. Однако могут быть созданы примеры (см. [205]), в которых “хорошие” или “интересные” совпадения пропущены только потому, что, по меньшей мере, один конец совпавших строк не соответствует вершине-источнику в G или \hat{G} .

Таким образом, данное здесь определение локальной оптимальности для приближенного совпадения не совсем обеспечивает все то, что предоставляет понятие непродолжаемости для точного совпадения; тем не менее это самое полезное определение, которым мы располагаем.

Алгоритм со временем выполнения $\Theta(n^3)$

В предыдущих трех подразделах мы представили общую задачу (13.16) вычисления приближенных совпадений, в частности, локально оптимальных совпадений между подстроками двух заданных строк x_1 и x_2 . Теперь обратимся к упомянутому ранее специальному случаю этой задачи, когда $x = x_1 = x_2$. Здесь найденные приближенные совпадения будут подстроками той же самой строки x и поэтому будут приближенными квадратами.

“Специальность” рассматриваемого случая позволяет сделать некоторые упрощения. В частности, в результате симметрии, присущей сравнению x с x , достаточно вычислить только верхнюю треугольную часть матрицы полезности. Кроме того, поскольку нас не интересуют совпадения $x[i'..i]$ с $x[i'..i]$, все диагональные элементы в b можно положить равными нулю.

Кажется, что алгоритм решения нашей задачи ясен: для заданной строки $x = x[1..n_1]$ вычисляем соответствующую верхнюю треугольную часть матрицы полезности b за время порядка $\Theta(n^2)$, используя для этого упрощенные версии формул (13.17)–(13.19), затем вычисляем \hat{b} , применяя ту же методологию к \hat{x} , и, наконец, идентифицируем локально оптимальные пары, которые определяют приближенные квадраты. На вычисления потребуется время порядка $\Theta(n^2)$. Какие здесь могут быть проблемы?

Увы, проблема имеется, и она указана в предыдущем подразделе. Критерий локальной оптимальности является слишком ограничительным: в x могут иметь место приближенные квадраты, которые не обнаруживаются из-за того, что, по меньшей мере, одна из вершин в паре вершин $([i', j'], [i, j])$, которые определяют квадрат, не соответствует источнику в графах G или \hat{G} . В работе [205] приводится пример, в котором обнаруживаются и выводятся различные расщепленные или покрывающиеся квадраты, в то время как более значительный последовательный квадрат пропущен.

В алгоритме Шмидта для преодоления этой трудности просто устраняется вычисление покрывающихся квадратов в целом. Это гарантирует, что будут выводиться все последовательные и расщепленные квадраты. Конечно, в случае точного совпадения, каждый покрывающийся квадрат подстроки u обязательно подразумевает существование последовательного квадрата собственного префикса подстроки u , так что, в некотором смысле, теряется немного информации, если не вычисляются покрывающиеся квадраты. Точно так же при приближенном сравнении покрывающиеся квадраты обычно подразумевают более короткие последовательные квадраты.

Поэтому для данной строки $x[1..n]$ в каждой позиции $t \in 1..n - 1$ алгоритм Шмидта ведет поиск таких значений $i \in 1..t$ и $i \in t + 1..n$, которые дают максимальный балл

$$S(x[i..t], x[t + 1..j]). \tag{13.21}$$

В терминах сеточного графа это означает определение всех путей с наивысшим баллом от вершин $[i, t + 1]$ в столбце $t + 1$ графа G до вершин $[i, t]$ в строке t . Непосредственно величины

$$\min_{\substack{1 \leq i \leq t \\ t+1 \leq j \leq n}} S(x[i..t], x[t + 1..j]) \tag{13.22}$$

можно вычислить за время порядка $\Theta(n^2)$ для $\Theta(n)$ значений t . Следовательно, для определения всех k -приближенных последовательных квадратов в x потребуется время порядка $\Theta(n^3)$.

Конечно, чтобы определить более интересное подмножество этих квадратов, состоящее только из локально оптимальных квадратов, необходимо снова сформулировать и решить проблему обратного графа \hat{G} (как было описано выше) так, чтобы можно было определить пути с наивысшим баллом от вершины-источника обратного графа \hat{G} к вершине-источнику графа G . Таким образом, локально оптимальные k -приближенные последовательные квадраты в x можно вычислить за время порядка $\Theta(n^3)$.

В алгоритме Шмидта предложен способ сокращения времени вычисления до величины порядка $O(n^2 \log n)$, который в общих чертах описан ниже. По существу алгоритм Шмидта разбивает исходную задачу на $n - 1$ подзадач, которые являются специальными случаями вычисления формулы (13.16), когда сравниваются строки x_1^t и x_2^t ($t = 1, 2, \dots, n - 1$), где

$$x_1^t = x[1..t], x_2^t = x_2^t[1..n - t] = x[t + 1..n] \quad (13.23)$$

В t -й подзадаче P^t вычисляется массив полезности $\mathbf{b}^t = \mathbf{b}^t[1..t, 1..n - t]$ и соответствующий сеточный граф G^t , в котором определяются максимальные баллы

$$\min_{\substack{1 \leq i \leq t \\ t+1 \leq j \leq n}} S(x_1^t[i..t], x_2^t[1..n - t]), \quad (13.24)$$

эквивалентные баллам (13.22). Алгоритм Шмидта решает каждую из этих подзадач на основе массивов, которые обладают специальными свойствами, дающими возможность вначале вычислить их элементы, а затем обеспечить к ним быстрый доступ.

В работе [205] показано, каким образом, не жертвуя ни объемом памяти, ни временем вычисления, этот подход может быть расширен для вычисления не только последовательных, но и расщепленных, локально оптимальных k -приближенных квадратов — однако здесь мы удовлетворимся схемой алгоритма вычисления последовательных квадратов.

Массивы Монжа

Сеточный граф G^t , который соответствует максимальным баллам (13.24), будет содержать пути с наивысшим баллом от вершин-источников $[i, 1]$ в первом столбце до вершин-стоков $[t, j]$ в последней строке, где $i \in 1..t, j \in 1..n - t$. В работах [1, 14] предложено ввести в рассмотрение специальные DIST-массивы¹ $\mathbf{d}^t = \mathbf{d}^t[1..t, 1..n - t]$, где $\mathbf{d}^t[i, j]$ — балл пути с наивысшим баллом от вершины $[i, 1]$

¹Название DIST происходит от DISTance — расстояние, что раскрывает смысл этих массивов. — *Примеч. пер.*

до вершины $[t, j]$ в графе G^t ($t = 1, 2, \dots, n - 1$), и поэтому балл пути совпадает с величиной, определенной формулой (13.24); если не существует ни одного такого пути, тогда предполагается, что $d^t[i, j] = -\infty$. Каждый DIST-массив, будучи вычисленным, содержит всю информацию, необходимую для решения подзадачи P^t (в соответствии с (13.24)), в то время как все $n - 1$ DIST-массивов обеспечивают основу для вычисления всех последовательных квадратов в x (в соответствии с (13.21) и (13.22)).

DIST-массивы обладают особым свойством, которое уменьшает суммарное время их вычисления и, кроме того, предоставляет эффективный доступ к вычисленным элементам.

Определение 13.3.4. Матрица $M[1..n_1, 1..n_2]$ называется матрицей **Монжа** (Monge), если выполняется одно из следующих условий:

- а) $M[i, j] + M[i + 1, j + 1] \geq M[i, j + 1] + M[i + 1, j]$ для всех $i \in 1..n_1 - 1$ и $j \in 1..n_2 - 1$;
- б) $M[i, j] + M[i + 1, j + 1] \leq M[i, j + 1] + M[i + 1, j]$ для всех $i \in 1..n_1 - 1$ и $j \in 1..n_2 - 1$. ■

Лемма 13.3.5. Каждый DIST-массив d^t является матрицей Монжа.

Доказательство. Покажем, что массив d^t удовлетворяет условию б) определения 13.3.4. Предположим, что для некоторых $i < t$ и $j < n - t$ $d^t[i, j] > 0$ и элемент $d^t[i, j]$ равен баллу пути $\pi_{i,j}$ с наивысшим баллом от вершины $[i, 1]$ до вершины $[t, j]$ в графе G^t . Если, кроме того, $d^t[i + 1, j + 1] > 0$ и равен баллу пути $\pi_{i+1,j+1}$ с наивысшим баллом от $[i + 1, 1]$ до $[t, j + 1]$ в G^t , тогда, как показано на рис. 13.8, эти пути должны пересечься в некоторой вершине $[i', j']$ графа G^t .

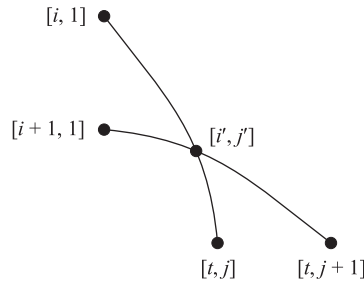


Рис. 13.8. Пересечение путей в сеточном графе G^t

Следовательно, должны существовать пути (не обязательно с наивысшим баллом) $\pi_{i,j+1}$ от вершины $[i, 1]$ до вершины $[t, j + 1]$ и $\pi_{i+1,j}$ — от $[i + 1, 1]$ до $[t, j]$,

сформированные из подпутей, которые проходят через вершину $[i', j']$. Обозначив как $|\pi|$ балл пути π , мы увидим, что

$$|\pi_{i,j}| + |\pi_{i+1,j+1}| = |\pi_{i,j+1}| + |\pi_{i+1,j}|,$$

где пути в левой части равенства являются по определению путями с наивысшим баллом, в то время как пути в правой части таковыми могут не быть. Этот вытекает из положительности величин $\mathbf{d}^t[i, j]$ и $\mathbf{d}^t[i + 1, j + 1]$.

Если хотя бы одна из величин $\mathbf{d}^t[i, j]$ и $\mathbf{d}^t[i + 1, j + 1]$ не положительна, то по определению она равна $-\infty$. В этом случае выполняется неравенство

$$\mathbf{d}^t[i, j] + \mathbf{d}^t[i + 1, j + 1] = -\infty \leq \mathbf{d}^t[i, j + 1] + \mathbf{d}^t[i + 1, j].$$

Лемма доказана. ■

Более важным для последующего анализа является аналогичный результат, который связывает значения в массиве \mathbf{d}^t со значениями в массиве \mathbf{d}^{t+1} .

Лемма 13.3.6. Для всех $t \in 1..n - 2$, $i \in 1..t$ и $j \in 1..n - t$

$$\mathbf{d}^t[i, j] + \mathbf{d}^{t+1}[i + 1, j] \geq \mathbf{d}^t[i + 1, j] + \mathbf{d}^{t+1}[i, j].$$

Доказательство основано на тех же аргументах, которые были использованы при доказательстве леммы 13.3.5. ■

Здесь мы заканчиваем подробное описание алгоритма Шмидта. Лемма 13.3.6 показывает, что мы можем определить массивы $X_j[t, i]$ ($j = 1, 2, \dots, n - t$), которые являются матрицами Монжа и связывают значения элементов массивов \mathbf{d}^{t+1} и \mathbf{d}^t . Это обеспечивает возможность применения эффективных итерационных процедур, а также эффективный доступ к элементам массивов X_j . Методы, используемые для выполнения этих вычислений, требуют построения и использования совокупности двоичных деревьев, которые “представляют” массивы \mathbf{d}^t . Объяснение и доказательство корректности этих методов являются длинными, сложными и специализированными, поэтому они не рассматриваются в этой книге. Однако приведем без доказательства два основных результата, касающиеся времени выполнения этих методов (формулировки следующих теорем менее общие по сравнению с оригинальными их версиями в [205]). Первая теорема говорит о вычислении матриц \mathbf{d}^t .

Теорема 13.3.7. Вычисление DIST-матриц \mathbf{d}^t ($t = 1, 2, \dots, n - 1$) может быть выполнено за время порядка $O(n^2 \log n)$ с использованием памяти порядка $O(n^2 \log n)$. Поиск любого отдельного элемента в любом массиве \mathbf{d}^t выполняется за время порядка $O(\log t)$, поиск любого элемента в любом столбце любого массива \mathbf{d}^t требует времени порядка $O(1)$. ■

Вторая теорема касается последующего вычисления максимумов строк и столбцов в матрицах d^t ; эти максимумы необходимы алгоритму Шмидта для вычисления путей с наивысшим баллом.

Теорема 13.3.8. Все максимумы строк и столбцов в DIST-матрицах d^t ($t = 1, 2, \dots, n - 1$) могут быть вычислены за время порядка $O(n^2 \log n)$ с использованием памяти порядка $O(n^2)$. ■

Модификация алгоритма Мейна–Лоренца

В этом разделе, избегая излишних подробностей, еще раз возвратимся к первоначальной формулировке задачи приближенных квадратов (13.22) и покажем, каким образом методику алгоритма Мейна–Лоренца определения точных квадратов (подраздел 12.1.2) можно применить в алгоритме Шмидта для определения приближенных квадратов.

Напомним, что алгоритм Мейна–Лоренца находит квадраты в любой подстроке u строки x , сначала расщепляя u на левую часть $u_1 = u[1.. \lfloor n/2 \rfloor]$ и на правую часть $u_2 = u[\lfloor n/2 \rfloor + 1..n]$ и вычисляя квадраты, покрывающие u_1 и u_2 , затем решая рекурсивные подзадачи для квадратов, распложенных или в u_1 , или в u_2 . Алгоритм Шмидта определяет пути с наивысшим баллом от вершин столбца $t + 1$ сеточного графа G к вершинам строки t . Далее вспомним, что поскольку наша проблема рефлексивна (мы сравниваем подстроки строки x с ее же подстроками), поэтому необходимо рассмотреть только верхнюю треугольную часть как массива полезности b , так и сеточного графа G . Таким образом, как показано на рис. 13.9, можно применить стратегию алгоритма Мейна–Лоренца для расщепления графа G на подграфы на основе положения путей с наивысшим баллом. Пути, которые пересекаются в “средней точке”, вычисляются сразу. Опишем другие возможные пути:

- а) непосредственно вычисляются и выводятся приближенные квадраты, соответствующие путям с наивысшим баллом не менее k , которые пересекают или столбец $\lfloor n/2 \rfloor + 1$, или строку $\lfloor n/2 \rfloor$;
- б) обрабатываются как рекурсивные подзадачи пути с наивысшим баллом, которые ограничиваются или верхним левым треугольником, или нижним правым треугольником;
- в) пути с наивысшим баллом, которые ограничиваются ближним квадратом, определяемым неравенствами $0 \leq i \leq \lfloor n/2 \rfloor$, $\lfloor n/2 \rfloor + 1 \leq j \leq n$, не могут быть решениями, поскольку нет такого значения t в этой области, которая удовлетворяет требованиям задачи. Таким образом, эта часть графа G может быть отброшена.

Как указано в теоремах 13.3.7 и 13.3.8, пути с наивысшим баллом типа a можно вычислить за время порядка $O(n^2 \log n)$. Таким образом, обозначив как t_n максимальное время, необходимое алгоритму Шмидта для вычисления приближенных

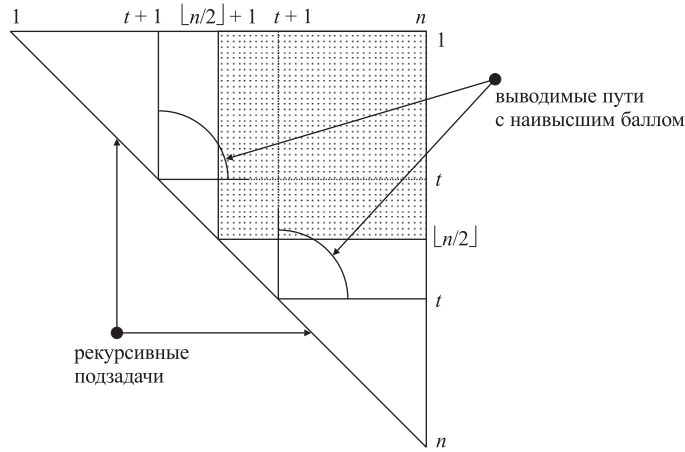


Рис. 13.9. Рекурсивное расщепление сеточного графа в соответствии с алгоритмом МЛ

последовательных квадратов в любой строке $x[1..n_1]$, с учетом вычисления путей типа \bar{b}) и \bar{v}), величину t_n можно определить путем решения рекуррентного соотношения

$$t_n \leq k_2 n^2 \log n + 2t_{\lfloor n/2 \rfloor}, \quad (13.25)$$

где k_2 — положительная константа. Поскольку существует положительная константа k_1 , такая, что $t_1 \leq k_1$, последняя задача является аналогом рекуррентной задачи (12.2) для алгоритма Мейна–Лоренца. Внимательный читатель может вспомнить решение несколько более общей версии соотношения (13.25) из упражнения 12.1.6, на основании которого можно утверждать, что временная сложность алгоритма Шмидта составляет $t_n \in O(n^2 \log n)$.

Детали алгоритма Шмидта, включая незатронутые здесь аспекты, можно найти в работе [205]. Сейчас, на основе теоремы 13.3.7 и предыдущего анализа, сформулируем конечный результат.

Теорема 13.3.9. Алгоритм Шмидта вычисляет локально оптимальные k -приближенные непокрывающиеся квадраты в заданной строке $x[1..n]$ за время порядка $O(n^2 \log n)$ с использованием памяти такого же порядка. ■

Обсуждение

В этом разделе мы постарались представить понятную схему алгоритма Шмидта, не слишком вдаваясь в его технические подробности. Мы выбрали для обсуждения этот алгоритм, потому что он самый быстрый и пригоден для решения проблем, представляющих практический интерес, т.е. алгоритм Шмидта является и полезным, и поучительным. Более раннему алгоритму [111] для вычисления

всех (не только локально оптимальных) непокрывающихся приближенных квадратов требовалось время порядка $O(n^2 \log^2 n)$ и объем памяти порядка $O(n^2 \log n)$; в работе [30] требуемый объем памяти был уменьшен до величины порядка $O(n^2)$. Алгоритм для вычисления всех “хороших” приближенных последовательных раппортов за время порядка $O(kn \log k \log(n/k))$ предложен в [145], но его можно использовать только для расстояния преобразования.

Одной из самых привлекательных особенностей алгоритма Шмидта является критерий локальной оптимальности, применяемый для ограничения всех раппортов до приближенных квадратов, представляющих наибольший интерес. Однако следует отметить, что этот критерий ставит проблемы, которые не присущи критериям непродолжаемости, используемым для точных раппортов. Возможно, переопределение понятия локальной оптимальности может помочь в поиске более удовлетворительного алгоритма для решения задачи приближенных квадратов.

Упражнения 13.3

1. Используя расстояние преобразования, создайте пример такой строки $x = u_1 u_2 u_3$, где $u_1 u_2$ и $u_2 u_3$ являются 1-приближенными квадратами, в то время как $u_1 u_3$ таковым не является. Затем обобщите этот пример для произвольного $k > 1$. Можно ли привести такой пример, где все эти строки имеют различную длину?
2. Предположим, что разрешаются обратимые операции редактирования. Создайте пример, в котором значение $S(x_1, x_2)$ становится сколь угодно большим.
3. В этом разделе был описан массив полезности из неотрицательных значений, базирующийся на положительных баллах для совпадений и, в основном, на отрицательных баллах для несовпадений. Покажите, что массив стоимостей, состоящий из неположительных значений, будет работать точно так же, и введите необходимые определения.
4. Докажите леммы 13.3.1 и 13.3.2.
5. Примеры, данные в этом разделе, имеют вершины-источники $[i', j']$ для путей только при $i' = 1$ и $j' = 1$. Приведите примеры, где вершины-источники находились бы в произвольном заданном месте в сеточном графе.
6. Нарисуйте сеточный граф \widehat{G} , соответствующий строкам $\widehat{x}_1 = stser$ и $\widehat{x}_2 = sserts$. Затем используйте его для определения всех локально оптимальных пар в графе G , соответствующем строкам $rests$ и $stress$.
7. Покажите, что при точном совпадении каждый покрывающийся квадрат подстроки u определяет последовательный квадрат префикса u . Можно ли создать пример для приближенного сравнения, в котором это не выполняется?

8. Сформулируйте задачу, обратную (13.22), и объясните, каким образом ее можно использовать для вычисления локально оптимальных k -приближенных последовательных квадратов за время порядка $\Theta(n^3)$.
9. Докажите лемму 13.3.6.

13.4 k -приближенные периоды — алгоритмы Сима–Илиопулоса–Парка–Смита

Первый алгоритм, описанный в этой книге, вычислял массив граней заданной строки x , определяя за линейное время каждый период, включая самый короткий, каждого префикса x . В этом заключительном разделе книги мы завершаем почти полный цикл. Теперь попробуем вычислить *приближенные* периоды строки x , что является более трудной задачей [209], особенно в тех случаях, когда необходимо, чтобы вычисленные элементы включали нечто, в некотором смысле, подобное самому короткому точному периоду, “наилучшему” или наиболее “содержательному”. Мы увидим, что проблемы, рассматриваемые в этом разделе, действительно находятся на грани познания: можно сформулировать несколько интересных проблем, но в настоящее время даже наиболее простые из них могут быть решены только частично.

Однако прежде чем вплотную заняться этими проблемами, необходимо ответить на следующий вопрос.

Что мы вкладываем в понятие “приближенный период” (или длина “приближенной образующей”), если приближенная образующая не определяется длиной периода?

Рассмотрим в качестве примера простую строку

$$x = abcabc \tag{13.26}$$

Предположим, что необходимо найти для расстояния преобразования d_E 1-приближенную образующую строки x , т.е. строку u , такую, что для некоторой декомпозиции $x = x_1x_2 \dots x_r$ ($r \geq 1$) строки x выполняются неравенства $d_E(u, x_j) \leq 1$ для всех $j \in 1..r$.

Для нашей строки (да и для любой другой), несомненно, “лучшее” решение этой задачи состоит в том, чтобы выбрать $u = x$ и $r = 1$, тогда $d_E(u, x) = d_E(u, x_1) = 0$. Если это решение исключить, то для нашей строки (как и для любой другой) можно выбрать $u = \lambda$ (λ — произвольная буква), тогда $d_E(u, x[i]) = d_E(\lambda, x[i]) = 1$ для всех позиций i в x . Продолжая этот процесс, обнаруживаем, что, по меньшей мере для расстояния преобразования, k -приближенные образующие любой строки x всегда будут включать саму строку x , так же как все

строки u длиной больше k . Это наблюдение вряд ли можно назвать полезным или содержательным результатом.

Заметим далее, что в нашем примере (13.26) можно также выбрать $u = ac$ и выполнить декомпозицию $x = x_1x_2x_3 = (ab)(ac)(bc)$, тогда $d_E(u, x_j) = d_E(ac, x_j) = 1$ для каждого $j \in 1..3$. Таким образом, ac — это 1-приближенная образующая, а не префикс строки (13.26). Или, если мы несколько изменим наш пример, взяв строку $y = abbcad$, то здесь (хотя у строки y нет подходящей подстроки u длиной 2, претендующей на роль 1-приближенной образующей строки y) снова можно выбрать строку $u = ac$ (не являющейся подстрокой строки y), и тогда, согласно декомпозиции $y = y_1y_2y_3 = (ab)(bc)(ad)$, $d_E(u, y_j) = d_E(ac, y_j) = 1$ для каждого $j \in 1..3$. Таким образом, ac также будет 1-приближенной образующей строки y .

В обоих этих примерах знание длины возможной образующей дает немного информации: что действительно необходимо знать — так это непосредственно саму образующую. Следовательно, наш первый шаг в формулировании проблемы k -приближенных периодов заключается в устранении любых ссылок на периоды! В этом разделе мы будем иметь дело с *приближенными образующими* заданной строки x .

В предыдущем разделе мы видели, что для того, чтобы понятие приближенного раппорта имело смысл, было необходимо тщательно продумать способ определения понятия “расстояние” — здесь мы также видим, что соответствующее определение расстояния является критическим предусловием для того, чтобы оперировать приближенными образующими. Далее в этом разделе будем использовать понятие *взвешенного расстояния* d_W (обобщение расстояния преобразования), где используется матрица весов W с неотрицательными значениями, как определено в разделе 2.3. При изменении значений матрицы W , d_W может преобразовываться в другие обычные виды расстояния: d_E , d_L и d_H . Следуя [209], вводим новое полезное понятие.

Определение 13.4.1. Для двух произвольных непустых строк x и y , определенных на алфавите A , и данной функции расстояния d_W , определенной на буквах алфавита A , значение

$$\delta_W(x, y) = \frac{2d_W(x, y)}{|x| + |y|}$$

называется **относительным расстоянием** между x и y . ■

Таким образом, относительное расстояние между двумя строками — это взвешенное расстояние (или *абсолютное расстояние*), деленное на среднюю длину строк. Расстояние δ_W сохраняет свойства метрики, если таковыми обладает d_W (см. упражнение 13.4.1). Но более важно заметить, что использование относительного расстояния (а не взвешенного) при вычислении расстояния между потенциальными периодами u и подстроками x_j строки x в значительной степени

изменяет полученные результаты. В нашем предыдущем примере (13.26) находим, что для строк длиной 1 значения всех видов расстояния неизменны.

$$\delta_E(\lambda, \mu) = d_E(\lambda, \mu) = 1, \text{ если } \lambda \neq \mu, \text{ иначе } \delta_E(\lambda, \mu) = d_E(\lambda, \mu) = 0. \quad (13.27)$$

Однако расстояния для некоторых строк длиной 2 уменьшаются:

$$\delta_E(ac, ab) = \delta_E(ac, bc) = \delta_E(ac, ad) = 1/2.$$

Таким образом, если бы мы искали образующую на основе относительного расстояния, вероятно, выбрали бы строку, такую, как ac , на которой достигается меньшее значение или абсолютного, или относительного расстояний, чем это может быть достигнуто на одной букве. Этот пример снова поднимает вопрос о том, какие критерии должны использоваться для выбора образующих. Кроме того, отметим, что выбор $u = x$ все еще приводит к результату $\delta_E(u, x) = 0$, и поэтому также необходимо “придумать” критерий, который ограничивал бы длину u .

Для удовлетворения этих требований дадим следующее определение приближенной образующей, немного изменив определение из [209].

Определение 13.4.2. Для вещественного числа $k \geq 0$ и целого числа $\rho \geq 1$, строка u будет (ρ, k) -**приближенной образующей** строки x , если существует декомпозиция $x = x_1 x_2 \dots x_r$ ($r \geq \rho$) на непустые факторы x_j ($j \in 1..r$), такие, что $\delta_W(u, x_j) \leq k$ для всех $j \in 1..r - 1$, и при этом $\delta_W(u', x_r) \leq k$ для некоторого непустого префикса u' строки u . ■

В “настойчивом” требовании этого определения, что u должна быть достаточно близка к *каждому* фактору x_j , фиксируется свойство “равномерности”, которому удовлетворяет точная образующая. Отметим, что это определение разрешает строке x иметь суффикс, который приближенно совпадает с префиксом строки u .

Чтобы узнать, какое влияние оказывает выбор параметра ρ на определение образующей, привлечем свою интуицию и рассмотрим пример строки $x = abcasaabaaca$. В табл. 13.4 приведены приближенные образующие, удовлетворяющие различным значениям параметров ρ и k . В этом примере, чтобы получить “интересные” приближенные образующие, следует выбрать $\rho \geq 3$, но не ясно, будет ли какая-либо польза при $\rho = 4$, поскольку при этом значении ρ исключается возможная интересная образующая $abcasa$. С другой стороны, могут быть такие случаи, где, например, x имеет длинную точную или приближенную образующую, повторенную несколько раз, — тогда желательно большее значение ρ . Отметим, что из этого примера также не ясно, будут ли образующие, соответствующие минимальному относительному расстоянию k (независимо от значения ρ), обязательно наиболее интересными!

Из определения 13.4.2 прямо вытекает определение понятия “минимальной” приближенной образующей.

Таблица 13.4. (ρ, k) -приближенные образующие для $x = abcacaabaaca$

ρ	k	Найденные образующие
1	0	$x, abcacaabaaca$
2	0	$abcacaabaaca$
3	1/3	$aba, abcaca, abaaca$
3	1/6	$abcaca, abaaca$
4	1/3	aba
4	1/6	—

Определение 13.4.3. (ρ, k) -приближенная образующая u строки x называется ρ -минимальной, если не существует (ρ, k') -приближенной образующей строки x , где $k' < k$. ■

Отметим, что в табл. 13.4 $abcaca$ и $abaaca$ являются 3-минимальными образующими, в то время как aba — 4-минимальная.

Теперь можно сформулировать три задачи, связанные с (ρ, k) -приближенными образующими [209], которые перечислим в порядке возрастания их сложности.

- P1.** Заданы непустые строки x и u , целое число ρ и функция относительного расстояния δ_W . Требуется найти такое минимальное значение k , при котором u будет (ρ, k) -приближенной образующей строки x .
- P2.** Заданы непустая строка x , целое число ρ и функция относительного расстояния δ_W . Требуется найти подстроку u строки x (если она существует), которая является ρ -минимальной образующей строки x .
- P3.** Заданы непустая строка x , целое число ρ и функция относительного расстояния δ_W . Требуется найти произвольную строку u (если она существует), которая является ρ -минимальной образующей строки x .

В следующих трех подразделах мы вкратце обсудим эти три задачи только для того, чтобы увидеть, что в настоящее время ни для одной из них не существует полного решения! Однако мы увидим, что для задач P1 и P2 имеются частичные решения, полученные путем непосредственного применения алгоритмов, рассмотренных в главе 9.

Задача P1 — вычисление минимального k

Для единообразия мы сформулировали эту задачу для функции относительного расстояния δ_W , но поскольку и x и u определены, ее можно точно также сформулировать для функции абсолютного расстояния d_W .

Описанный в этом подразделе первый алгоритм Сима–Илиопулоса–Парка–Смита (Sim–Iliopoulos–Park–Smyth, сокращенно — алгоритм СИПС1) решает зада-

чу P1 только для *некоторого* значения ρ , не заданного заранее. Другими словами, он вычисляет такое минимальное значение k , что для некоторого ρ \mathbf{u} является (ρ, k) -приближенной образующей строки \mathbf{x} . Выполнение алгоритма СИПС1 делится на два этапа. На первом этапе применяется повторяемое вычисление расстояний между строками, на втором — итерационно вычисляется минимальное пороговое значение k .

На первом этапе вычисляются n стандартных массивов стоимостей (раздел 9.1)

$$\mathbf{c}^{(i)}[0..n - i + 1, 0..m], i = 1, 2, \dots, n, \quad (13.28)$$

в соответствии с расстояниями $\delta_W(\mathbf{x}[i..n], \mathbf{u})$. Эти вычисления определяют расстояния $\delta_W(\mathbf{x}[i..h], \mathbf{u})$ между $\mathbf{x}[i..h]$ и \mathbf{u} для каждого $h \in i..n - 1$; эти расстояния представлены в m -х столбцах массивов $\mathbf{c}^{(i)}$, $i = 1, 2, \dots, n - 1$. Поскольку конечное вхождение \mathbf{u} в \mathbf{x} должно быть только непустым префиксом \mathbf{u} , мы можем искать $(n - i + 1)$ -ю (последнюю) строку каждого массива $\mathbf{c}^{(i)}$ для минимального значения, назовем его $\text{minpref}(\mathbf{u}, i, n)$ — это будет наименьшим расстоянием между любым префиксом \mathbf{u} и $\mathbf{x}[i..n]$. Таким образом, используя последний столбец и один элемент из последней строки каждого массива $\mathbf{c}^{(i)}$, можно сформировать верхний треугольный массив стоимостей $\mathbf{c}[1..n, 1..n]$, где для $i \in 1..n$ $\mathbf{c}[i, h] = \delta_W(\mathbf{x}[i..h], \mathbf{u})$, если $h \in i..n - 1$, и $\mathbf{c}[i, h] = \text{minpref}(\mathbf{u}, i, n)$, если $h = n$.

Вычисление массивов $\mathbf{c}^{(i)}$ с использованием алгоритма Вагнера–Фишера (раздел 9.2) требует времени порядка $\Theta(mn^2)$, вычисление массива \mathbf{c} требует времени порядка $\Theta(n^2)$ — следовательно, общее время выполнения первого этапа алгоритма СИПС1 должно иметь порядок $\Theta(mn^2)$.

Теперь рассмотрим второй этап алгоритма. Предположим, что \mathbf{u} является приближенной образующей строки $\mathbf{x}[1..i]$. Разобьем $\mathbf{x}[1..i]$ на факторы, у которых наибольшее расстояние (вычисляемое как δ_W) до \mathbf{u} равно, скажем, k_i . Тогда для каждого $h \in i + 1..n - 1$ можно определить соответствующее расстояние k_h путем вычисления

$$\max\{k_i, \delta_W(\mathbf{x}[i + 1..h], \mathbf{u})\} = \max\{k_i, \mathbf{c}[i + 1, h]\},$$

при этом для $h = n$ вычисляется

$$\max\{k_i, \text{minpref}(\mathbf{u}, i, n)\} = \max\{k_i, \mathbf{c}[i + 1, h]\}.$$

Если, кроме того, для каждого $i \in 1..h - 1$ факторы подстроки $\mathbf{x}[1..i]$ выбраны таким способом, чтобы минимизировать k_i по всем выборам факторов, то мы видим, что вычисление

$$k_h = \min_{1 \leq i \leq h} \{\max\{k_i, \mathbf{c}[i + 1, h]\}\} \quad (13.29)$$

должно приводить к минимальному расстоянию k_h для $x[1..h]$ по всем выборам факторов. Это классический подход динамического программирования, который можно реализовать за время порядка $\Theta(n^2)$. Следовательно, для двух этапов алгоритма СИПС1 справедлива теорема.

Теорема 13.4.4. Для данной функции расстояния δ_W и непустых строк x и u алгоритм СИПС1 вычисляет минимальное значение k , такое, что для некоторого ρ u является (ρ, k) -приближенной образующей строки x . Для выполнения этого алгоритма требуется время порядка $\Theta(mn^2)$ и память порядка $\Theta(n^2)$. ■

Ввиду явного подобия между первым этапом алгоритма СИПС1 и вычислением массивов полезности b^t алгоритма Шмидта, кажется вполне вероятным, что здесь также можно использовать DIST-массивы для уменьшения времени порядка $\Theta(mn^2)$, необходимого для выполнения алгоритма на первом этапе, до времени порядка $O(mn \log n)$, следовательно, для уменьшения полного времени до величины порядка $\Theta(n^2)$. Эта увлекательная возможность оставлена в качестве исследовательского проекта упражнения 13.4.4.

Пока мы предполагали, что функция расстояния δ_W является обобщенной, основанной на произвольной матрице баллов. Однако если вместо расстояния преобразования используются расстояния δ_E или d_E , то результатом будет существенная экономия времени и памяти. Принимая во внимание (13.27), не обязательно вычислять расстояние преобразования между u и подстроками x длиной больше $2m$: такие расстояния обязательно превысили бы значение m и поэтому не могли бы способствовать вычислению минимального значения k_i . Таким образом, каждый из n массивов стоимостей (13.28) будет иметь максимальный размер $2m \times m$, в то время как пересмотренный массив стоимостей c будет иметь размер $n \times 2m$. Время, необходимое для выполнения первого этапа, уменьшается, соответственно, до величины порядка $\Theta(m^2n)$, а на втором этапе нужно рассмотреть не более $2m$ элементов массива c для каждого значения $h = 1, 2, \dots, n$, что приведет ко времени порядка $O(mn)$. Следовательно, использование расстояния преобразования уменьшает полное время, необходимое для выполнения алгоритма СИПС1, до $\Theta(m^2n)$, а необходимый объем памяти — до $\Theta(mn)$.

Но другой подход, реализованный в алгоритме СИПС2, уменьшает это время еще в большей степени. Вспомним, что в конце раздела 9.5, была сделана ссылка на “пошаговый” алгоритм вычисления расстояний между строками [144], который для данного расстояния $d_E(x, u)$ может вычислить любое расстояние $d_E(\lambda x, u)$ и $d_E(x\lambda, u)$ для любой буквы λ за время $O(|x| + |u|)$. Теперь предположим, что такой алгоритм применен для вычисления расстояния преобразования между u и следующей последовательностью $n + 1$ строк; каждая длиной не более $2m$ сформирована из предыдущей путем присоединения слева одной буквы:

$$\varepsilon, x[n], x[n - 1..n], \dots, x[n - 2m + 1..n], x[n - 2m..n - 1], \dots, x[1..2m].$$

Эти вычисления делают доступными все расстояния между u и каждой подстрокой x длиной не более $2m$. Поскольку $|u| = m$, время, требуемое этому алгоритму для каждого вычисления, будет равно $O(2m + m) = O(m)$, следовательно, полное время имеет порядок $O(mn)$. Далее матрица c формируется, как и в алгоритме СИПС1, за время порядка $\Theta(mn)$, затем на втором этапе обрабатывается за время порядка $\Theta(mn)$. Таким образом, имеем следующую теорему.

Теорема 13.4.5. Используя расстояние преобразования и данные непустые строки x и u , алгоритм СИПС2 вычисляет минимальное значение k , такое, что для некоторого ρ u является (ρ, k) -приближенной образующей строки x . На его выполнение требуется время порядка $\Theta(mn)$ и память объемом порядка $\Theta(m^2)$. ■

Задача P2 — вычисление ρ -минимальной подстроки u

Один из очевидных способов определить подстроку u строковой последовательности x (для некоторого ρ , не выбранного заранее) как ρ -минимальную образующую строки x состоит в исследовании всех возможных подстрок строки x путем использования алгоритмов СИПС1 или СИПС2 и в последующем выборе подстроки u , который приводит к минимальному пороговому значению k в качестве периода. Напомним из упражнения 1.2.10, что в строковой последовательности x имеется $O(n^2)$ различных подстрок средней длины $O(n)$. Таким образом, или алгоритм СИПС1 (для общего расстояния δ_W), или СИПС2 (для расстояния преобразования δ_E) может быть выполнен $O(n^2)$ раз с использованием образующей, средняя длина которой $m \in O(n)$. Следовательно, имеет место такой результат.

Теорема 13.4.6. Для некоторой, заранее не определенной величины ρ задача P2 может быть решена за время порядка $O(n^5)$ с использованием расстояния δ_W и за время порядка $O(n^4)$ с использованием расстояния преобразования δ_E . ■

В работе [209] для случая расстояния преобразования приводится другой метод решения задачи P2, но со временем выполнения все еще порядка $O(n^4)$.

Задача P3 — вычисление произвольной ρ -минимальной строки u

Эта задачу мы особенно хотели бы решить с приемлемой эффективностью, но, как оказывается, в такой формулировке эта задача является NP-полной²! Для получения более подробной информации см. [209].

Замечание

В этом разделе мы выделили алгоритмы, которые касаются новой области анализа строковых последовательностей, где проблемы еще недостаточно хорошо

²NP-полной называется задача, которая эквивалентна задаче полного перебора всех возможных вариантов ее решения. — *Примеч. ред.*

определены, а решений еще меньше. За курс из 13 глав мы достигли этого пункта длинным, но, тем не менее, по структуре довольно естественным и понятным маршрутом. Однако остается тот факт, что проблемы, только что рассмотренные, находятся среди наиболее простых обобщений вычисления периодичности, алгоритм которого первым рассматривался в этой книге. Маршрут был длинным, а покрытое расстояние коротким.

Упражнения 13.4

1. Покажите, что δ_W удовлетворяет условиям (2.2)–(2.5) метрики тогда и только тогда, когда d_W удовлетворяет этим же условиям.
2. Покажите для расстояния преобразования, что в общем случае $0 \leq \delta_E(\mathbf{x}, \mathbf{y}) < 2$, а в случае $|\mathbf{x}| = |\mathbf{y}|$ $\delta_E(\mathbf{x}, \mathbf{y}) \leq 1$.
3. Объясните, почему в алгоритме СИПС1 необходимый объем памяти достигает величины порядка $\Theta(n^2)$, а не порядка $\Theta(mn^2)$.
4. **Исследовательский проект.** Исследуйте возможность использования DIST-массивов (раздел 13.3) для уменьшения времени выполнения алгоритма СИПС1 до времени порядка $\Theta(n^2)$.
5. Разработайте вариант алгоритма СИПС1, который для расстояния Хемминга d_H решает задачу P1 за время порядка $\Theta(n)$.

Литература

- [1] *Aggarwal Alok, Park James*. Notes on searching in multidimensional monotone arrays // Proc. 29th Annual IEEE Symp. Foundations of Computer Science. — 1988. — P. 497–512.
- [2] *Aho Alfred V*. Algorithms for finding patterns in strings // Handbook of Theoretical Computer Science. J. van Leeuwen (ed.). Elsevier. — 1990. — P. 255–300.
- [3] *Aho Alfred V., Corasick Margaret J*. Efficient string matching: an aid to bibliographic search // *CACM*, 18–6, 1975. — P. 333–340.
- [4] *Aho Alfred V., Hirschberg Dan S., Ullman Jeffrey D*. Bounds on the complexity of the longest common subsequence problem // *J. ACM*, 23–1, 1976. — P. 1–12.
- [5] *Aho Alfred V., Hopcroft John E., Ullman Jeffrey D*. The Design & Analysis of Computer Algorithms // Addison-Wesley. — 1974. (Русский перевод: *Ахо А., Хоркрофт Дж., Ульман Дж.* Структуры данных и алгоритмы. — Издательский дом “Вильямс”, 2000.)
- [6] *Aho Alfred V., Sethi Ravi, Ullman Jeffrey D*. Compilers: Principles, Techniques & Tools // Addison-Wesley. — 1986. (Русский перевод: *Ахо А.В., Рэви С., Ульман Дж.Д.* Компиляторы: принципы, технологии и инструментарий. — Издательский дом “Вильямс”, 2001.)
- [7] *Allison Lloyd and Dix Trevor I*. A bit string longest subsequence algorithm // *IPL*, 23–6, 1986. — P. 305–310.
- [8] *Allouche Jean-Paul, Astoorian Dan, Randall Jim and Shallit Jeffrey*. Morphisms, square-free strings and the Tower of Hanoi puzzle // *Amer. Math. Monthly*, 101–7, 1994. — P. 651–658.
- [9] *Allouche Jean-Paul*. Sur la complexite des suites infinies // *Bull. Belg. Math. Soc.* 1. 1994. — P. 133–143.
- [10] *Amoux Pierre, Rauzy Gerard*. Representation geometriques de suites de complexite $2n + \ell$ // *Bull. Soc. Math. France*, 119, 1991. — P. 199–215.

- [11] *Andersson Arne, Larsson N. Jesper, Swanson Kurt.* Suffix trees on words // *Algorithmica*, 23–3, 1999. — P. 246–260.
- [12] *Apostolico Albano.* The myriad virtues of subword trees // *Combinatorial Algorithms on Words (NATO ASI Series F12)*. Albano Apostolico & Zvi Galil (eds.), Springer-Verlag, 1985. — P. 85–96.
- [13] *Apostolico Alberto and Crochemore Maxime.* Optimal canonization of all substrings of a string // *Information & Computation*, 95–1, 1991. — P. 76–95.
- [14] *Apostolico Alberto, Atallah Mikhail J., Larmore Lawrence and McFaddin Scott.* Efficient parallel algorithms for string editing problems // *SIAM J. Computing*, 19–5, 1990. — P. 968–988.
- [15] *Apostolico Alberto, Ehrenfeucht Andrzej.* Efficient detection of quasiperiodicities in strings // *TCS*, 119, 1993. — P. 247–265.
- [16] *Apostolico Alberto, Giancarlo Raffaele.* The Boyer-Moore-Galil string searching strategies revisited // *SIAM J. Comput.*, 15–1, 1986. — P. 98–105.
- [17] *Apostolico Alberto, Guerra C.* The longest common subsequence problem revisited // *Algorithmica*, 2, 1987. — P. 315–336.
- [18] *Apostolico Alberto, Preparata Franco P.* Optimal off-line detection of repetitions in a string // *TCS*, 22, 1983. — P. 297–315.
- [19] *Apostolico Alberto, Farach Martin, Iliopoulos Costas S.* Optimal superprimitivity testing for strings // *IPL*, 39, 1991. — P. 17–20.
- [20] *Arlazarov V. L., Dinic E. A., Kronrod M. A., Faradzev I. A.* On economical construction of the transitive closure of a directed graph // *Soviet Math. Dokl*, 11–5, 1970. — P. 1209–1210.
- [21] *Baeza-Yates Ricardo A.* Algorithms for string searching: a survey // *ACM SIGIR Forum*, 25-3/4, 1989. — P. 34–58.
- [22] *Baeza-Yates Ricardo A.* String searching algorithms revisited // *Algorithms & Data Structures. Lecture Notes in Computer Science 382*, Frank Dehne, Jörg-Rüdiger Sack & Nicola Santoro (eds.). Springer-Verlag, 1989. — P. 75–96.
- [23] *Baeza-Yates Ricardo A., Chris H. Perleberg.* Fast and practical approximate string matching // *Proc. Third Annual Symp. Combinatorial Pattern Matching*. Alberto Apostolico, Maxime Crochemore, Zvi Galil & Udi Manber (eds.) *Lecture Notes in Computer Science*, 644, Springer-Verlag, 1992. — P. 185–192.
- [24] *Baeza-Yates Ricardo A., Gaston Gonnet H.* A new approach to text searching // *CACM*, 35–70, 1992. — P. 74–82.
- [25] *Baeza-Yates Ricardo A., Gonnet Gaston H., Mireille Regnier.* Analysis of Boyer-Moore-type string searching algorithms // *Proc. First Annual ACM-SIAM Symp. Discrete Algs.*, 1990. — P. 328–343.

- [26] *Baeza-Yates Ricardo A., Gonzaio Navarro.* A faster algorithm for approximate string matching // Proc. Seventh Annual Symp. Combinatorial Pattern Matching. Dan S. Hirschberg & Gene W. Myers (eds.). Lecture Notes in Computer Science, 1075, Springer-Verlag, 1996. — P. 1–23.
- [27] *Baeza-Yates Ricardo A., Gonzaio Navarro.* Multiple approximate string matching // Proc. Fifth Annual Workshop on Algorithms & Data Structures. F. Dehne et al. (eds.), 1997. — P. 174–184.
- [28] *Baeza-Yates Ricardo A., Regnier Mireille.* Average running time of the Boyer-Moore-Horspool algorithm // TCS, 92, 1992. — P. 19–31.
- [29] *Bean Dwight R., Ehrenfeucht Andrzej, McNulty George F.* Avoidable patterns in strings of symbols // Pacific J. Math., 85–2, 1979. — P. 261–294.
- [30] *Benson Gary.* A space efficient algorithm for finding the best non-overlapping alignment score // Proc. Fifth Annual Symp. Combinatorial Pattern Matching. Maxime Crochemore & Dan Gusfield (eds.). Lecture Notes in Computer Science, 807, Springer-Verlag, 1994. — P. 1–14.
- [31] *Bentley Jon L., Sedgewick Robert.* Fast algorithms for sorting and searching strings // Proc. Eighth Annual ACM-SIAM Symp. Discrete Algs., 1997. — P. 360–369.
- [32] *Berstel Jean, Luc Boasson.* Partial words and a theorem of Fine and Wilf // TCS, 218, 1999. — P. 135–141.
- [33] *Berstel Jean, Seebold Patrice.* A characterization of Sturmian morphisms // The Mathematical Foundations of Computer Science. A. Borzyszkowski & S. Sokolowski (eds.), Springer-Verlag, 1993. — P. 281–290.
- [34] *Berstel Jean.* Recent results in Sturmian words // Developments in Language Theory. World Scientific Press, 1996. — P. 13–24.
- [35] *Blumer Anselm, Blumer Janet, Haussler David, Ehrenfeucht Andrzej, Chen M. T., Seiferas Joel.* The smallest automaton recognizing the subwords of a text // TCS, 40–1, 1985. — P. 31–55.
- [36] *Booth Kellogg S.* Lexicographically least circular substrings // IPL, 10–4/5, 1980. — P. 240–242.
- [37] *Boshemitzan M., Praenkel Aviezri S.* A linear algorithm for nonhomogeneous spectra of numbers // J. Algs., 5, 1984. — P. 187–198.
- [38] *Boyer Robert S., Moore J. Strother.* A fast string searching algorithm // CACM, 20–10, 1977. — P. 762–772.
- [39] *Braunholtz C. H.* Solution to Problem 5030: an infinite sequence of 3 symbols with no adjacent repeats // American Mathematical Monthly, 70, 1963. — P. 675–676.

- [40] *Breslauer Daily*. An on-line string superprimitivity test // IPL, 44, 1992. — P. 345–347.
- [41] *Brodal Gerth S., Lyngso Rune B., Pedersen Christian N. S., Stoye Jens*. Finding maximal pairs with bounded gap // J. Discrete Algs., 1, 2000. — P. 77–103.
- [42] *Brodal Gerth S., Pedersen Christian N. S.* Finding maximal quasiperiodicities in strings // Proc. Eleventh Annual Symp. Combinatorial Pattern Matching. Raffaele Giancarlo & David Sankoff (eds.). Lecture Notes in Computer Science, 1848, Springer-Verlag, 2000. — P. 397–411.
- [43] *Castelli M. Gabriella, Mignosi Filippo, Restivo Antonio*. Fine & Wilf’s theorem for three periods and a generalization of Sturmian words // TCS, 218. — 1999. — P. 83–94.
- [44] *Chang William I., Lampe Jordan*. Theoretical and empirical comparisons of approximate string matching algorithms // Proc. Third Annual Symp. Combinatorial Pattern Matching. Alberto Apostolico, Maxime Crochemore, Zvi Galil & Udi Manber (eds.). Lecture Notes in Computer Science, 644, Springer-Verlag, 1992. — P. 172–181.
- [45] *Chang William I., Lawler Eugene L.* Approximate string matching in sublinear expected time // Proc. 31st Annual IEEE Symp. Foundations of Computer Science. — Vol. I. — 1990. — P. 116–124.
- [46] *Chang William I., Lawler Eugene L.* Sublinear approximate string matching and biological applications // Algorithmica, 12–4/5, 1994. — P. 327–344.
- [47] *Charras Christian, Lecroq Thierry*. Exact String Matching Algorithms. — Laboratoire d’Informatique, Universite de Rouen, 1997: <http://www.igm.univ-mlv.fr/~lecroq/string/index.html>
- [48] *Chen K. T., Fox R. H., Lyndon Roger C.* Free differential calculus IV // Ann. Math., 68–1, 1958. — P. 81–95.
- [49] *Chen M. T., Seiferas Joel*. Efficient and elegant subword-tree construction // Combinatorial Algorithms on Words (NATO ASI Series F12). Alberto Apostolico & Zvi Galil (eds.). Springer-Verlag, 1985. — P. 97–107.
- [50] *Cole Richard, Hariharan Ramesh, Paterson Michael S., Zwick Uri*. Tighter lower bounds on the exact complexity of string matching // SIAM J. Comput., 24–1, 1995. — P. 30–45.
- [51] *Cole Richard, Hariharan Ramesh*. Faster approximate string matching // Proc. Ninth Annual ACM-SIAM Symp. Discrete Algs., 1998. — P. 463–472.
- [52] *Cole Richard, Hariharan Ramesh*. Tighter bounds on the exact complexity of string matching // Proc. 33rd Annual IEEE Symp. Foundations of Computer Science. 1992. — P. 600–609.

- [53] *Cole Richard*. Tight bounds on the complexity of the Boyer-Moore string matching algorithm // *SIAM J. Comput.*, 23–5, 1994. — P. 1075–1091.
- [54] *Colussi Livio, Galil Zvi, Giancarlo Raffaele*. On the exact complexity of string matching // *Proc. 31st Annual IEEE Symp. on Foundations of Computer Science*. — Vol I. — 1990. — P. 135–143.
- [55] *Colussi Livio*. Correctness and efficiency of pattern matching algorithms // *Information & Computation*, 95, 1991. — P. 225–251.
- [56] *Colussi Livio*. Fastest pattern matching in strings // *J. Algs.*, 16–2, 1994. — P. 163–189.
- [57] *Commentz-Waller Beate*. A string matching algorithm fast on the average // *Proc. Sixth Internat. Colloquium on Automata. Languages & Programming*. Hermann A. Maurer (ed.), *Lecture Notes in Computer Science*, 71, Springer-Verlag, 1979. — P. 118–132.
- [58] *Crawford Tim, Iliopoulos Costas S., Rajeev Raman*. String matching techniques for musical similarity and melodic recognition // *Computing in Musicology*, 11, 1998. — P. 73–100.
- [59] *Crochemore Maxime, Czumaj Artur, Gasieniec Leszek, Lecroq Thierry, Plandowski Wojciech, Rytter Wojciech*. Fast practical multi-pattern matching // *IPL*, 71–3/4, 1999. — P. 107–113.
- [60] *Crochemore Maxime, Hancart Christophe, Lecroq Thierry*. A unifying look at the Apostolico-Giancarlo string matching algorithm // *J. Discrete Algs.*, 2–1 (в печати).
- [61] *Crochemore Maxime, Hancart Christophe, Thierry Lecroq*. *Algorithmique du Texte*, Vuibert. — Paris. — 2001.
- [62] *Crochemore Maxime, Holub Jan*. On the implementation of compact DAWGs // *Proc. Seventh Internat. Conf. Implementation & Application of Automata*. — 2002. — P. 293–298.
- [63] *Crochemore Maxime, Iliopoulos Costas S., Pinzon Yoan J*. Speeding up Hirschberg and Hunt-Szymanski LCS algorithms // *Proc. Eighth IEEE Internat. Symp. String Processing & Information Retrieval*. — IEEE Computer Science Press, 2001. — P. 59–67.
- [64] *Crochemore Maxime, Iliopoulos Costas S., Pinzon Yoan J., Reid James F*. A fast and practical bit-vector algorithm for the longest common subsequence problem // *IPL*, 80–6, 2001. — P. 279–285.
- [65] *Crochemore Maxime, Lecroq Thierry*. Tight bounds on the complexity of the Apostolico-Giancarlo algorithm // *IPL*, 63–4, 1997. — P. 195–203.
- [66] *Crochemore Maxime, Perrin Dominique*. Two-way string matching // *JACM*, 38–3, 1991. — P. 651–675.

- [67] *Crochemore Maxime, Rytter Wojciech. Text Algorithms.* — Oxford University Press. 1994, 412 p.
- [68] *Crochemore Maxime, Verin Renaud.* Direct construction of compact directed acyclic word graphs // Proc. Eighth Annual Symp. Combinatorial Pattern Matching. Alberto Apostolico & Jotun Hein (eds.). Lecture Notes in Computer Science, 1264, Springer-Verlag, 1997. — P. 116–129.
- [69] *Crochemore Maxime.* An optimal algorithm for computing all the repetitions in a word // IPL, 12–5, 1981. — P. 244–248.
- [70] *Crochemore Maxime.* Linear searching for a square in a word // Bull. EATCS, 24, 1984. — P. 66–72.
- [71] *Crochemore Maxime.* Reducing space for index implementation // TCS (в печати).
- [72] *Crochemore Maxime.* Transducers and repetitions // TCS, 45–1, 1986. — P. 63–86.
- [73] *Crochemore Maxime, Czumaj Artur, Gasieniec Leszek, Jarominek Stefan, Lecroq Thierry, Plandowski Wojciech, Rytter Wojciech.* Speeding up two string-matching algorithms // Algoritmica, 12, 1994. — P. 247–267.
- [74] *Currie James D.* There are ternary circular square-free words of length n for $n \geq 18$ // *Electronic J. Combinatorics* (в печати).
- [75] *Currie James D.* Which graphs allow infinite nonrepetitive walks? // *Discrete Math.*, 87, 1991. — P. 249–260.
- [76] *Davies G., Bowsher S.* Algorithms for pattern matching // *Software — Practice & Experience*, 16–6, 1986. — P. 575–601.
- [77] *Domolki Balint.* A universal computer system based on production rules // BITS, 1968. — P. 262–275.
- [78] *Domolki Balint.* Algorithms for the recognition of properties of symbol strings (Russian) // *J. Computing Math. & Math. Physics*, 5–1, 1965.
- [79] *Duval Jean-Pierre, Lecroq Thierry, Arnaud Lefebvre.* Border array on bounded alphabet // Proc. Prague Stringology Conf. '02, 2002.
- [80] *Duval Jean-Pierre.* Factorizing words over an ordered alphabet // *J. Algs.*, 4, 1983. — P. 363–381.
- [81] *Ehrenfeucht Andrzej, Haussler David.* A new distance metric on strings computable in linear time // *Discrete Appl. Math.*, 20, 1988. — P. 191–203.
- [82] *El-Mabrouk Nadia, Maxime Crochemore.* Boyer-Moore strategy to efficient approximate string matching // Proc. Seventh Annual Symp. Combinatorial Pattern Matching. Dan S. Hirschberg & Gene W. Myers (eds.), Lecture Notes in Computer Science, 1075, Springer-Verlag, 1996. — P. 24–28.

- [83] *Erdos Pal*. Some unsolved problems // Magyar Tud. Akad. Mat. Kutato Intez.et Kozl., 6, 1961. — P. 221–254.
- [84] *Erickson Bruce W., Sellers Peter H.* Recognition of patterns in genetic sequences // Time Warps, String Edits & Macromolecules: The Theory & Practice of Sequence Comparison. David Sankoff & Joseph B. Kruskal (eds.), Addison-Wesley, 1983. — P. 55–91.
- [85] *Farach Martin, Muthukrishnan S.* Optimal logarithmic time randomized suffix tree construction // Proc. 23rd Internat. Colloquium Automata, Languages & Programming, 1996.
- [86] *Farach Martin*. Optimal suffix tree construction with large alphabets // Proc. 38th Annual IEEE Symp. Foundations of Computer Science, 1997. — P. 137–143.
- [87] *Fine N. J., Wilf Herbert S.* Uniqueness theorems for periodic functions // Proc. Amer. Math. Soc., 16, 1965. — P. 109–114.
- [88] *Fraenkel Aviezri S., Simpson R. Jamie*. How many squares can a string contain? // J. Combinatorial Theory (A), 82, 1998. — P. 112–120.
- [89] *Fraenkel Aviezri S., Simpson R. Jamie*. The exact number of squares in Fibonacci words // TCS, 218–1, 1999. — P. 83–94.
- [90] *Franek Frantisek, Gao Shudi, Lu Weilin, Ryan Patrick J., Smyth W. F., Sun Yu, Yang Lu*. Verifying a Border Array in Linear Time // J. Combinatorial Math. & Combinatorial Comput., 42, 2002.
- [91] *Franek Frantisek, Jiang Jiandong, Lu Weilin, Smyth W. F.* Two-pattern strings // Proc. Thirteenth Annual Symposium on Combinatorial Pattern Matching. Alberto Apostolico & Masayuki Takeda (eds.). Lecture Notes in Computer Science, Springer-Verlag, 2002. — P. 76–84.
- [92] *Franek Frantisek, Karaman Ayse, Smyth W. F.* Repetitions in Sturmian strings // TCS, 249, 2000. — P. 289–303.
- [93] *Franek Frantisek, Smyth W. F., Tang Yudong*. Computing all repeats using suffix arrays // Proc. Thirteenth Australasian Workshop on Combinatorial Algorithms. Elizabeth Billington, Diane Donovan & Abdollah Khodkar (eds.), 2002. — P. 171–184.
- [94] *Franek Frantisek, Smyth W. F., Xiao Xiangdong*. A note on Crochemore’s repetitions algorithm — a fast space-efficient approach // Proc. Prague Stringology Conf. ’02. M. Balik & M. Simanek (eds.), 2002. — P. 36–43.
- [95] *Galil Zvi, Giancarlo Raffaele*. Improved string matching with k mismatches // SIGACT News, 17, 1986. — P. 52–54.
- [96] *Galil Zvi, Giancarlo Raffaele*. On the exact complexity of string matching: lower bounds // SIAM J. Comput., 20–6, 1991. — P. 1008–1020.

- [97] *Galil Zvi, Giancarlo Raffaele*. On the exact complexity of string matching: upper bounds // *SIAM J. Comput.*, 21–3, 1992. — P. 407–437.
- [98] *Galil Zvi, Park Kunsoo*. An improved algorithm for approximate string matching // *SIAM J. Comput.* 19–8, 1990. — P. 989–999.
- [99] *Galil Zvi, Seiferas Joel*. Time-space optimal string matching // *J. Computer & System Sci.* 26–3, 1983. — P. 280–294.
- [100] *Galil Zvi*. On improving the worst case running time of the Boyer-Moore string matching algorithm // *CACM*, 22–9, 1979. — P. 505–508.
- [101] *Gao Shudi*. New Properties of Borders & Covers of Strings // M. Sc. thesis, Department of Computing & Software, McMaster University, 2001.
- [102] *Giegerich Robert, Kurtz Stefan*. A comparison of imperative and purely functional suffix tree constructions // *Science of Computer Programming*, 25-2/3, 1995. — P. 187–218.
- [103] *Gonnet Gaston H., Baeza-Yates Ricardo A.* Handbook of Algorithms & Data Structures in Pascal and C, 2nd edition // Addison–Wesley, 1991.
- [104] *Guibas Leo J., Odlyzko Andrew M.* A new proof of the linearity of the Boyer-Moore string searching algorithm // *SIAM J. Comput.* 9–4, 1980. — P. 672–682.
- [105] *Guibas Leo J., Odlyzko Andrew M.* Periods in strings // *J. Combinatorial Theory (A)* 30, 1981. — P. 19–42.
- [106] *Guibas Leo J., Odlyzko Andrew M.* String overlaps, pattern matching and non-transitive games // *J. Combinatorial Theory (A)* 30, 1981. — P. 183–208.
- [107] *Gusfield Dan, Stoye Jens*. Linear Time Algorithms for Finding & Representing all the Tandem Repeats in a String // *Tech. Rep. CSE-98-4*, Computer Science Department, University of California, Davis, 1998.
- [108] *Gusfield Dan*. Algorithms on Strings, Trees, & Sequences. — Cambridge University Press, 1997, 534 p.
- [109] *Hamming Richard*. Coding & Information Theory // Prentice Hall, 1982.
- [110] *Hancart Christophe*. Analyse Exacte et Moyenne d’Algorithmes de Recherche d’un Motif dans un Texte // Ph. D. thesis, Université de Paris, 1993, 96 p.
- [111] *Hannan Sampath K., Myers Gene W.* An algorithm for locating nonoverlapping regions of maximum alignment score // *Proc. Fourth Annual Symp. Combinatorial Pattern Matching*. Alberto Apostolico, Maxime Crochemore, Zvi Galil & Udi Manber (eds.). Lecture Notes in Computer Science 684, Springer-Verlag, 1993. — P. 74–86.
- [112] *Harel Dov, Tarjan Robert Endre*. Fast algorithms for finding nearest common ancestors // *SIAM J. Computing* 13-2, 1984. — P. 338–355.

- [113] *Haton Jean-Paul*. Contribution a l’Analyse, Parametrisation et la Reconnaissance Automatique de la Parole // Doctoral thesis, Universite de Nancy, 1973.
- [114] *Hirschberg Dan S*. A linear space algorithm for computing maximal common subsequences // CACM 18-6, 1975. — P. 341–343.
- [115] *Hirschberg Dan S*. Algorithms for the longest common subsequence problem // JACM 24-4, 1977. — P. 664–675.
- [116] *Hirschberg Dan S*. An information-theoretic lower bound for the longest common subsequence problem // IPL 7-1, 1978. — P. 40–42.
- [117] *Hopcroft John E*. An $n \log n$ algorithm for minimizing states in a finite automaton // Theory of Machines & Computations, Z. Kohavi & A. Paz (eds.). Academic Press, 1971.
- [118] *Hopcroft John E., Ullman Jeffrey D*. Introduction to Automata Theory. Languages, & Computation // Addison–Wesley, 1979.
- [119] *Horspool R. Nigel*. Practical fast searching in strings // Software — Practice & Experience 70-6, 1980. — P. 501–506.
- [120] *Huang Xiaoqiu*. A lower bound for the edit-distance problem under arbitrary cost function // IPL 27-6, 1988. — P. 319–321.
- [121] *Hume Andrew, Sunday Daniel*. Fast string searching // Software — Practice & Experience 21-11, 1991. — P. 1221–1248.
- [122] *Hunt James W., Szymanski Thomas G*. A fast algorithm for computing longest common subsequences // CACM 20-5, 1977. — P. 350–353.
- [123] *Iliopoulos Costas S., Moore Dennis, Smyth W. F*. A characterization of the squares in a Fibonacci string // TCS 172, 1997. — P. 281–291.
- [124] *Iliopoulos Costas S., Mouchard Laurent*. Fast local covers // Препринт.
- [125] *Iliopoulos Costas S., Pinzon Yoan J*. Recovering an LCS in $O(n^2/w)$ time and space // Proc. Twelfth Australasian Workshop on Combinatorial Algorithms. Edi T, Baskoro (ed.), Institute of Technology Bandung, 2001. — P. 106–117.
- [126] *Iliopoulos Costas S., Smyth W. F*. A fast average case algorithm for Lyndon decomposition // Internal. 3. Computer Math. 57-3/4, 1995. — P. 15–31.
- [127] *Inenaga S., Hoshino H., Shmohara A., Takeda Masayuki, Arikawa S., Mauri G., Pavesi G*. On-line construction of compact directed acyclic word graphs // Proc. Twelfth Annual Symp. Combinatorial Pattern Matching. Amihood Amir & Gad M. Landau (eds.) Lecture Notes in Computer Science 2087, Springer–Verlag, 2001. — P. 169–180.
- [128] *Jacobson Guy, Vo Kiem-Phong*. Heaviest increasing/common subsequence problems // Proc. Third Annual Symp. Combinatorial Pattern Matching. Alberto Apostolico, Maxime Crochemore, Zvi Galil & Udi Manber (eds.) Lecture Notes in Computer Science 644, Springer–Verlag, 1992. — P. 52–66.

- [129] *Jennings Christopher G.* A Linear-Time Algorithm for Fast Exact Pattern Matching in Strings // M. Sc. thesis. Department of Computing & Software, McMaster University, 2002.
- [130] *Karp Richard M.* Reducibility among combinatorial problems // Complexity of Computer Computations. Raymond E. Miller & James W. Thatcher (eds.), IBM Research Symposia Series, 1972. — P. 85–103.
- [131] *Karp Richard M., Rabin Michael O.* Efficient randomized pattern-matching algorithms // IBM J. Res. Develop. 31-2, 1987. — P. 249–260.
- [132] *Keranen Veikko.* Abelian squares are avoidable on 4 letters // Automata, Languages & Programming. Lecture Notes in Computer Science 623, Springer-Verlag, 1992. — P. 41–52.
- [133] *Kim Sung-Ryul, Park Kunsoo.* A dynamic edit distance table // Proc. Eleventh Annual Symp. Combinatorial Pattern Matching. Raffaele Giancarlo & David Sankoff (eds.), Lecture Notes in Computer Science 1848, Springer-Verlag, 2000. — P. 60–68.
- [134] *Knuth Donald E.* Big omichron, big omega, and big theta // ACM SIGACT News. April, 1976. — P. 18–24.
- [135] *Knuth Donald E.* The Art of Computer Programming I: Fundamental Algorithms, 2nd edition // Addison–Wesley, 1973, 722 p. (Русский перевод 3-го издания: *Кнут Д.Э.* Искусство программирования. Т.1. Основные алгоритмы. — Издательский дом “Вильямс”, 2000.)
- [136] *Knuth Donald E.* The Art of Computer Programming III: Sorting & Searching // Addison-Wesley, 1973, 634 p. (Русский перевод 2-го издания: *Кнут Д.Э.* Искусство программирования. Т.3. Сортировка и поиск. — Издательский дом “Вильямс”, 2000.)
- [137] *Knuth Donald E., Morris James H., Pratt Vaughan R.* Fast pattern matching in strings // SIAM J. Comput. 6-2, 1977. — P. 323–350.
- [138] *Kolpakov Roman, Kucherov Gregory.* Maximal Repetitions in Words, or How to Find All Squares in Linear Time // Rapport Inteme LORIA 98-R-227, Laboratoire Lorrain de Recherche en Informatique et ses Applications, France, 1998.
- [139] *Kolpakov Roman, Kucherov Gregory.* On maximal repetitions in words // J. Discrete Algs. 7, 2000. — P. 159–186.
- [140] *Kolpakov Roman, Kucherov Gregory.* On the Sum of Exponents of Maximal Repetitions in a Word // Rapport Inteme LORIA 99-R-034, Laboratoire Lorrain de Recherche en Informatique et ses Applications, France, 1999.
- [141] *Kosaraju S. Rao.* Computation of squares in a string // Proc. Fifth Annual Symp. Combinatorial Pattern Matching. Maxime Crochemore & Dan Gusfield (eds.). Lecture Notes in Computer Science 807, Springer-Verlag, 1994. — P. 146–150.

- [142] *Kowalski G., Meltzer A.* New multi-term high speed text search algorithms // Proc. First IEEE Internat. Conf. on Computer Applications, 1984. — P. 514–522.
- [143] *Kurtz Stefan.* Reducing the space requirement of suffix trees // Software — Practice & Experience 29-13, 1999. — P. 1149–1171.
- [144] *Landau Gad M., Myers Gene W., Schmidt Jeanette P.* Incremental string comparison // SIAMJ. Comput. 27-2, 1998. — P. 557–582.
- [145] *Landau Gad M., Schmidt Jeanette P.* An algorithm for approximate tandem repeats // Proc. Fourth Annual Symp. Combinatorial Pattern Matching. Alberto Apostolico, Maxime Crochemore, Zvi Galil & Udi Manber (eds.), Lecture Notes in Computer Science 684, Springer-Verlag, 1993. — P. 120–133.
- [146] *Landau Gad M., Vishkin Uzi.* Efficient string matching in the presence of errors // Proc. 26th Annual IEEE Symp. Foundations of Computer Science, 1985. — P. 126–136.
- [147] *Landau Gad M., Vishkin Uzi.* Efficient string matching with k mismatches // TCS 43, 1986. — P. 239–249.
- [148] *Landau Gad M., Vishkin Uzi.* Fast parallel and serial approximate string matching // J.Algs. 10, 1989. — P. 157–169.
- [149] *Landau Gad M., Vishkin Uzi.* Fast string matching with k differences // J. Computer & System Sci. 37, 1988. — P. 63–78.
- [150] *Landau Gad M., Vishkin Uzi.* Introducing efficient parallelism into approximate string matching and a new serial algorithm // Proc. Eighteenth ACM Symp. Theory of Comput., 1986. — P. 220–230.
- [151] *Larsson N. Jesper.* Extended application of suffix trees to data compression // Proc. IEEE Data Compression Conference, 1996. — P. 190–199.
- [152] *Lecroq Thierry.* A variation on the Boyer–Moore algorithm // TCS 92, 1992. — P. 119–144.
- [153] *Lecroq Thierry.* Experimental results on string matching algorithms // Software — Practice & Experience 25-7, 1995. — P. 727–765.
- [154] *Lecroq Thierry.* New Experimental Results on Exact String-Matching // Rapport LIFAR 2000.03, Universite de Rouen, 2000.
- [155] *Leech John.* A problem on strings of beads // Math. Gazette 4J, 1957. — P. 277–278.
- [156] *Lempel Abraham, Ziv Jacob.* On the complexity of finite sequences // IEEE Trans. Information Theory 22, 1976. — P. 75–81.
- [157] *Lentin Andre, Schutzenberger Marcel P.* A combinatorial problem in the theory of free monoids // Combinatorial Mathematics & Its Applications. R. C. Bose & T. A. Dowling (eds.), University of North Carolina Press, 1969. — P. 128–144.

- [158] *Levenshtein V. I.* Binary codes capable of correcting deletions, insertions and reversals // *Cybernetics & Control Theory* 10-8, 1966. — P. 707–710.
- [159] *Lewis Harry R., Papadimitriou Christos H.* Elements of the Theory of Computation, 2nd edition // Prentice-Hall, 1998.
- [160] *Li Yin, Smyth W. F.* Computing the cover array in linear time // *Algorithmica* 32-1, 2001. — P. 95–106.
- [161] *Li Yin.* A Windows-Based String Algorithm Testing System // Undergraduate Computer Science Project, Department of Comp. Sci. & Systems, McMaster University, 1996.
- [162] *Lothaire M.* (ed.). Algebraic Combinatorics on Words // Encyclopedia of Mathematics and its Applications. — Vol.90. — Cambridge University Press. — 2002. См. также <http://www-igm.univ-mlv.fr/~\,berstel/>
- [163] *Lothaire M.* (ed.). Combinatorics on Words, 2nd edition // Cambridge University Press, 1997.
- [164] *Lothaire M.* (ed.). Combinatorics on Words // Addison–Wesley, 1983.
- [165] *Lowrance Roy, Wagner Robert A.* An extension of the string-to-string correction problem // *JACM* 22-2, 1975. — P. 177–183.
- [166] *Luca de Aldo, Filippo Mignosi.* Some combinatorial properties of Sturmian words // *TCS* 136, 1994. — P. 361–385.
- [167] *Lyndon Roger C., Schützenberger Marcel P.* The equation $a^M = b^N c^P$ in a free group // *Michigan Math. J.* 9, 1962. — P. 289–298.
- [168] *Main Michael G.* Detecting leftmost maximal periodicities // *Discrete Applied Maths.* 25, 1989. — P. 145–153.
- [169] *Main Michael G., Lorentz Richard J.* An $O(n \log n)$ Algorithm for Recognizing Repetition // Tech. Rep. CS-79-056, Computer Science Department, Washington State University, 1979.
- [170] *Main Michael G., Lorentz Richard J.* Linear time recognition of squarefree strings // *Combinatorial Algorithms on Words (NATO ASI Series F12)*. Alberto Apostolico & Zvi Galil (eds.), Springer–Verlag, 1985. — P. 271–278.
- [171] *Manber Udi, Myers Gene W.* Suffix arrays: a new method for on-line string searches // *Proc. First Annual ACM-SIAM Symp. Discrete Algs.*, 1990. — P. 319–327.
- [172] *Manber Udi, Myers Gene W.* Suffix arrays: a new method for on-line string searches // *SIAM J. Comput.* 22-5, 1993. — P. 935–948.
- [173] *Masek William J., Paterson Michael S.* A faster algorithm for computing string-edit distances // *J. Computer & System Sci.* 20-1, 1980. — P. 18–31.

- [174] *McCreight Edward M.* A space-economical suffix tree construction algorithm // JACM 23-2, 1976. — P. 262–272.
- [175] *McNaughton Robert F., Yamada H.* Regular expressions and state graphs for automata // IRE Trans. Electronic Computers EC9-1, 1960. — P. 39–47.
- [176] *Michael G. Main & Richard J. Lorentz.* An $O(n \log n)$ algorithm for finding all repetitions in a string // J. Algs., 5, 1984. — P. 422–432.
- [177] *Michailidis Panagiotis D., Margaritis Konstantinos G.* On-line string matching algorithms: survey and experimental results // Internal. J. Computer Math. 76-4, 2001. — P. 411–434.
- [178] *Mignosi Filippo, Seebold Patrice.* Morphismes sturmiens et regles de Rauzy // J. Theorie de Nombres de Bordeaux 5, 1993. — P. 221–233.
- [179] *Mignosi Filippo.* Infinite words with linear subword complexity // TCS 65, 1989. — P. 221–242.
- [180] *Mignosi Filippo.* On the number of factors of Sturmian words // TCS 82, 1989. — P. 221–242.
- [181] *Moore Dennis, Smyth W. F.* A correction to: An optimal algorithm to compute all the covers of a string // IPL 54, 1995. — P. 101–103.
- [182] *Moore Dennis, Smyth W. F.* An optimal algorithm to compute all the covers of a string // IPL 50-5, 1994. — P. 239–246.
- [183] *Moore Dennis, Smyth W. F., Miller Dianne.* Counting Distinct Strings // Aig-orthmica 13-1, 1999. — P. 1–13.
- [184] *Morris James H., Pratt Vaughan R.* A Linear Pattern-Matching Algorithm // Tech. Rep. 40, University of California, Berkeley, 1970.
- [185] *Morrison Donald R.* PATRICIA — practical algorithm to retrieve information coded in alphanumeric // JACM 15-4, 1968. — P. 514–534.
- [186] *Muth Robert, Manber Udi.* Approximate multiple string search // Proc. Seventh Annual Symp. Combinatorial Pattern Matching. Dan S. Hirschberg & Gene W. Myers (eds.), Lecture Notes in Computer Science 1075, Springer-Verlag, 1996. — P. 75–86.
- [187] *Myers Gene W.* A fast bit-vector algorithm for approximate string matching based on dynamic programming // JACM 46-3, 1999. — P. 395–415.
- [188] *Myers Gene W.* A Four Russians Algorithm for Regular Expression Pattern Matching // Technical Report 88-34, Dept. of Computer Science, University of Arizona, 1988.
- [189] *Myers Gene W.* An $O(ND)$ difference algorithm and its variations // Algorith-mica I, 1986. — P. 251–266.

- [190] *Nagy Benedek*, частное сообщение.
- [191] *Navarro Gonzalo*. A guided tour to approximate string matching // ACM Computing Surveys 33-1, 2001. — P. 31–88.
- [192] *Needleman Saul B., Wunsch Christian D.* A general method applicable to the search for similarities in the amino-acid principle of two proteins // J. Molecular Biology 48, 1970. — P. 443–453.
- [193] *Parikh Rohit J.* On context-free languages // JACM 13-4, 1966. — P. 570–581.
- [194] *Pleasant Peter A. B.* Non-repetitive sequences // Proc. Camb. Phil. Soc. 68, 1970. — P. 267–274.
- [195] *Reichert Thomas A., Cohen Donald N., Wong Andrew K. C.* An application of information theory to genetic mutations and the matching of polypeptide sequences // J. Theoretical Biology 42, 1973. — P. 245–261.
- [196] *Rivest Ronald L.* On the worst-case behaviour of string-searching algorithms // SIAMJ. Comput. 6-4, 1977. — P. 669–674.
- [197] *Rote Gunter.* Sequences with subword complexity $2n$ // J. Number Theory 46-2, 1994. — P. 196–213.
- [198] *Ryan Patrick J., Smyth W. F.* An approach to asymptotic complexity // Mathematics & Computer Education 26-2, 1992. — P. 135–146.
- [199] *Rytter Wojciech.* A correct preprocessing algorithm for Boyer-Moore string searching // SIAMJ. Comput. 9-3, 1980. — P. 509–512.
- [200] *Sadakane Kunihiko.* A fast algorithm for making suffix arrays and for Burrows-Wheeler transformation // Proc. IEEE Data Compression Conference, 1998. — P. 129–138.
- [201] *Sakoe Hiroaki, Chiba Seibi.* A dynamic programming approach to continuous speech recognition // Proc. Internal. Congress of Acoustics, Budapest, 1971, paper 20C13.
- [202] *Sankoff David, Kruskal Joseph B.* (eds.). Time Warps, Siring Edits, & Macromolecules: the Theory & Practice of Sequence Comparison // Addison-Wesley, 1983.
- [203] *Sankoff David.* Matching sequences under deletion-insertion constraints // Proc. USA National Acad. Sci., 69, 1972. — P. 4–6.
- [204] *Schensted Craige.* Largest increasing and decreasing subsequences // Canadian J. Math., 73, 1961. — P. 179–191.
- [205] *Schmidt Jeanette P.* All highest scoring paths in weighted grid graphs and their application to finding all approximate repeats in strings // SIAM J. Comput. 27-4, 1998. — P. 972–992.

- [206] *Setubal Joao, Meidanis Joao*. Introduction to Computational Molecular Biology // PWS Publishing, 1997, 296 p.
- [207] *Shiloach Yossi*. A fast equivalence-checking algorithm for circular lists // IPL 8-5, 1979. — P. 236–238.
- [208] *Shiloach Yossi*. Fast canonization of circular strings // J. Algs., 2, 1981. — P. 107–121.
- [209] *Sim Jeong Seop, Iliopoulos Costas S., Park Kunsoo, Smyth W. F.* Approximate periods of strings // TCS 262, 2001. — P. 557–568.
- [210] *Simpson R. Jamie*. Covers and the periodicity lemma (частное сообщение).
- [211] *Sloane N. J. A., Plouffe Simon*. The Encyclopedia of Integer Sequences // Academic Press, 1995. См. также <http://www.research.att.com/~njas/sequences/>
- [212] *Smit G. de V.* A comparison of three string matching algorithms // Software — Practice & Experience 12, 1982. — P. 57–66.
- [213] *Smith Temple F., Waterman Michael S.* Identification of common molecular subsequences // J. Molecular Biology 747, 1981. — P. 195–197.
- [214] *Stephen Graham A.* String Searching Algorithms // World Scientific Publishing, 1994, 243 p.
- [215] *Stolarsky Kenneth B.* Beatty sequences, continued fractions, and certain shift operators // Canad. Math. Bull. 19-4, 1976. — P. 473–481.
- [216] *Stoye Jens, Gusfield Dan*. Simple and flexible detection of contiguous repeats using a suffix tree // Proc. Ninth Annual Symp. Combinatorial Pattern Matching. Martin Farach-Colton (ed.), Lecture Notes in Computer Science 1448, Springer-Verlag, 1998. — P. 140–152.
- [217] *Sunday Daniel M.* A very fast substring search algorithm // CACM 33-8, 1990. — P. 132–142.
- [218] *Tarhio Jorma, Ukkonen Esko*. Approximate Boyer–Moore string matching // SIAM J. Comput. 22-2, 1993. — P. 243–260.
- [219] *Thompson Ken*. Regular expression search algorithm // CACM 11-6, 1968. — P. 419–422.
- [220] *Thue Axel*. Uber unendliche zeichenreihen // Norske Vid. Selsk. Skr. I. Mat. Nat. Kl. Christiana 7, 1906. — P. 1–22.
- [221] *Tolkien J. R. R.* The Lord of the Rings, Part II: The Two Towers // Houghton Mifflin, 1955.
- [222] *Ukkonen Esko, Wood Derick*. Approximate string matching with suffix automata // Algorithmica 10, 1993. — P. 353–364.

- [223] *Ukkonen Esko*. Algorithms for approximate string matching // Information & Control, 64, 1985. — P. 100–118.
- [224] *Ukkonen Esko*. Approximate string-matching with q -grams and maximal matches // TCS 92, 1992. — P. 191–211.
- [225] *Ukkonen Esko*. Constructing suffix trees on-line in linear time // Proc. IFIP 92. — Vol. I. — 1992. — P. 484–492.
- [226] *Ukkonen Esko*. Finding approximate patterns in strings // J.Algs. 6, 1985. — P. 132–137.
- [227] *Velichko V. M., Zagoruyko N. G.* Automatic recognition of 200 words // Internal. J. Man-Machine Studies 2, 1970. — P. 223–234.
- [228] *Vintsyuk T. K.* Speech discrimination by dynamic programming // Cybernetics 4-1, 1968. — P. 52–57.
- [229] *Wagner Robert A., Fischer Michael J.* The string-to-string correction problem // JACM 21-1, 1974. — P. 168–173.
- [230] *Weiner Peter*. Linear pattern matching algorithms // Proc. 14th Annual IEEE Symp. Switching & Automata Theory, 1973. — P. 1–11.
- [231] *Wright Alden H.* Approximate string matching using within-word parallelism // Software — Practice & Experience 24, 1994. — P. 337–362.
- [232] *Wu Sun, Manber Udi, Myers Gene W.* A subquadratic algorithm for approximate limited expression matching // Algorithmica 15, 1996. — P. 50–67.
- [233] *Wu Sun, Manber Udi*. Fast text searching allowing errors // CACM 35-10, 1992. — P. 83–91.
- [234] *Yao Andrew Chi-Chih*. The complexity of pattern matching for a random string // SIAM J. Comput. 8-3, 1979. — P. 368–387.
- [235] *Ziv Jacob, Lempel Abraham*. A universal algorithm for sequential data compression // IEEE Trans. Information Theory 23, 1977. — P. 337–343.

ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ

- С**
CFL, 192
- L**
LCS, 69, 285, 310
- M**
MSP, 51
- S**
s-факторизация, 63, 211, 406, 407
- A**
Алгоритм
ACFA, 361
BMFAST, 246
BYN, 369, 373
CWFA, 364
N DFA, 351
Бойера–Мура (БМ), 224, 319, 360
 применение, 230
Бойера–Мура–Гелила (БМГ), 258, 259
Бойера–Мура–Гелила–Смита (БМГС),
 258, 259
Бойера–Мура–Санди
 второй (БМС2), 254
 первый (БМС1), 252
Бойера–Мура–Хоспула (БМХ), 248
Вагнера–Фишера (ВФ), 283, 307, 310
 модифицированный, 313
Ву–Менбера (ВМ), 334
 модифицированный, 355, 367
- вычисления
 NE-дерева, 433
 NE-массива, 438
 максимальных суффиксов, 208
 массива граней, 35
 минимальных суффиксов, 205
Демелки–Бейза–Ятса–Гоннета (ДБГ),
 240, 242, 309
Дюваля
 второй, 198
 первый, 196
Зива–Лемпела, 213
Карпа–Рабина (КР), 236, 238
Кнута–Морриса–Пратта (КМП), 217,
 219, 266, 360
 модифицированный, 221
Кнута–Морриса–Пратта–Ханката
 (КМПХ), 266, 268
Колпакова–Кучерова (КК), 406, 413, 414
Крочемора, 386, 390
Ландау–Вишкина, 318
Ли–Смита (ЛС), 418, 427
Майерса, 325
 модифицированный (ММ), 333
Мейна, 407, 409
Мейна–Лоренца (МЛ), 397
 модифицированный, 456
онлайновый, 39
ПДБГ, 316
Сима–Илиопулоса–Парка–Смита

- второй (СИПС2), 464
- первый (СИПС1), 462
- Турбо-БМ, 263, 265
- Укконена, 322
- Укконена–Майерса (УМ), 299, 305, 307, 310, 322
- Франека–Смита–Танга (ФСТ), 430
- Ханта–Шиманского (ХШ), 292, 296, 307, 310
- Хешберга, 310
 - Н1, 287
 - Н2, 288
 - Н3, 290, 307
- Шмидта, 441, 457
- Алгоритмы
 - временная сложность, 114
 - динамического программирования, 280
 - для регулярных выражений, 346
 - корректность, 114
 - локализации паттернов,
 - классификация, 272
 - независимость от алфавита, 116
 - онлайновость, 115
 - оценки сложности, 274
 - приближенного сравнения с паттерном, 341
 - пространственная сложность, 114
 - сравнения с множественными паттернами, 359
 - характеристики, 113
- Алфавит, 23
 - бинарный, 23
 - индексированный, 117
 - кватернарный, 23
 - неупорядоченный, 116
 - стандартный, 121
 - тернарный, 23
 - упорядоченный, 116
- В**
- Выражения регулярные, 70
- Г**
- Грань, 27
 - наибольшая, 27
 - определение, 27
- Граф
 - зависимостей, 300
 - сеточный, 442, 447
- Д**
- Декомпозиция, 27
 - линдонская, 50, 63
- Дерево
 - NE, 432
 - NRE, 431
 - Patricia, 58
 - граней, 57, 420
 - компактное синтаксическое, 58
 - оболочек, 418, 420
 - периодов, 430
 - подслов, 59
 - синтаксическое, 58, 348
 - суффиксов, 59, 345, 396
- З**
- Задача
 - k*-несовпадений, 314
 - k*-разностей, 320
- И**
- Исключения, 86
 - слабые, 86
- К**
- Каноническая форма строковых петель, 48
- Кодирование, 118
- Конечный автомат, 344
 - ACFA, 360
 - BYNFA, 369
 - CWFA, 364
 - DFA, 345, 352
 - NDFFA, 345, 346
 - Ахо–Корасика, 360
 - детерминированный, 345
 - Комменца–Вальтера, 364
 - недетерминированный, 345
 - суффиксный, 345
- Крочемора тройки, 76
- Л**
- Лакуна, 78
- Лемма периодичности, 40
- Линдонская декомпозиция, 50, 192

алгоритм Дюваля, 192
 применение, 203
 Линдонское слово, 48

М

Массив
 DIST, 454
 NE, 435
 граней, 57, 127
 Монжа, 453
 оболочек, 57, 418, 419, 426
 периодов, 429
 полезности, 441, 444, 446
 стоимостей, 280, 325
 Массив граней, 33, 57
 правый, 43
 Матрица
 баллов, 445
 весов, 67, 460
 Монжа, 454
 переходная, 354
 разностей, 325
 Метасимволы, 70
 Множество регулярное, 72
 Морфизм, 87
 код, 88
 обратный, 88
 тождественный, 87

Н

Надстрока
 определение, 26
 собственная, 26
 Нормальная форма, 29, 57

О

Образующая, 28
 приближенная, 461
 Оптимальность локальная, 442, 450

П

Палиндром, 96
 Паттерны
 внутренние, 56
 множественные, 72
 характеристические, 75
 частные, 64

Период, 28
 минимальный, 29
 Подпоследовательность, 69
 наибольшая общая, 69
 Подстрока
 определение, 26
 собственная, 26
 Порядок, 28
 лексикографический, 26
 максимальный, 29
 Префикс, 27
 собственный, 27

Р

Раппорт, 77, 418
 NLE, 430
 NRE, 430
 аппроксимирующий, 82
 непродолжаемый (NE), 80
 оболочка, 78, 419
 покрытие, 78
 полный NE, 430
 последовательный, 78
 приближенный, 441
 продолжаемый
 влево (LE), 80
 вправо (RE), 80
 расщепление, 78
 смешанный, 78
 Расстояние
 абсолютное, 460
 взвешенное, 68, 280, 460
 Левенштейна, 66, 280, 311
 относительное, 460
 преобразования, 67, 280
 преобразования асимметричное, 74
 удаления, 74
 Хемминга, 66, 280, 421

С

Серия, 80, 406
 Сигнатура, 236
 Символы замещения, 70, 338
 Сравнение со стражем, 250
 Строка
b-каноническая, 128

- p -каноническая, 121
 p -свободная, 85
 бесконечная, 21
 исключаемая паттерном, 85
 квадрат, 29
 кратная, 29
 куб, 29
 линейная, 21, 22, 24
 непродолжаемая, 92
 нормальная форма, 29
 определение, 20
 определенная на алфавите, 23
 палиндром, 32
 периодическая, 29
 примитивная, 29
 пустая, 22, 24
 сильно периодическая, 29
 слабо кратная, 84
 слабо периодическая, 29
 сопряженная, 90
 циклическая, 21
- Строки
 b -эквивалентные, 128
 p -эквивалентные, 120
 лексикографический порядок, 26
 сравнение, 26
 Tue, 88, 89, 97
- Фибоначчи, 88, 100, 386
 бесконечные, 101
 обобщенные, 103
 обратные, 111
 Штурма, 110
 Строковая петля, 21, 45
 бесконечная, 21
 каноническая форма, 48
 некратная, 54
 подстрока, 46
 Строковая последовательность, *см.* Строка
 Суффикс, 27
 собственный, 27
- Ф**
- Факторизация, *см.* Декомпозиция
- Ц**
- Циклический сдвиг, 45
- Ч**
- Числа
 Белла, 123
 Стирлинга, 123
- Э**
- Эквивалентность относительно паттерна,
 120

Научно-популярное издание

Билл Смит

Методы и алгоритмы вычислений на строках

Литературный редактор *Т.П. Кайгородова*

Верстка *А.Н. Полинчик*

Художественные редакторы *В.Г. Павлютин,*

Т.А. Тараброва

Корректоры *Т.А. Корзун,*

В.В. Смоляр

Издательский дом “Вильямс”

101509, г. Москва, ул. Лесная, д. 43, стр. 1

Подписано в печать 21.08.2006. Формат 70×100/16.

Гарнитура Times. Печать офсетная.

Усл. печ. л. 40. Уч.-изд. л. 26,5.

Тираж 3000 экз. Заказ №0000.

Отпечатано по технологии StP

в ОАО “Печатный двор” им. А. М. Горького
197110, Санкт-Петербург, Чкаловский пр., 15.