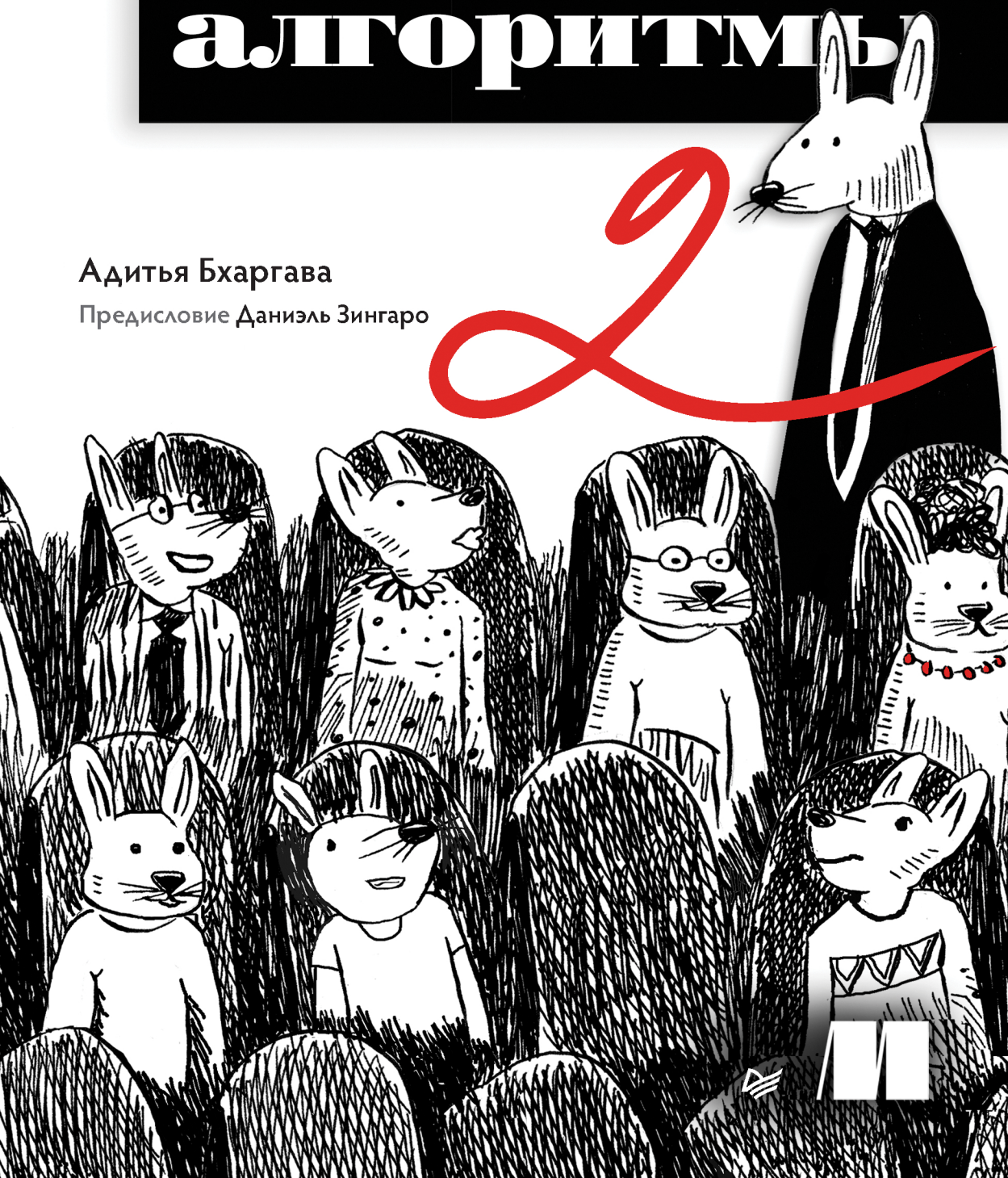


прокаем второе издание

алгоритмы

Адитья Бхаргава

Предисловие Даниэль Зингаро



grokking algorithms

Second Edition

Aditya Y. Bhargava

Foreword by Daniel Zingaro



MANNING
SHELTER ISLAND

прокаем алгоритмы

Второе издание

Адитья Бхаргава

Предисловие Даниэля Зингаро



Санкт-Петербург • Москва • Минск

2024

ББК 32.973-018
УДК 004.421
Б94

Бхаргава Адитья

Б94 Грокаем алгоритмы. 2-е изд. — СПб.: Питер, 2024. — 352 с.: ил. — (Серия «Библиотека программиста»).

ISBN 978-5-4461-4172-2

Алгоритмы — это пошаговые инструкции решения задач, большинство из которых уже были кем-то решены, протестированы и доказали свою эффективность. Второе издание «Грокаем алгоритмы» упрощает изучение, понимание и использование алгоритмов. В этой книге вы найдете простые и внятные объяснения, более 400 забавных иллюстраций и десятки примеров. Ее чтение — лучший способ раскрыть всю мощь алгоритмов и подготовиться к интервью по программированию. Глубоких знаний математики не требуется!

Вы узнаете о главных алгоритмах, позволяющих ускорить работу программ, упростить код и решить распространенные проблемы программирования. Начните с сортировки и поиска, а затем развивайте свои навыки для решения сложных задач, таких как сжатие данных и искусственный интеллект. Научитесь сравнивать эффективность различных алгоритмов.

Во втором издании даны новые, более подробные описания деревьев, NP-полные задачи, а код примеров обновлен на Python 3.

Пора прокатать алгоритмы по-новому!

16+ (В соответствии с Федеральным законом от 29 декабря 2010 г. № 436-ФЗ.)

ББК 32.973-018
УДК 004.421

Права на издание получены по соглашению с Manning Publications. Все права защищены. Никакая часть данной книги не может быть воспроизведена в какой бы то ни было форме без письменного разрешения владельцев авторских прав.

Информация, содержащаяся в данной книге, получена из источников, рассматриваемых издательством как надежные. Тем не менее, имея в виду возможные человеческие или технические ошибки, издательство не может гарантировать абсолютную точность и полноту приводимых сведений и не несет ответственности за возможные ошибки, связанные с использованием книги. В книге возможны упоминания организаций, деятельность которых запрещена на территории Российской Федерации, таких как Meta Platforms Inc., Facebook, Instagram и др. Издательство не несет ответственности за доступность материалов, ссылки на которые вы можете найти в этой книге. На момент подготовки книги к изданию все ссылки на интернет-ресурсы были действующими.

ISBN 978-1633438538 англ.

Authorized translation of the English edition © 2024 Manning Publications.
This translation is published and sold by permission of Manning Publications, the owner of all rights to publish and sell the same.

ISBN 978-5-4461-4172-2

© Перевод на русский язык ООО «Прогресс книга», 2024
© Издание на русском языке, оформление ООО «Прогресс книга», 2024
© Серия «Библиотека программиста», 2024



Краткое содержание

Глава 1. Знакомство с алгоритмами	27
Глава 2. Сортировка выбором	49
Глава 3. Рекурсия	70
Глава 4. Быстрая сортировка	86
Глава 5. Хеш-таблицы	111
Глава 6. Поиск в ширину	139
Глава 7. Деревья	163
Глава 8. Сбалансированные деревья	182
Глава 9. Алгоритм Дейкстры	207
Глава 10. Жадные алгоритмы	238
Глава 11. Динамическое программирование	253
Глава 12. Алгоритм k ближайших соседей	283
Глава 13. Что дальше?	303
Приложение А. Производительность AVL-деревьев	322
Приложение Б. NP-трудные задачи	324
Приложение В. Ответы к упражнениям	336



Оглавление

Отзывы о первом издании	13
Предисловие	15
Введение	17
Благодарности	19
О книге	22
Как работать с этой книгой	23
Для кого предназначена эта книга	23
Структура книги	24
О коде в книге	25
Об авторе	26
О научном редакторе	26
От издательства	26
Глава 1. Знакомство с алгоритмами	27
Введение	27
Что вы узнаете об эффективности алгоритмов	28
Что вы узнаете о решении задач	28

Бинарный поиск	29
Более эффективный поиск	32
<i>Упражнения</i>	36
Время выполнения	37
«О-большое»	38
Время выполнения алгоритмов растет с разной скоростью	38
Наглядное представление «О-большое»	41
«О-большое» определяет время выполнения в худшем случае	43
Типичные примеры «О-большого»	44
<i>Упражнения</i>	45
Задача о коммивояжере	46
Шпаргалка	48
Глава 2. Сортировка выбором	49
Как работает память	50
Массивы и связанные списки	52
Связанные списки	54
Массивы	55
Терминология	56
<i>Упражнение</i>	57
Вставка в середину списка	58
Удаление	60
Какая структура данных используется чаще: массивы или списки?	60
<i>Упражнения</i>	63
Сортировка выбором	64
Пример кода	68
Шпаргалка	69
Глава 3. Рекурсия	70
Рекурсия	71
Базовый случай и рекурсивный случай	74

Стек	76
Стек вызовов	77
Упражнение	79
Стек вызовов с рекурсией	80
Упражнение	84
Шпаргалка	85
Глава 4. Быстрая сортировка	86
«Разделяй и властвуй»	87
Упражнения	96
Быстрая сортировка	96
Снова об «О-большом»	103
Сортировка слиянием и быстрая сортировка	104
Средний и худший случаи	106
Упражнения	109
Шпаргалка	110
Глава 5. Хеш-таблицы	111
Хеш-функции	114
Упражнения	119
Примеры использования	120
Использование хеш-таблиц для поиска	120
Исключение дубликатов	122
Использование хеш-таблицы как кэша	124
Шпаргалка	128
Коллизии	128
Быстродействие	131
Коэффициент заполнения	134
Хорошая хеш-функция	136
Упражнения	137
Шпаргалка	138

Глава 6. Поиск в ширину	139
Знакомство с графами	140
Что такое граф?	142
Поиск в ширину	144
Поиск кратчайшего пути	147
Очереди	148
Упражнения	149
Реализация графа	150
Реализация алгоритма	153
Время выполнения	158
Упражнения	159
Шпаргалка	162
Глава 7. Деревья	163
Ваше первое дерево	164
Каталоги файлов	165
Космическая одиссея: поиск в глубину	168
Правильное определение дерева	172
Бинарные деревья	172
Код Хаффмана	174
Шпаргалка	181
Глава 8. Сбалансированные деревья	182
Балансировка	183
Деревья повышают скорость вставки	184
Короткие деревья работают быстрее	187
AVL-деревья: разновидность сбалансированных деревьев	191
Повороты	192
Как AVL-дерево узнает, что требуется поворот?	195
Косые деревья	201

В-деревья	203
Какие преимущества есть у В-деревьев?	204
Шпаргалка	206
Глава 9. Алгоритм Дейкстры	207
Работа с алгоритмом Дейкстры	208
Терминология	213
История одного обмена	216
Ребра с отрицательным весом	223
Реализация	226
Упражнение	237
Шпаргалка	237
Глава 10. Жадные алгоритмы	238
Задача составления расписания	239
Задача о рюкзаке	241
Упражнения	243
Задача о покрытии множества	243
Приближенные алгоритмы	245
Шпаргалка	252
Глава 11. Динамическое программирование	253
И снова задача о рюкзаке	253
Простое решение	254
Динамическое программирование	255
Задача о рюкзаке: вопросы	264
Что произойдет при добавлении элемента?	264
Упражнение	267
Что произойдет при изменении порядка строк?	267
Можно ли заполнять таблицу по столбцам, а не по строкам?	268
Что произойдет при добавлении меньшего элемента?	268

Можно ли взять часть предмета?	268
Оптимизация туристического маршрута	270
Взаимозависимые элементы	271
Может ли оказаться, что решение требует более двух «подрюкзак»?	272
Возможно ли, что при лучшем решении в рюкзаке остается пустое место?	273
<i>Упражнение</i>	273
Самая длинная общая подстрока	273
Построение таблицы	275
Заполнение таблицы	276
Решение	277
Самая длинная общая подпоследовательность	279
Самая длинная общая подпоследовательность — решение	280
<i>Упражнение</i>	282
Шпаргалка	282

Глава 12. Алгоритм k ближайших соседей 283

Апельсины и грейпфруты	283
Построение рекомендательной системы	286
Извлечение признаков	287
<i>Упражнения</i>	292
Регрессия	293
Выбор признаков	296
<i>Упражнение</i>	297
Знакомство с машинным обучением	297
OCR	297
Построение спам-фильтра	298
Прогнозы на биржевых торгах	299
Тренировка модели МО: общий обзор	300
Шпаргалка	302

Глава 13. Что дальше?	303
Линейная регрессия	303
Инвертированные индексы	305
Преобразование Фурье	306
Параллельные алгоритмы	307
map/reduce	308
Фильтры Блума и HyperLogLog	309
Фильтры Блума	310
HyperLogLog	311
HTTPS и обмен ключами Диффи — Хеллмана	312
Локально-чувствительное хеширование	317
Минимальные кучи и приоритетные очереди	318
Линейное программирование	320
Эпилог	321
Приложение А. Производительность AVL-деревьев	322
Приложение Б. NP-трудные задачи	324
Задачи разрешимости	325
Задачи выполнимости	326
Трудно решить, легко проверить	329
Сведение	331
NP-трудность	332
NP-полнота	333
Шпаргалка	335
Приложение В. Ответы к упражнениям	336

Отзывы о первом издании

Этой книге удалось невозможное: сделать математику простой и интересной!

Сандер Россел (Sander Rossel), COAS Software Systems

А вам бы хотелось изучать алгоритмы так же, как вы читаете свой любимый роман? Если да, то это именно та книга, которая вам нужна!

Санкар Раманатан (Sankar Ramanathan), IBM Analytics

В современном мире нет ни одной сферы жизни, которую бы не оптимизировал какой-нибудь алгоритм. Те, кому требуется понятное введение в тему алгоритмов, должны в первую очередь прочесть эту книгу.

Амит Ламба (Amit Lamba), Tech Overture, LLC

Алгоритмы — это вовсе не скучно! Эта книга была интересной и полезной для меня и моих студентов.

Кристофер Хаунт (Christopher Haupt), Mobirobo, Inc

.....

*Посвящается моим родителям —
Сангити и Йогешу*



Предисловие

В наши дни количество людей, желающих научиться программировать, больше, чем когда-либо. Разумеется, многие профессии напрямую связаны с программированием (например, разработчик или веб-разработчик). Но оно проникает (или начнет проникать) и в другие сферы, традиционно, казалось бы, не связанные с ним. Программирование также помогает людям понять технологичный мир, в котором они живут.

К сожалению, не у всех имеются равные возможности доступа к программированию. Например, в США среди изучающих программирование очень мало женщин и представителей некоторых этнических/расовых групп. Очень важно, чтобы программирование и computer science получили распространение в разнообразных группах. Решение этой проблемы потребует изменений на многих направлениях, включая преодоление предрассудков, подготовку большего количества преподавателей и обеспечение разнообразного опыта обучения. Нужно помочь «включиться» большему количеству людей.

Книга Бхаргавы мне очень понравилась, потому что она помогает по-новому взглянуть на алгоритмы, которые лежат в основе эффективного программирования. Кто-то скажет, что есть только один способ изучать алгоритмы — найти толстый том по математике, прочитать его и (вроде бы) понять все написанное. Но такой подход хорош для тех, кто может учиться подобным образом, у кого есть на это время, и прежде всего — кому это нравится. Этот

метод также основан на том, что мы знаем, для чего хотим изучать алгоритмы. Но, откровенно говоря, это слишком смелое предположение.

На всякий случай оговорюсь, что некоторые из моих любимых книг по computer science именно такие, с математическим уклоном. Мне они подходят. Они подходят и многим профессорам в нашей области. Но, может быть, проблема именно в этом: слишком велик соблазн считать, что другие учатся так же, как мы. Нам нужны обучающие ресурсы по разным областям computer science для разных аудиторий.

Книга Бхаргавы предназначена для тех, кто хочет изучать алгоритмы без лишней математики. Больше всего меня впечатлило даже не то, что Бхаргава включил в книгу, а то, что он решил *не* включать. В такой книге охватить все не получится — она окажется слишком объемной, да это и не нужно.

Благодаря своему опыту преподавания Бхаргаве удалось уместить много полезного в небольшом объеме текста. Например, меня поразило, как в главе «Динамическое программирование» Бхаргава аккуратно дает ответы на многие возможные вопросы читателей, на которые другие книги, посвященные алгоритмам, не отвечают.

Надеюсь, эта книга поможет вам учиться — и тем, кто только пытается освоить алгоритмы, и тем, кому до сих пор не удавалось найти подходящий ресурс. Удачного грокания!

Даниэль Зингаро, Университет Торонто



Введение

Сначала программирование было для меня простым увлечением. Я изучил азы по книге «Visual Basic для чайников», а потом стал читать другие книги, чтобы узнать больше. Но алгоритмы мне никак не давались. Помню, как я смаковал оглавление своей первой книги по алгоритмам и думал: «Наконец-то я все узнаю!» Но материал оказался слишком сложным, и я сдался через несколько недель. Только благодаря хорошему преподавателю теории алгоритмов я понял, насколько простые и элегантные идеи заложены в ее основу.

Я написал свое первое иллюстрированное сообщение в блоге в 2012 году. Сам я визуал, поэтому мне нравится наглядный стиль изложения. С тех пор я создал немало иллюстрированных материалов по функциональному программированию, Git, машинному обучению и параллелизму. Кстати говоря, в начале своей карьеры я писал довольно посредственно. Объяснять научные концепции трудно. Чтобы придумать хорошие примеры, требуется время, чтобы объяснить сложную концепцию — тоже. Проще всего умолчать о сложных моментах. Я думал, что у меня все хорошо получается, пока после одной из моих популярных публикаций ко мне не обратился коллега со словами: «Я прочитал твой материал, но все равно ничего не понял». Мне еще предстояло многое узнать о том, как пишутся научные тексты.

В самом разгаре работы над иллюстрированными публикациями в блоге ко мне обратилось издательство *Manning* с предложением написать иллюстрированную книгу. Оказалось, что редакторы *Manning* хорошо умеют объяснять научные концепции, и они показали мне, как следует учить других. У меня была совершенно определенная цель: мне хотелось создать книгу, которая объясняла бы сложные научные темы и легко читалась.

Первое издание книги вышло в 2016 году. С тех пор книгу прочитало более 100 000 человек. Мне приятно, что очень многим подошел ее наглядный стиль обучения.

Во втором издании моя цель осталась прежней. В этой книге я использую рисунки и легко запоминающиеся примеры, чтобы материал лучше закреплялся в памяти. Книга написана для читателей, которые умеют программировать и хотят узнать больше об алгоритмах без обязательного знания математики.

Второе издание заполняет некоторые пробелы первого. Многие читатели просили объяснить концепцию деревьев. Поэтому в книге появились две главы, посвященные деревьям. Также был расширен раздел, где я рассказываю о NP-полноте. Концепция NP-полноты весьма абстрактна, и мне хотелось привести объяснение, придающее ей больше конкретности. Если вам тоже этого не хватало, надеюсь, раздел о NP-полноте заполнит пробел в ваших знаниях.

С того первого поста в блоге я прошел долгий путь. Надеюсь, что мне удалось сделать книгу простой и содержательной.



Благодарности

Спасибо издательству Manning, которое дало мне возможность написать эту книгу и предоставило большую творческую свободу в ходе работы. Я благодарен издателю Марджану Бейсу (Marjan Bace), Майку Стивенсу (Mike Stephens) за то, что он ввел меня в курс дела, и Иэну Хоу (Ian Hough) — невероятно отзывчивому редактору, всегда готовому прийти на помощь. Спасибо всем участникам производственной группы Manning: Полу Веллсу (Paul Wells), Дебби Хольмгрен (Debbie Holmgren) и всем остальным. Кроме того, я хочу поблагодарить всех, кто читал рукопись и делился своим мнением: Даниэля Зингаро (Daniel Zingaro), Бена Вайнегара (Ben Vinegar), Александра Мэннинга (Alexander Manning) и Мэгги Венгер (Maggie Wenger). Спасибо Дэвиду Айзенштату (David Eisenstat), моему научному редактору, и Тони Холдройду (Tony Holdroyd), техническому корректору Manning, за исправление множества моих ошибок.

Спасибо всем, кто помог мне в достижении цели: Берту Бейтсу (Bert Bates), который научил меня писать на научные темы, сотрудникам *Flashkit*, научившим меня программировать; многочисленным друзьям, которые помогали мне в работе — рецензировали главы, делились советами и предлагали разные варианты объяснений. Это были Бен Вайнгер (Ben Vinegar), Карл Пьюзон (Karl Puzon), Алекс Мэннинг (Alex Manning), Эстер Чан (Esther Chan), Аниш Бхатт (Anish Bhatt), Майкл Гласс (Michael Glass), Никрад

Махди (Nikrad Mahdi), Чарльз Ли (Charles Lee), Джаред Фридман (Jared Friedman), Хема Маникавасагам (Hema Manickavasagam), Хари Раджа (Hari Raja), Мурали Гудипати (Murali Gudipati), Шриниваса Варадан (Srinivas Varadan) и другие; также спасибо Джерри Брэди (Gerry Brady), моему учителю по теории алгоритмов. Отдельное большое спасибо таким классикам алгоритмов, как CLRS¹, Кнут и Стрэнг; безусловно, я стою на плечах гигантов.

Папа, мама, Приянка и все родные: спасибо за вашу неустанную поддержку. Огромное спасибо моей жене Мэгги и моему сыну Йоги. Впереди у нас много прекрасных моментов, и мне уже не придется проводить вечер пятницы за переписыванием книги.

Спасибо всем рецензентам:

Абишеку Косервалю (Abhishek Koserwal), Алексу Лукасу (Alex Lucas), Андресу Сакко (Andres Sacco), Аруну Сахе (Arun Saha), Бекки Хьют (Becky Huett), Сезару Аугусто Ороско Манотасу (Cesar Augusto Orozco Manotas), Кристиану Саттону (Christian Sutton), Диогинешу Гольдони (Diogines Goldoni), Дирку Гомесу (Dirk Gomez), Эду Бахеру (Ed Bacher), Эдер Андрес Авила Ниньо (Eder Andres Avila Nino), Франсу Оилинки (Frans Oilinki), Ганешу Свамнатану (Ganesh Swaminathan), Джампиеро Гранателле (Giampiero Granatella), Глену Ю (Glen Yu), Грегу Крейтеру (Greg Kreiter), Явиду Асгарову (Javid Asgarov), Жоао Ферейре (Joao Ferreira), Жобинешу Пурушотаману (Jobinesh Purushothaman), Джо Куэвасу (Joe Cuevas), Джошу Макадамсу (Josh McAdams), Кришне Анипинди (Krishna Anipindi), Кшиштофу Камычеку (Krzysztof Kamyczek), Кирилу Калиниченко (Kyrylo Kalinichenko), Лакшминараяану АС (Lakshminarayanan AS), Лоду Бентилу (Laud Bentil), Маттео Батисте (Matteo Battista), Микаэлю Бистрому (Mikael Bystrom), Нику Ракоши (Nick Rakochy), Нинославу Черкесу (Ninoslav Cerkez), Оливеру Кортену (Oliver Korten), Оои Кван Сану (Ooi Kuan San), Пабло Вареле (Pablo Varela), Патрику Регану (Patrick Regan), Патрику Ванъяу (Patrick Wanjau), Филиппу Конраду (Philipp Konrad), Петру Пинделу (Piotr Pindel), Раджешу Монахану (Rajesh Mohanan), Ранжиту Сахаи (Ranjit Sahai), Рохини Уппулари

¹ Авторы классической книги по алгоритмам: Кормен, Лейзерсон, Ривест, Штайн. — *Примеч. пер.*

(Rohini Uppuluri), Роману Левченко, Самбарану Хазре (Sambaran Hazra), Сету Макферсону (Seth MacPherson), Шанкару Свами (Shankar Swamy), Шрихари Шридхарану (Srihari Sridharan), Тобиасу Копфу (Tobias Kopf), Вивеку Веераппану (Vivek Veerappan), Вильяму Джамиру Сильве (William Jamir Silva) и Сян Бо Мао (Xiangbo Mao) — ваши предложения помогли сделать эту книгу лучше.

Наконец, я хочу поблагодарить всех читателей, которые заинтересовались книгой, и тех, кто поделился своим мнением на форуме книги. Благодаря вам она действительно стала лучше.



О книге

Я прежде всего стремился к тому, чтобы книга легко читалась. Я избегаю неожиданных поворотов; каждый раз, когда в книге упоминается новая концепция, я либо объясняю ее сразу, либо говорю, где буду объяснять. Основные концепции подкрепляются упражнениями и повторными объяснениями, чтобы вы могли проверить свои предположения и убедиться в том, что не потеряли нить изложения.

В книге приводится множество примеров. Моя цель — не вывалить на читателя кучу невразумительных формул, а упростить наглядное представление этих концепций. Я также считаю, что мы лучше всего учимся тогда, когда можем вспомнить что-то уже известное, а примеры помогают освежить память. Так, когда вы вспоминаете, чем массивы отличаются от связанных списков (глава 2), просто вспомните, как ищите места для компании в кинотеатре. Наверное, вы уже поняли, что я сторонник визуального стиля обучения — в книге полно рисунков.

Содержимое книги было тщательно продумано. Нет смысла писать книгу, в которую будут включены все алгоритмы сортировки — для этого есть такие источники, как Википедия и *Khan Academy*. Все алгоритмы, описанные в книге, имеют практическую ценность. Я применял их в своей работе программиста, и они закладывают хорошую основу для изучения более сложных тем.

Приятного чтения!

Как работать с этой книгой

Порядок изложения и содержимое книги были тщательно продуманы. Если вас очень сильно интересует какая-то тема — переходите прямо к ней. В противном случае читайте главы по порядку, они логически переходят одна в другую.

Я настоятельно рекомендую самостоятельно выполнять код всех примеров. Вы не поверите, насколько это важно. Просто введите мои примеры кода «с листа» (или загрузите их по адресу <https://www.manning.com/books/grokking-algorithms-second-edition> или https://github.com/egonschiele/grokking_algorithms) и выполните. Так у вас в памяти останется гораздо больше, чем просто при чтении.

Также я рекомендую выполнить упражнения, приведенные в книге. Упражнения не займут много времени — обычно задачи решаются за минуту или две, иногда за 5–10 минут. Упражнения помогут проверить правильность понимания материала. Если вы где-то сбились с пути, то узнаете об этом, не заходя слишком далеко.

Для кого предназначена эта книга

Эта книга предназначена для читателей, которые владеют азами программирования и хотят разобраться в алгоритмах. Может быть, вы уже столкнулись с задачей программирования и пытаетесь найти алгоритмическое решение. А может, вы хотите понять, где вам могут пригодиться алгоритмы. Ниже приведен короткий и неполный список людей, которым может пригодиться книга:

- программисты-самоучки;
- студенты, начавшие изучать программирование;
- выпускники, желающие освежить память;
- специалисты по физике/математике/другим дисциплинам, интересующиеся программированием.

Структура книги

В первых трех главах закладываются основы.

- **Глава 1:** вы изучите свой первый нетривиальный алгоритм: бинарный поиск. Также здесь рассматриваются основы анализа скорости алгоритмов с применением «О-большое». Эта запись часто используется в книге для описания относительной быстроты выполнения алгоритмов.
- **Глава 2:** вы познакомитесь с двумя основополагающими структурами данных: массивами и связанными списками. Эти структуры данных часто встречаются в книге и используются для создания более сложных структур данных, например хеш-таблиц (глава 5).
- **Глава 3:** вы узнаете о рекурсии — удобном приеме, используемом многими алгоритмами (например, алгоритмом быстрой сортировки, о котором рассказано в главе 4).

По моему опыту, темы «О-большое» и рекурсии сложны для новичков, поэтому в соответствующих разделах я снижаю темп изложения и привожу более подробные объяснения.

В оставшейся части книги представлены алгоритмы, часто применяемые в разных областях.

- **Методы решения задач** рассматриваются в главах 4, 10 и 11. Если вы столкнулись со сложной задачей и не знаете, как эффективно ее решить, воспользуйтесь стратегией «разделяй и властвуй» (глава 4) или методом динамического программирования (глава 11). А если вы поняли, что эффективного решения не существует, попробуйте получить приближенный ответ с использованием жадного алгоритма (глава 10).
- **Хеш-таблицы** рассматриваются в главе 5. Хеш-таблицы — исключительно полезная структура данных, предназначенная для хранения пар ключей и значений (например имени человека и адреса электронной почты или имени пользователя и пароля). Трудно переоценить практическую полезность хеш-таблиц. Приступая к решению задачи, я обычно прежде всего задаю себе два вопроса: можно ли здесь воспользоваться хеш-таблицей и можно ли смоделировать задачу в виде графа.
- **Графы и деревья** рассматриваются в главах 6, 7, 8 и 9. Графы используются для моделирования сетей: социальных, дорожных, нейронных

или любых других совокупностей связей. Поиск в ширину (глава 6) и алгоритм Дейкстры (глава 9) предназначены для поиска кратчайшего расстояния между двумя точками сети: с их помощью можно вычислить кратчайший маршрут к точке назначения или количество промежуточных знакомых у двух людей в социальной сети. Деревья — разновидность графа. Они используются в базах данных (обычно B-деревья), браузерах (DOM-деревья) или файловой системе.

- **Алгоритм k ближайших соседей** рассматривается в главе 12. Это простой алгоритм машинного обучения; с его помощью можно построить рекомендательную систему, механизм оптического распознавания текста, систему прогнозирования курсов акций — словом, всего, что требует прогнозирования значений («Мы думаем, что Адит поставит этому фильму четыре звезды») или классификации объектов («Это буква Q»).
- **Следующий шаг:** в главе 13 представлен еще ряд алгоритмов, которые хорошо подойдут для дальнейшего изучения темы.

О коде в книге

Во всех примерах в книге используется Python 3. Весь программный код оформлен моноширинным шрифтом, чтобы его можно было отличить от обычного текста. Некоторые листинги сопровождаются аннотациями, подчеркивающими важные концепции.

Исполняемые фрагменты кода можно загрузить из версии liveBook (электронной) по адресу <https://livebook.manning.com/book/grokkingalgorithms-second-edition>. Полный код примеров книги доступен для загрузки на сайте Manning <https://www.manning.com> и по адресу https://github.com/egonschiele/grokking_algorithms.

Я считаю, что мы лучше всего учимся тогда, когда нам это нравится, — так что получайте удовольствие от процесса... и запускайте примеры кода!

ФОРУМ LIVEBOOK

Приобретая книгу, вы получаете бесплатный доступ к закрытому веб-форуму издательства Manning (на английском языке), на котором можно оставлять комментарии о книге, задавать технические вопросы и получать помощь от автора и других пользователей. Чтобы получить доступ к фору-

му, откройте страницу <https://livebook.manning.com/book/grokking-algorithms-second-edition>. За информацией о форумах Manning и правилах поведения на них обращайтесь по адресу <https://livebook.manning.com/discussion>.

В рамках своих обязательств перед читателями издательство Manning предоставляет ресурс для содержательного общения читателей и авторов. Эти обязательства не подразумевают конкретную степень участия автора, которое остается добровольным (и неоплачиваемым). Задавайте автору хорошие вопросы, чтобы он не терял интереса к происходящему! Форум и архивы обсуждений доступны на веб-сайте издательства, пока книга продолжает издаваться.

Об авторе



Адितья Бхаргава работает программистом. Он получил степень магистра по computer science в Чикагском университете и ведет популярный иллюстрированный технический блог *adit.io*.

О научном редакторе

Дэвид Айзенштат (David Eisenstat) — инженер-исследователь программного обеспечения. Он получил степень PhD в области computer science в Университете Брауна.

От издательства

Ваши замечания, предложения, вопросы отправляйте по адресу

`comp@piter.com`

(издательство «Питер», компьютерная редакция).

Мы будем рады узнать ваше мнение!

На веб-сайте издательства www.piter.com вы найдете подробную информацию о наших книгах.

1

Знакомство с алгоритмами



В этой главе

- ✓ Закладываются основы для остальных глав книги.
- ✓ Вы напишете свой первый алгоритм поиска (бинарный поиск).
- ✓ Вы узнаете, как описывается время выполнения алгоритма («О-большое»).

Введение

Алгоритмом называется набор инструкций для выполнения некоторой задачи. В принципе любой фрагмент программного кода можно назвать алгоритмом, но в этой книге рассматриваются более интересные темы. Когда я отбирал алгоритмы для этой книги, то следил за тем, чтобы они были быстрыми или решали интересные задачи... или и то и другое сразу. Вот лишь несколько примеров.

- В главе 1 речь пойдет о бинарном поиске и о том, как алгоритмы могут ускорить работу кода. В одном примере алгоритм сокращает количество необходимых действий с 4 миллиардов до 32!

- Устройство GPS использует алгоритмы из теории графов (об этом в главах 6 и 9) для вычисления кратчайшего пути к точке назначения.
- При помощи методов динамического программирования (см. главу 11) можно создать алгоритм для игры в шашки.

В каждом случае я опишу алгоритм и приведу пример. Затем мы обсудим время выполнения алгоритма в понятиях «О-большое». В завершение будут рассмотрены типы задач, которые могут решаться с применением того же алгоритма.

Что вы узнаете об эффективности алгоритмов

А теперь хорошая новость: скорее всего, реализация каждого алгоритма в этой книге уже доступна на вашем любимом языке программирования и вам не придется писать каждый алгоритм самостоятельно! Но любая реализация будет бесполезной, если вы не понимаете ее плюсов и минусов. В этой книге вы научитесь сравнивать сильные и слабые стороны разных алгоритмов: из каких соображений выбирать между сортировкой слиянием и быстрой сортировкой? Что использовать — массив или список? Даже выбор другой структуры данных может оказать сильное влияние на результат.

Что вы узнаете о решении задач

Вы освоите методы решения задач, которые вам сейчас, возможно, неизвестны. Примеры:

- Если вы любите создавать видеоигры, вы можете написать систему на базе искусственного интеллекта, моделирующую действия пользователя с применением алгоритмов из теории графов.
- Вы узнаете, как построить рекомендательную систему на базе k ближайших соседей.
- Некоторые проблемы не решаются за разумное время! В части книги, посвященной NP-полноте задач, рассказано о том, как идентифициро-

вать такие задачи и построить алгоритм для получения приближенного ответа.

А если брать шире, к концу этой книги вы освоите некоторые широко применяемые алгоритмы. После этого вы сможете воспользоваться новыми знаниями для изучения более специализированных алгоритмов из области искусственного интеллекта, баз данных и т. д. или взяться за решение более сложных задач в практической работе.

ЧТО НЕОБХОДИМО ЗНАТЬ

Чтобы читать эту книгу, необходимо знать базовую алгебру. Например, возьмем следующую функцию: $f(x) = x \times 2$. Чему равен результат $f(5)$? Если вы ответили «10» — читайте спокойно.

Кроме того, вам будет проще понять эту главу (и всю книгу), если вы владеете хотя бы одним языком программирования. Все приведенные примеры написаны на Python. Если вы не знаете ни одного языка программирования, но хотите изучить — выбирайте Python: это отличный язык для начинающих. Если вы уже знаете другой язык (скажем, JavaScript) — все в порядке.

Бинарный поиск

Предположим, вы ищете фамилию человека в телефонной книге (какая древняя технология!). Она начинается с буквы «К». Конечно, можно начать с самого начала и перелистывать страницы, пока вы не доберетесь до буквы «К». Но для ускорения поиска лучше раскрыть книгу на середине: ведь буква «К» должна находиться где-то ближе к середине телефонной книги.

Или предположим, что вы ищете слово в словаре и оно начинается с буквы «О». И снова лучше начать с середины.



Теперь допустим, что вы вводите свои данные при входе на Facebook. При этом Facebook необходимо проверить, есть ли у вас учетная запись на сайте. Для этого ваше имя пользователя нужно найти в базе данных. Допустим, вы выбрали себе имя пользователя «karlmageddon». Facebook может начать с буквы «А» и проверять все подряд, но разумнее будет начать с середины.



Перед нами типичная задача поиска. И во всех этих случаях для решения задачи можно применить один алгоритм: *бинарный поиск*.

Бинарный поиск — это алгоритм; на входе он получает отсортированный список элементов (позднее я объясню, почему он должен быть отсортирован). Если элемент, который вы ищете, присутствует в списке, то бинарный поиск возвращает ту позицию, в которой он был найден. В противном случае бинарный поиск возвращает `null`.

Например:



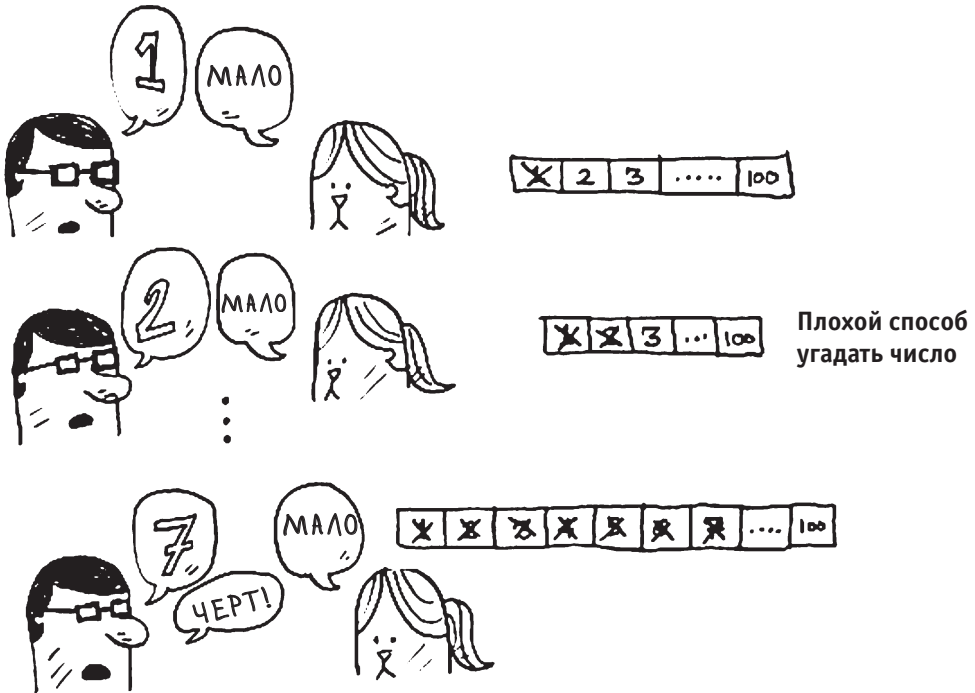
Ищем компанию
в телефонной книге
с применением
бинарного поиска

Рассмотрим пример того, как работает бинарный поиск. Сыграем в простую игру: я загадал число от 1 до 100.



Вы должны отгадать мое число, используя как можно меньше попыток. При каждой попытке я буду давать один из трех ответов: «мало», «много» или «угадал».

Предположим, вы начинаете перебирать все варианты подряд: 1, 2, 3, 4... Вот как это будет выглядеть.



Это пример *простого поиска* (возможно, термин «*тупой поиск*» был бы уместнее). При каждой догадке исключается только одно число. Если я загадал число 99, то, чтобы добраться до него, потребуется 99 попыток!

Более эффективный поиск

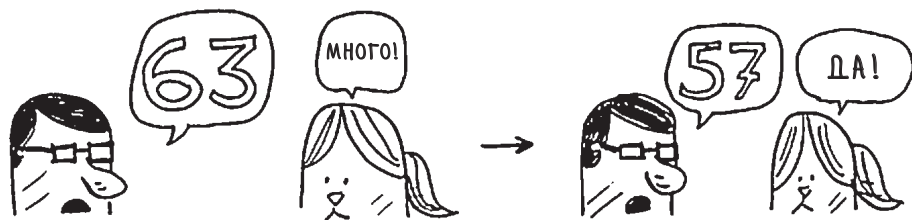
Существует другой, более эффективный способ. Начнем с 50.



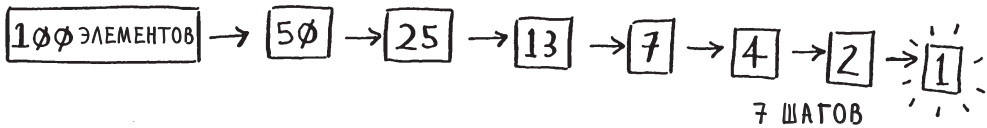
Слишком мало... но вы только что исключили *половину* чисел! Теперь вы знаете, что все числа 1–50 меньше загаданного. Следующая попытка: 75.



На этот раз перелет... Но вы снова исключили половину оставшихся чисел! *С бинарным поиском вы каждый раз загадываете число в середине диапазона и исключаете половину оставшихся чисел.* Следующим будет число 63 (по середине между 50 и 75).



Так работает бинарный поиск. А вы только что узнали свой первый алгоритм! Попробуем поточнее определить, сколько чисел будет исключаться каждый раз.



При бинарном поиске каждый раз исключается половина чисел

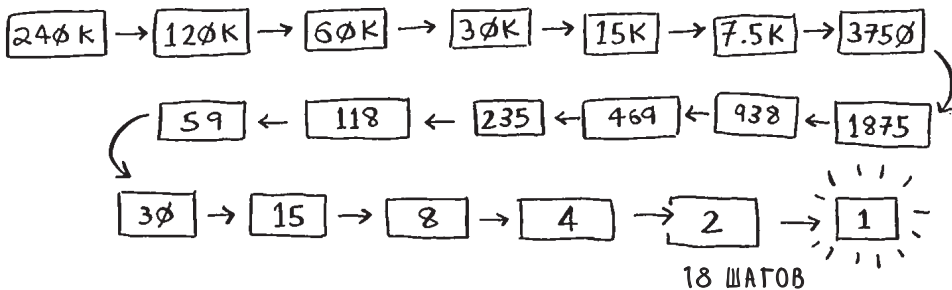
Какое бы число я ни задумал, вы гарантированно сможете угадать его не более чем за 7 попыток, потому что с каждой попыткой исключается половина оставшихся чисел!

Предположим, вы ищете слово в словаре с 240 000 словами. Как вы думаете, сколько попыток вам понадобится в худшем случае?

ПРОСТОЙ ПОИСК: ____ шагов

БИНАРНЫЙ ПОИСК: ____ шагов

При простом поиске может потребоваться 240 000 попыток, если искомое слово находится на самой последней позиции в книге. С каждым шагом бинарного поиска количество слов сокращается вдвое, пока не останется только одно слово.



Итак, бинарный поиск потребует 18 шагов — заметная разница! В общем случае для списка из n элементов бинарный поиск выполняется за $\log_2 n$ шагов, тогда как простой поиск будет выполнен за n шагов.

ЛОГАРИФМЫ

Возможно, вы уже забыли, что такое логарифм, но наверняка помните, что такое возведение в степень. $\log_{10} 100$, по сути, означает, сколько раз нужно перемножить 10, чтобы получить 100. Правильный ответ — 2: 10×10 . Итак, $\log_{10} 100 = 2$. Логарифм по смыслу противоположен возведению в степень.

$$\begin{array}{l} 10^2 = 100 \quad \leftrightarrow \quad \log_{10} 100 = 2 \\ \hline 10^3 = 1000 \quad \leftrightarrow \quad \log_{10} 1000 = 3 \\ \hline 2^3 = 8 \quad \leftrightarrow \quad \log_2 8 = 3 \\ \hline 2^4 = 16 \quad \leftrightarrow \quad \log_2 16 = 4 \\ \hline 2^5 = 32 \quad \leftrightarrow \quad \log_2 32 = 5 \end{array}$$

Логарифм — операция, обратная возведению в степень

Когда я в этой книге упоминаю «О-большое» (об этом чуть позднее), \log всегда означает \log_2 . Когда вы ищете элемент с применением простого поиска, в худшем случае вам придется проверить каждый элемент. Итак, для списка из 8 чисел понадобится не больше 8 проверок. Для бинарного поиска в худшем случае потребуется не более \log_n проверок. Для списка из 8 элементов $\log 8 = 3$, потому что $2^3 = 8$. Итак, для списка из 8 чисел вам придется проверить не более 3 чисел. Для списка из 1024 элементов $\log 1024 = 10$, потому что $2^{10} = 1024$. Следовательно, для списка из 1024 чисел придется проверить не более 10 чисел.

ПРИМЕЧАНИЕ

В этой книге я буду много говорить о логарифмах, поэтому нужно четко представлять, что это такое. Если вы пока не очень хорошо разбираетесь в логарифмах, посмотрите отличное обучающее видео в Академии Хана (<https://khanacademy.org>).

ПРИМЕЧАНИЕ

Бинарный поиск работает только в том случае, если список отсортирован. Например, имена в телефонной книге хранятся в алфавитном порядке, и вы можете воспользоваться бинарным поиском. А что произойдет, если имена не будут отсортированы?

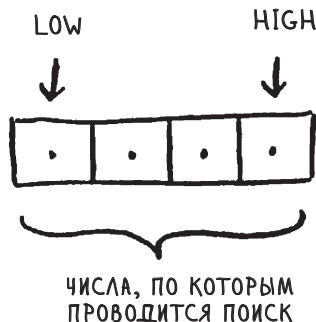
Посмотрим, как написать реализацию бинарного поиска на Python. В следующем примере кода используется массив. Если вы не знаете, как работают массивы, не беспокойтесь: эта тема рассматривается в следующей главе. Пока достаточно знать, что серию элементов можно сохранить в непрерывной последовательности ячеек, которая называется массивом. Нумерация ячеек начинается с 0: первая ячейка находится в позиции с номером 0, вторая — в позиции с номером 1 и т. д.

ПРИМЕЧАНИЕ

Поскольку в Python массивы называются списками, я буду использовать термины "список" и "массив" как взаимозаменяемые.

Функция `binary_search` получает отсортированный массив и значение. Если значение присутствует в массиве, то функция возвращает его позицию. При этом мы должны следить за тем, в какой части массива проводится поиск. Вначале это весь массив:

```
low = 0  
high = len(list) - 1
```



Каждый раз алгоритм проверяет средний элемент:

```
mid = (low + high) / 2
guess = list[mid]
```

Если значение (low+high) нечетно, то Python автоматически округляет значение mid в меньшую сторону

Если названное число было слишком мало, то переменная low обновляется соответственно:

```
if guess < item:
    low = mid + 1
```

А если догадка была слишком велика, то обновляется переменная high. Полный код выглядит так:

```
def binary_search(list, item):
    low = 0
    high = len(list) - 1

    while low <= high:
        mid = (low + high) // 2
        guess = list[mid]
        if guess == item:
            return mid
        elif guess > item:
            high = mid - 1
        else:
            low = mid + 1
    return None
```

В переменных low и high хранятся границы той части списка, в которой выполняется поиск

Пока эта часть не сократится до одного элемента...
...проверяем средний элемент

Значение найдено

Много

Мало

Значения не существует

my_list = [1, 3, 5, 7, 9] А теперь протестируем функцию!

```
print binary_search(my_list, 3) # => 1
print binary_search(my_list, -1) # => None
```

Вспомните: нумерация элементов начинается с 0. Второй ячейке соответствует индекс 1

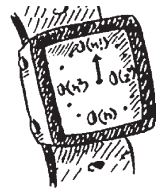
"None" в Python означает "ничто". Это признак того, что элемент не найден

УПРАЖНЕНИЯ

- 1.1 Имеется отсортированный список из 128 имен, и вы ищете в нем значение методом бинарного поиска. Какое максимальное количество проверок для этого может потребоваться?
- 1.2 Предположим, размер списка увеличился вдвое. Как изменится максимальное количество проверок?

Время выполнения

Каждый раз, когда мы будем рассматривать очередной алгоритм, я буду обсуждать время его выполнения. Обычно следует выбирать самый эффективный алгоритм, будь то оптимизация по времени или памяти.



Вернемся к бинарному поиску. Сколько времени сэкономит его применение? В первом варианте мы последовательно проверяли каждое число, одно за другим. Если список состоит из 100 чисел, может потребоваться до 100 попыток. Для списка из 4 миллиардов чисел потребуется до 4 миллиардов попыток. Таким образом, максимальное количество попыток совпадает с размером списка. Такое время выполнения называется *линейным*.

С бинарным поиском дело обстоит иначе. Если список состоит из 100 элементов, потребуется не более 7 попыток. Для списка из 4 миллиардов элементов потребуется не более 32 попыток. Впечатляет, верно? Бинарный поиск выполняется за *логарифмическое время*. В следующей таблице приводится краткая сводка результатов.

ПРОСТОЙ ПОИСК	БИНАРНЫЙ ПОИСК	
100 ЭЛЕМЕНТОВ	100 ЭЛЕМЕНТОВ	
↓	↓	
100 ПОПЫТОК	7 ПОПЫТОК	БОЛЬШАЯ ЭКОНОМИЯ!
4 000 000 000 ЭЛЕМЕНТОВ	4 000 000 000 ЭЛЕМЕНТОВ	
↓	↓	
4 000 000 000 ПОПЫТОК	32 ПОПЫТКИ	БОЛЬШАЯ ЭКОНОМИЯ!
$O(n)$	$O(\log n)$	
↑	↑	
ЛИНЕЙНОЕ ВРЕМЯ	ЛОГАРИФМИЧЕСКОЕ ВРЕМЯ	

Время выполнения алгоритмов поиска



«О-большое»

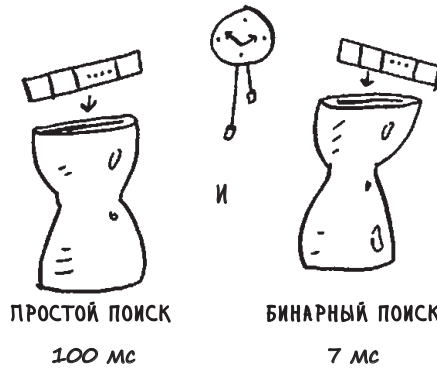
Специальная нотация «О-большое» описывает скорость работы алгоритма. Зачем вам это? Время от времени вам придется использовать чужие алгоритмы, а потому неплохо было бы понимать, насколько быстро или медленно они работают. В этом разделе я объясню, что представляет собой «О-большое», и приведу список самых распространенных вариантов времени выполнения для некоторых алгоритмов.

Время выполнения алгоритмов растет с разной скоростью

Боб пишет алгоритм поиска для NASA. Его алгоритм заработает, когда ракета будет подлетать к Луне, и поможет вычислить точку посадки.

Это один из примеров того, как время выполнения двух алгоритмов растет с разной скоростью. Боб пытается выбрать между простым и бинарным поиском. Его алгоритм должен работать быстро и правильно. С одной стороны, бинарный поиск работает быстрее. У Боба есть всего 10 секунд, чтобы выбрать место посадки; если он не уложится в это время, то момент для посадки будет упущен. С другой стороны, простой поиск пишется проще и вероятность ошибок в нем ниже... Конечно, Боб *совершенно* не хочет допустить ошибку в коде посадки ракеты. И тогда для пущей уверенности Боб решает измерить время выполнения обоих алгоритмов для списка из 100 элементов.

Допустим, проверка одного элемента занимает 1 миллисекунду (мс). При простом поиске Бобу придется проверить 100 элементов, поэтому поиск займет 100 мс. С другой стороны, при бинарном поиске достаточно проверить всего 7 элементов ($\log_2 100$ равен приблизительно 7), а поиск займет 7 мс. Но реальный список может содержать более миллиарда элементов. Сколько времени в таком случае потребуется для выполнения простого поиска? А при бинарном поиске? Обязательно ответьте на оба вопроса, прежде чем продолжить чтение.



Время выполнения простого и бинарного поиска для списка из 100 элементов

Боб проводит бинарный поиск с 1 миллиардом элементов, и на это уходит 30 мс ($\log_2 1\,000\,000\,000$ равен приблизительно 30). «Тридцать миллисекунд! — думает Боб. — Бинарный поиск в 15 раз быстрее простого, потому что простой поиск для 100 элементов занял 100 мс, а бинарный поиск занял 7 мс. Значит, простой поиск займет $30 \times 15 = 450$ мс, верно? Гораздо меньше отведенных 10 секунд». И Боб выбирает простой поиск. Верен ли его выбор?

Нет, Боб ошибается. Глубоко ошибается. Время выполнения для простого поиска с 1 миллиардом элементов составит 1 миллиард миллисекунд, а это 11 дней! Проблема в том, что время выполнения для бинарного и простого поиска *растет с разной скоростью*.

	ПРОСТОЙ ПОИСК	БИНАРНЫЙ ПОИСК
100 ЭЛЕМЕНТОВ	100 мс	7 мс
10 000 ЭЛЕМЕНТОВ	10 секунд	14 мс
1 000 000 000 ЭЛЕМЕНТОВ	11 дней	30 мс

Время выполнения растет с совершенно разной скоростью!

Другими словами, с увеличением количества элементов бинарный поиск занимает чуть больше времени. А простой поиск займет *гораздо* больше времени. Таким образом, с увеличением списка бинарный список внезап-

но начинает работать *гораздо* быстрее простого. Боб думал, что бинарный поиск работает в 15 раз быстрее простого, но это не так. Если список состоит из 1 миллиарда элементов, бинарный поиск работает приблизительно в 33 миллиона раз быстрее. Вот почему недостаточно знать, сколько времени должен работать алгоритм, — необходимо знать, как возрастает время выполнения с ростом размера списка. Здесь-то вам и пригодится «О-большое».



«О-большое» описывает, насколько быстро работает алгоритм. Предположим, имеется список размера n . Простой поиск должен проверить каждый элемент, поэтому ему придется выполнить n операций. Время выполнения «О-большое» имеет вид $O(n)$. Постойте, но где же секунды? А их здесь нет — «О-большое» не сообщает скорость в секундах, а *позволяет сравнить количество операций*. Оно указывает, насколько быстро возрастает время выполнения алгоритма.

А теперь другой пример. Для проверки списка размером n бинарному поиску потребуется $\log n$ операций. Как будет выглядеть «О-большое»? $O(\log n)$. В общем случае «О-большое» выглядит так:

$\text{«О-БОЛЬШОЕ»} \rightarrow O(n) \leftarrow \text{КОЛИЧЕСТВО ОПЕРАЦИЙ}$

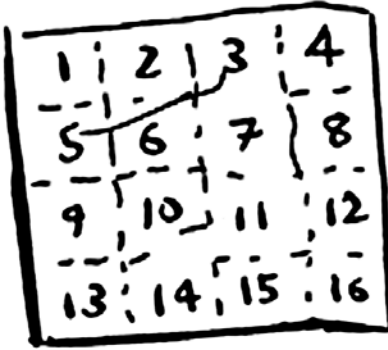
Как записывается «О-большое»

Такая запись сообщает количество операций, которые придется выполнить алгоритму. Она называется «О-большое», потому что перед количеством операций ставится символ «О» (а большое — потому что в верхнем регистре).

Теперь рассмотрим несколько примеров. Попробуйте самостоятельно оценить время выполнения этих алгоритмов.

Наглядное представление «О-большое»

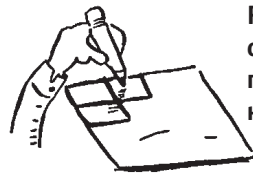
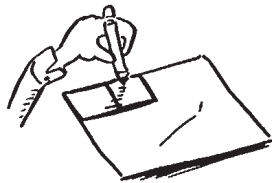
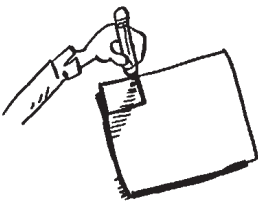
Чтобы повторить следующий практический пример, достаточно иметь несколько листов бумаги и карандаш. Допустим, вы должны построить сетку из 16 квадратов.



Как должен выглядеть хороший алгоритм для построения этой сетки?

Алгоритм 1

Как вариант можно нарисовать 16 квадратов, по одному за раз. Напоминаю: «О-большое» подсчитывает количество операций. В данном примере рисование квадрата считается одной операцией. Нужно нарисовать 16 квадратов. Сколько операций по рисованию одного квадрата придется выполнить?



Рисуем сетку по одному квадрату

Чтобы нарисовать 16 квадратов, потребуется 16 шагов. Как выглядит время выполнения этого алгоритма?

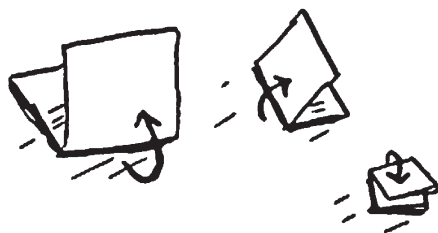
Алгоритм 2

А теперь попробуем иначе. Сложите лист пополам.



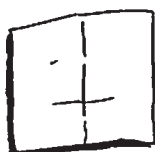
На этот раз операцией считается сложение листка. Получается, что одна операция создает сразу два прямоугольника!

Сложите бумагу еще раз, а потом еще и еще.

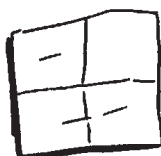


Разверните листок после четырех сложений — получилась замечательная сетка! Каждое сложение удваивает количество прямоугольников. За 4 операции вы создали 16 прямоугольников!

1 СЛОЖЕНИЕ



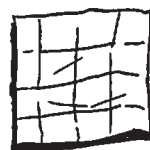
2 СЛОЖЕНИЯ



3 СЛОЖЕНИЯ



4 СЛОЖЕНИЯ



Построение сетки за 4 сложения

При каждом складывании количество прямоугольников увеличивается вдвое, так что 16 прямоугольников строятся за 4 шага. Как записать время выполнения этого алгоритма? Напишите время выполнения обоих алгоритмов, прежде чем двигаться дальше.

Ответы: алгоритм 1 выполняется за время $O(n)$, а алгоритм 2 — за время $O(\log n)$.

«О-большое» определяет время выполнения в худшем случае

Предположим, вы используете простой поиск для поиска фамилии в телефонной книге. Вы знаете, что простой поиск выполняется за время $O(n)$, то есть в худшем случае вам придется просмотреть каждую без исключения запись в телефонной книге. Но представьте, что искомая фамилия начинается на букву «А» и этот человек стоит на самом первом месте в вашей телефонной книге. В общем, вам не пришлось просматривать все записи — вы нашли нужную фамилию с первой попытки. Отработал ли алгоритм за время $O(n)$? А может, он занял время $O(1)$, потому что результат был получен с первой попытки?

Простой поиск все равно выполняется за время $O(n)$. Просто в данном случае вы нашли нужное значение моментально; это лучший возможный случай. Однако «О-большое» описывает *худший* возможный случай. Фактически вы утверждаете, что *в худшем случае* придется просмотреть каждую запись в телефонной книге по одному разу. Это и есть время $O(n)$. И это дает определенные гарантии — вы знаете, что простой поиск никогда не будет работать медленнее $O(n)$.

ПРИМЕЧАНИЕ

Наряду с временем худшего случая также полезно учитывать среднее время выполнения. Тема худшего и среднего времени выполнения обсуждается в главе 4.

Типичные примеры «О-большого»

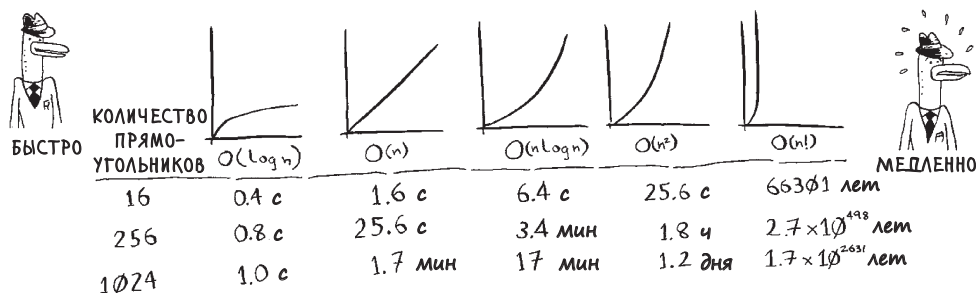
Ниже перечислены пять разновидностей «О-большого», которые будут встречаться вам особенно часто, в порядке убывания скорости выполнения.

- $O(\log n)$, или *логарифмическое время*. Пример: бинарный поиск.
- $O(n)$, или *линейное время*. Пример: простой поиск.
- $O(n * \log n)$. Пример: эффективные алгоритмы сортировки (быстрая сортировка — но об этом в главе 4).
- $O(n^2)$. Пример: медленные алгоритмы сортировки (сортировка выбором — см. главу 2).
- $O(n!)$. Пример: очень медленные алгоритмы (задача о коммивояжере — о ней будет рассказано в следующем разделе).

Предположим, вы снова строите сетку из 16 квадратов и можете выбрать для решения этой задачи один из пяти алгоритмов. При использовании первого алгоритма сетка будет построена за время $O(\log n)$. В секунду выполняются до 10 операций. С временем $O(\log n)$ для построения сетки из 16 квадратов потребуются 4 операции ($\log 16$ равен 4). Итак, сетка будет построена за 0,4 секунды. А если бы было нужно построить 1024 квадрата? На это бы потребовалось $\log 1024 = 10$ операций, или 1 секунда. Напомню, что эти числа получены при использовании первого алгоритма.

Второй алгоритм работает медленнее: за время $O(n)$. Для построения 16 прямоугольников потребуется 16 операций, а для построения 1024 прямоугольников — 1024 операции. Сколько это составит в секундах?

Ниже показано, сколько времени потребуется для построения сетки с остальными алгоритмами, от самого быстрого до самого медленного.



Существуют и другие варианты времени выполнения, но эти пять встречаются чаще всего.

Помните, что эта запись является упрощением. На практике «О-большое» не удастся легко преобразовать в количество операций с такой точностью, но пока нам хватит и этого. Мы вернемся к «О-большому» в главе 4 после рассмотрения еще нескольких алгоритмов. А пока перечислим основные результаты.

- Скорость алгоритмов измеряется не в секундах, а в темпе роста количества операций.
- По сути, формула описывает, насколько быстро возрастает время выполнения алгоритма с увеличением размера входных данных.
- Время выполнения алгоритмов выражается как «О-большое».
- Время выполнения $O(\log n)$ быстрее $O(n)$, а с увеличением размера списка, в котором происходит поиск значения, оно становится *намного* быстрее.

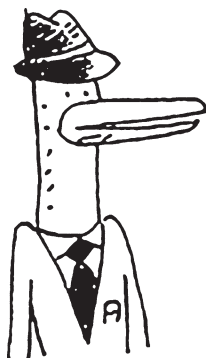
УПРАЖНЕНИЯ

Приведите время выполнения «О-большое» для каждого из следующих сценариев.

- 1.3 Известна фамилия, нужно найти номер в телефонной книге.
- 1.4 Известен номер, нужно найти фамилию в телефонной книге. (Подсказка: вам придется провести поиск по всей книге!)
- 1.5 Нужно прочитать телефоны всех людей в телефонной книге.
- 1.6 Нужно прочитать телефоны всех людей, фамилии которых начинаются с буквы «А». (Вопрос с подвохом! В нем задействованы концепции, которые более подробно рассматриваются в главе 4. Прочитайте ответ — скорее всего, он вас удивит!)

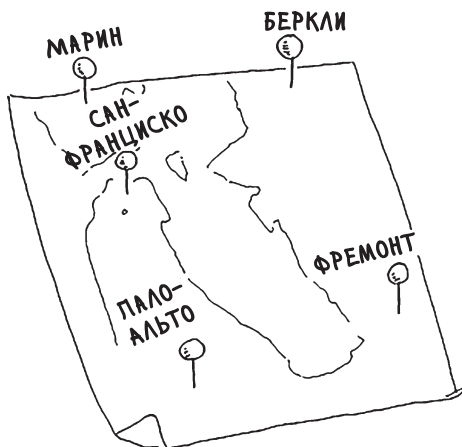
Задача о коммивояжере

Наверное, после прочтения предыдущего раздела вы подумали: «Уж мне-то точно не попадется алгоритм с временем $O(n!)$ ». Ошибаетесь, и я это сейчас докажу! Мы рассмотрим алгоритм с очень, очень плохим временем выполнения. Это известная задача из области теории вычислений, в которой время выполнения растет просто с ужасающей скоростью, и некоторые очень умные люди считают, что с этим ничего не поделать. Она называется *задачей о коммивояжере*.

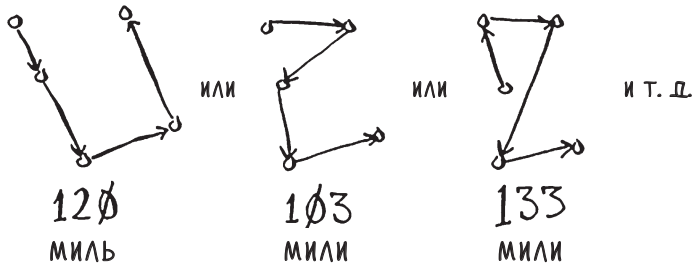


Это коммивояжер.

Он должен объехать 5 городов.



Коммивояжер хочет побывать в каждом из 5 городов так, чтобы при этом проехать минимальное общее расстояние. Одно из возможных решений: нужно перебрать все возможные комбинации порядка объезда городов.



Все расстояния суммируются, после чего выбирается путь с кратчайшим расстоянием. Для 5 городов можно создать 120 перестановок, поэтому решение задачи для 5 городов потребует 120 операций. Для 6 городов количество операций увеличивается до 720 (существует 720 возможных перестановок). А для 7 городов потребуется уже 5040 операций!

ГОРОДА	ОПЕРАЦИИ
6	720
7	5040
8	40320
...	...
15	1307674368000
...	...
30	2652528598121910686363084800000000

Количество операций стремительно растет

В общем случае для вычисления результата при n элементах потребуется $n!$ (n -факториал) операций. А значит, время выполнения составит $O(n!)$ (такое время называется *факториальным*). При любом сколько-нибудь серьезном размере списка количество операций будет просто огромным. Скажем, если вы попытаетесь решить задачу для 100+ городов, сделать это вовремя не удастся — Солнце погаснет раньше.

Какой ужасный алгоритм! Значит, коммивояжер должен найти другое решение, верно? Но у него ничего не получится. Это одна из знаменитых нерешенных задач в области теории вычислений. Для нее не существует известного быстрого алгоритма, и ученые считают, что найти более эффективный алгоритм для этой задачи в принципе невозможно. В лучшем случае для нее можно поискать приближенное решение; за подробностями обращайтесь к главе 10.

Шпаргалка

- Бинарный поиск работает намного быстрее простого.
- Время выполнения $O(\log n)$ быстрее $O(n)$, а с увеличением размера списка, в котором происходит поиск значения, оно становится намного быстрее.
- Скорость алгоритмов не измеряется в секундах.
- Время выполнения алгоритма описывается *ростом* количества операций.
- Время выполнения алгоритмов выражается как «О-большое».

2

Сортировка выбором



В этой главе

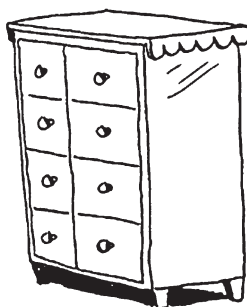
- ✓ Вы познакомитесь с массивами и связанными списками — двумя основными структурами данных, которые используются буквально везде. Мы уже применяли массивы в главе 1 и будем говорить о них почти в каждой главе книги. Массивы чрезвычайно важны, уделите им внимание! Впрочем, иногда вместо массива лучше воспользоваться связанным списком. В этой главе объясняются плюсы и минусы обеих структур данных, чтобы вы могли решить, какой вариант лучше подходит для вашего алгоритма.
- ✓ Вы изучите свой первый алгоритм сортировки. Многие алгоритмы работают только с отсортированными данными. Помните бинарный поиск? Он применяется только к предварительно отсортированному списку. В большинстве языков существуют встроенные алгоритмы сортировки, так что вам редко приходится писать свою версию «с нуля». Однако алгоритм сортировки выбором поможет перейти к алгоритму быстрой сортировки, описанному в главе 4. Алгоритм быстрой сортировки очень важен, и вам будет проще разобраться в нем, если вы уже знаете хотя бы один алгоритм сортировки.

ЧТО НЕОБХОДИМО ЗНАТЬ

Чтобы понять ту часть этой главы, которая относится к анализу эффективности, необходимо понимать смысл понятия «О-большое» и логарифмов. Если вы совершенно не разбираетесь в этих вопросах, лучше вернуться и прочитать главу 1. «О-большое» будет использоваться в оставшихся главах книги.

Как работает память

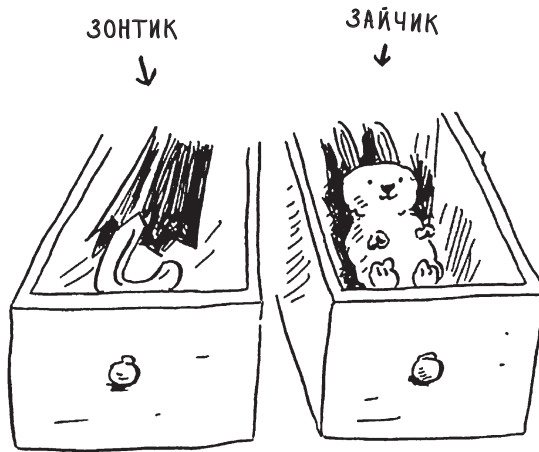
Представьте, что вы пришли в театр и хотите оставить свои личные вещи в гардеробе. Для хранения вещей есть специальные ящики.



В каждом ящике помещается один предмет. Вы хотите сдать на хранение две вещи, поэтому требуете выделить вам два ящика.

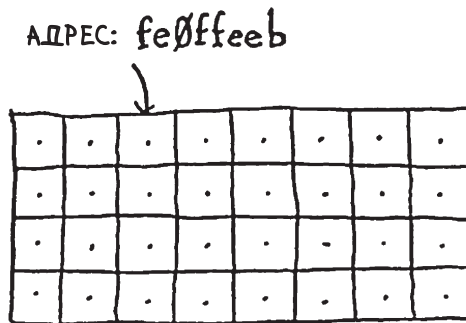


И вы оставляете свои две вещи.



Готово, можно идти на спектакль!

В сущности, именно так работает память вашего компьютера. Она представляет собой нечто вроде огромного шкафа со множеством ящиков, и у каждого ящика есть адрес.



fe0ffeeb — адрес ячейки памяти.

Каждый раз, когда вы хотите сохранить в памяти отдельное значение, вы запрашиваете у компьютера место в памяти, а он выдает адрес для сохранения значения. Если же вам понадобится сохранить несколько элементов, это можно сделать двумя основными способами: воспользоваться массивом или списком. В следующем разделе мы обсудим массивы и списки, их достоинства и недостатки. Не существует единственно верного способа сохранения данных на все случаи жизни, поэтому вы должны знать, в чем отличие разных способов.

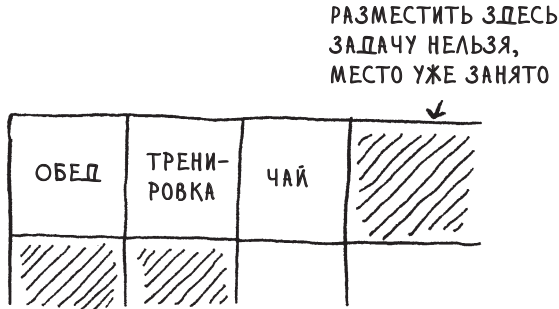
Массивы и связанные списки

Иногда в памяти требуется сохранить список элементов. Предположим, вы пишете приложение для управления текущими делами. Описания задач должны храниться в виде списка в памяти.

Что использовать — массив или связанный список? Для начала попробуем сохранить задачи в массиве, потому что этот способ более понятен. При использовании массива все задачи хранятся в памяти непрерывно (то есть рядом друг с другом).



Теперь предположим, что вы захотели добавить четвертую задачу. Но следующий ящик уже занят — там лежат чужие вещи!



Представьте, что вы пошли в кино с друзьями и нашли места для своей компании, но тут приходит еще один друг, и ему сесть уже некуда. Приходится искать новое место, где смогут разместиться все. В нашем случае вам придется запросить у компьютера другой блок памяти, в котором поместятся все четыре задачи, а потом переместить все свои задачи туда.

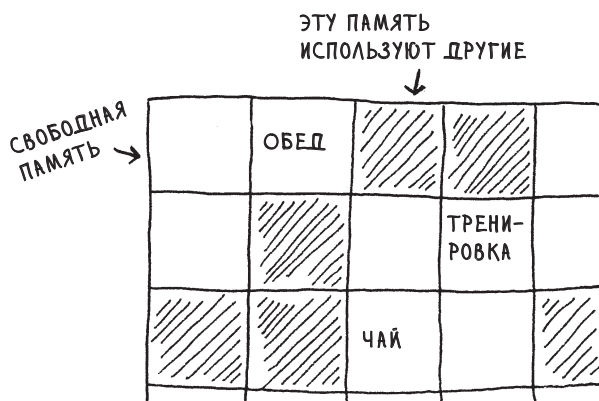
Если вдруг придет еще один друг, места опять не хватит и вам всем придется перемещаться снова! Сплошная суета. Кроме того, добавление новых элементов в массив станет серьезной проблемой. Если свободного места нет и вам каждый раз приходится перемещаться в новую область в памяти, операция добавления нового элемента будет выполняться очень медленно. Простейшее решение — «бронирование мест»: даже если список состоит всего из 3 задач, вы запрашиваете у компьютера место на 10 позиций... просто на всякий случай. Тогда в список можно будет добавить до 10 задач и ничего перемещать не придется. Это неплохое обходное решение, но у него есть пара недостатков.

- Лишнее место может не понадобиться, и тогда память будет расходоваться неэффективно. Вы ее не используете, однако никто другой ее использовать тоже не может.
- Если в список будет добавлено более 10 задач, перемещать все равно придется.

В общем, прием неплохой, но его нельзя назвать идеальным. Связанные списки решают проблему добавления новых элементов.

Связанные списки

При использовании связанного списка элементы могут размещаться где угодно в памяти.



В каждом элементе хранится адрес следующего элемента списка. Таким образом, набор произвольных адресов памяти объединяется в цепочку.



Все как в игре «Найди клад». Вы приходите по первому адресу, там написано: «Следующий элемент находится по адресу 123». Вы идете по адресу 123, там написано: «Следующий элемент находится по адресу 847» и т. д. Добавить новый элемент в связанный список проще простого: просто разместите его по любому адресу памяти и сохраните этот адрес в предыдущем элементе.

Со связанными списками ничего перемещать в памяти не нужно. Также сама собой решается другая проблема: допустим, вы пришли в кино с пятью друзьями. Вы пытаетесь найти место на шестерых, но кинотеатр уже забит, и найти шесть соседних мест невозможно. Нечто похожее происходит и с массивами. Допустим, вы пытаетесь найти для массива блок на 10 000 элементов. В памяти можно найти место для 10 000 элементов, но только не смежное. Для массива не хватает места! При хранении данных в связанном списке вы фактически говорите: «Ладно, тогда садимся на свободные места и смотрим кино». Если необходимое место есть в памяти, вы сможете сохранить данные в связанном списке.

Если связанные списки так хорошо справляются со вставкой, то чем тогда хороши массивы?

Массивы

На сайтах со всевозможными хит-парадами и «первыми десятками» применяется жульническая тактика для увеличения количества просмотров. Вместо того чтобы вывести весь список на одной странице, они размещают по одному элементу на странице и заставляют вас нажимать кнопку Next для перехода к следующему элементу. Например, «Десятка лучших злодеев в сериалах» не выводится на одной странице. Вместо этого вы начинаете с № 10 (Ньюман из «Сайнфелда») и нажимаете Next на каждой странице, пока не доберетесь до № 1 (Густаво Фринг из «Во все тяжкие»). В результате сайту удастся показать вам рекламу на целых 10 страницах, но нажимать Next 9 раз для перехода к первому месту скучно. Было бы гораздо лучше, если бы весь список помещался на одной странице, а вы могли бы просто щелкнуть на имени человека для получения дополнительной информации.



Похожая проблема существует и у связанных списков. Допустим, вы хотите получить последний элемент связанного списка. Просто прочитать нужное значение не удастся, потому что вы не знаете, по какому адресу оно хранится. Вместо этого придется сначала обратиться к элементу № 1 и узнать адрес элемента № 2, потом обратиться к элементу № 2 и узнать адрес элемента № 3... и так далее, пока не доберетесь до последнего элемента. Связанные списки отлично подходят в тех ситуациях, когда данные должны читаться

последовательно: сначала вы читаете один элемент, по адресу переходите к следующему элементу и т. д. Но если вы намерены прыгать по списку туда-сюда, держитесь подальше от связанных списков.

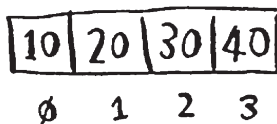
С массивами дело обстоит совершенно иначе. Работая с массивом, вы заранее знаете адрес каждого его элемента. Допустим, массив содержит пять элементов и вы знаете, что он начинается с адреса 00. По какому адресу хранится пятый элемент?



Простейшая математика дает ответ: это адрес 04. Массивы прекрасно подходят для чтения элементов в произвольных позициях, потому что обращение к любому элементу в массиве происходит мгновенно. В связанном списке элементы не хранятся рядом друг с другом, поэтому мгновенно определить позицию i -го элемента в памяти невозможно — нужно обратиться к первому элементу, чтобы получить адрес второго элемента, затем обратиться ко второму элементу для получения адреса третьего — и так далее, пока вы не доберетесь до i -го.

Терминология

Элементы массива пронумерованы, причем нумерация начинается с 0, а не с 1. Например, в этом массиве значение 20 находится в позиции 1.



А значение 10 находится в позиции 0. Неопытных программистов этот факт обычно вводит в ступор. Тем не менее выбор нулевой начальной позиции

упрощает написание кода по работе с массивами, поэтому программисты остановились на этом варианте. Почти во всех языках программирования нумерация элементов массива начинается с 0. Вскоре вы к этому привыкнете.

Позиция элемента называется его *индексом*. Таким образом, вместо того чтобы говорить «Значение 20 находится в позиции 1», правильно сказать «Значение 20 имеет индекс 1». В этой книге термин «индекс» означает то же, что и «позиция».

Ниже приведены примеры времени выполнения основных операций с массивами и списками.

	МАССИВЫ	СПИСКИ
ЧТЕНИЕ	$O(1)$	$O(n)$
ВСТАВКА	$O(n)$	$O(1)$

$O(n)$ = ЛИНЕЙНОЕ ВРЕМЯ
 $O(1)$ = ПОСТОЯННОЕ ВРЕМЯ

Вопрос: почему вставка элемента в массив требует времени $O(n)$? Предположим, вы хотите вставить элемент в начало массива. Как бы вы это сделали? Сколько времени на это потребуется? Ответы на эти вопросы вы найдете в следующем разделе!

УПРАЖНЕНИЕ

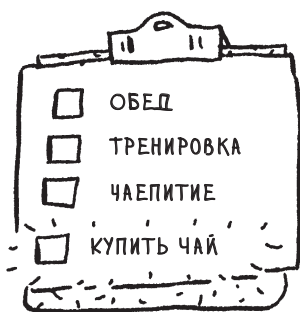
2.1 Допустим, вы строите приложение для управления финансами.

1. ПРОДУКТЫ
2. КИНО
3. ВЕЛОСИПЕДНЫЙ КЛУБ

Ежедневно вы записываете все свои траты. В конце месяца вы анализируете расходы и вычисляете, сколько денег было потрачено. При работе с данными выполняется множество операций вставки и относительно немного операций чтения. Какую структуру использовать — массив или список?

Вставка в середину списка

Предположим, вы решили, что список задач должен больше напоминать календарь. Прежде данные добавлялись только в конец списка, а теперь они должны добавляться в порядке их выполнения.

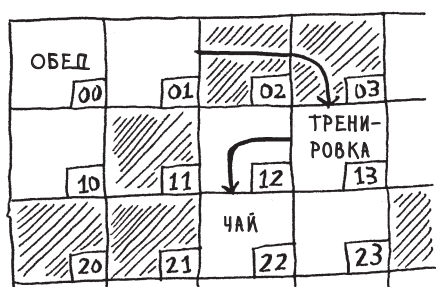


Неупорядоченный

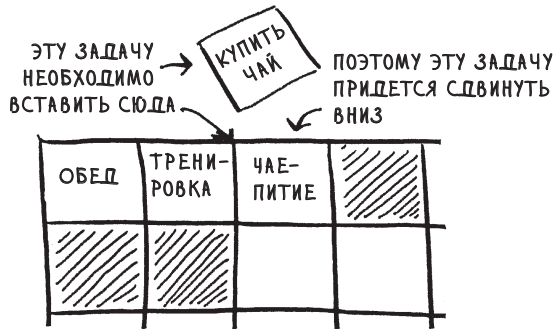


Упорядоченный

Что лучше подойдет для вставки элементов в середину: массивы или списки? Со списком задача решается изменением указателя в предыдущем элементе.



А при работе с массивом придется сдвигать все остальные элементы.



А если свободного места не осталось, все данные придется скопировать в новую область памяти! В общем, списки лучше подходят для вставки элементов в середину.

УКАЗАТЕЛИ

Я часто говорю о том, что каждый элемент связанного списка указывает на следующий элемент списка. Но как именно он это делает? С помощью указателей.



В каждом элементе связанного списка выделяется некоторый объем памяти для хранения адреса следующего элемента. Эта часть элемента называется указателем.

Периодически вы будете встречать этот термин, особенно если вы пишете на низкоуровневом языке, таком как С. Поэтому полезно знать, что он означает.

Удаление

Что, если вы захотите удалить элемент? И снова список лучше подходит для этой операции, потому что в нем достаточно изменить указатель в предыдущем элементе. В массиве при удалении элемента все последующие элементы нужно будет сдвинуть вверх.

В отличие от вставки, удаление возможно всегда. Попытка вставки может быть неудачной, если в памяти не осталось свободного места. С удалением подобных проблем не бывает.

Ниже приведены примеры времени выполнения основных операций с массивами и связанными списками.

	МАССИВЫ	СПИСКИ
ЧТЕНИЕ	$O(1)$	$O(n)$
ВСТАВКА	$O(n)$	$O(1)$
УДАЛЕНИЕ	$O(n)$	$O(1)$

Заметим, что вставка и удаление выполняются за время $O(1)$ только в том случае, если вы можете мгновенно получить доступ к удаляемому элементу. На практике обычно сохраняются ссылки на первый и последний элементы связанного списка, поэтому время удаления этих элементов составит всего $O(1)$.

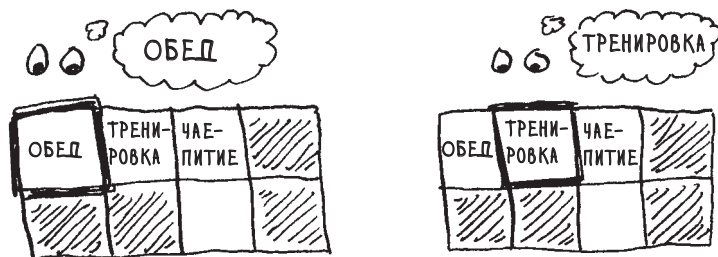
Какая структура данных используется чаще: массивы или списки?

Массивы используются часто, поскольку имеют множество преимуществ перед связанными списками. Во-первых, их проще читать. Во-вторых, массивы поддерживают произвольный доступ.

Всего существуют два вида доступа: произвольный и последовательный. При последовательном доступе элементы читаются по одному, начиная с первого. Связанные списки поддерживают только последовательный

доступ. Если вы захотите прочитать 10-й элемент связанного списка, вам придется прочитать первые 9 элементов и перейти по ссылкам к 10-му элементу. Многие реальные ситуации требуют произвольного доступа, поэтому массивы часто применяются на практике.

Однако даже безотносительно произвольного доступа массивы работают быстрее, поскольку могут использовать кэширование. Возможно, вы представляете себе процесс чтения как восприятие одного элемента за раз.



Но на самом деле компьютеры читают целыми разделами, поскольку это значительно ускоряет переход к следующему элементу.

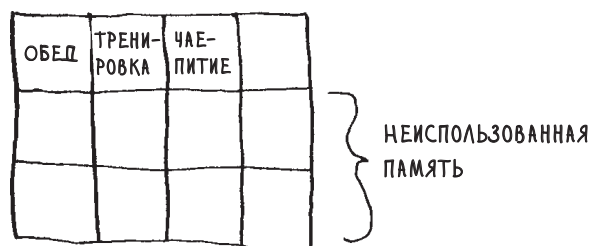


С массивами можно сделать то, что не получится со связанными списками. С помощью массива можно прочитать целый раздел элементов. В связанном списке местоположение следующего элемента неизвестно. Придется прочитать элемент, узнать, где находится следующий, а затем прочитать и его.

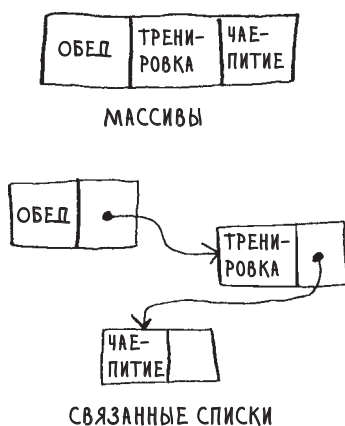
Таким образом, массивы обеспечивают не только произвольный, но и более быстрый последовательный доступ!

Массивы проще читать. А как насчет эффективного использования памяти? Помните, я говорил, что в случае с массивами обычно выделяют больше

памяти, чем нужно, и если в конечном итоге она не задействуется, то ее ресурс расходуется впустую?



На самом деле объем такой неиспользованной памяти получается не очень большим. С другой стороны, когда вы применяете связанный список, вы выделяете дополнительную память для каждого элемента, поскольку в ней сохраняется адрес следующего элемента. Таким образом, связанные списки занимают больше места, если их элементы относительно невелики. На рисунке показано, сколько места занимает одинаковый объем информации в массиве и в связанном списке. Видно, что связанный список гораздо больше.



Конечно, если элементы массива крупные, то даже один слот неиспользованной памяти может оказаться очень большим, а объем дополнительной памяти, используемой для хранения указателей, по сравнению с ним будет незначительным.

Таким образом, массивы используются чаще, чем связанные списки, за исключением особых случаев.

УПРАЖНЕНИЯ

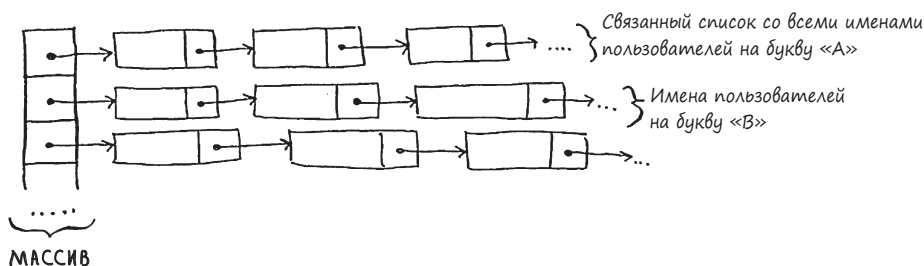
- 2.2** Допустим, вы пишете приложение для приема заказов от посетителей ресторана. Приложение должно хранить список заказов. Официанты добавляют заказы в список, а повара читают заказы из списка и выполняют их. Заказы образуют очередь: официанты добавляют заказы в конец очереди, а повар берет первый заказ из очереди и начинает готовить.



Какую структуру данных вы использовали бы для реализации этой очереди: массив или связанный список? (Подсказка: связанные списки хорошо подходят для вставки/удаления, а массивы — для произвольного доступа к элементам. Что из этого понадобится в данном случае?)

- 2.3** Проведем мысленный эксперимент. Допустим, Facebook хранит список имен пользователей. Когда кто-то пытается зайти на сайт Facebook, система пытается найти имя пользователя. Если имя входит в список имен зарегистрированных пользователей, то вход разрешается. Пользователи приходят на Facebook достаточно часто, поэтому поиск по списку имен пользователей будет выполняться часто. Будем считать, что Facebook использует бинарный поиск для поиска в списке. Бинарному поиску необходим произвольный доступ — алгоритм должен мгновенно обратиться к среднему элементу текущей части списка. Зная это обстоятельство, как бы вы реализовали список пользователей: в виде массива или в виде связанного списка?
- 2.4** Пользователи также довольно часто создают новые учетные записи на Facebook. Предположим, вы решили использовать массив для хранения списка пользователей. Какими недостатками обладает массив для выполнения вставки? Допустим, вы используете бинарный поиск для нахождения учетных данных. Что произойдет при добавлении новых пользователей в массив?

2.5 В действительности Facebook не использует ни массив, ни связанный список для хранения информации о пользователях. Рассмотрим гибридную структуру данных: массив связанных списков. Имеется массив из 26 элементов. Каждый элемент содержит ссылку на связанный список. Например, первый элемент массива указывает на связанный список всех имен пользователей, начинающихся на букву «А». Второй элемент указывает на связанный список всех имен пользователей, начинающихся на букву «В», и т. д.

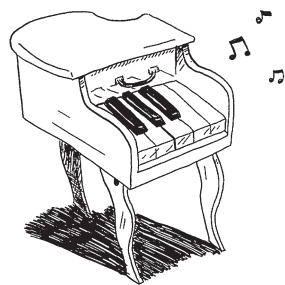


Предположим, пользователь с именем «Adit В» регистрируется на Facebook и вы хотите добавить его в список. Вы обращаетесь к элементу 1 массива, находите связанный список элемента 1 и добавляете «Adit В» в конец списка. Теперь предположим, что зарегистрировать нужно пользователя «Zakhir Н». Вы обращаетесь к элементу 26, который содержит связанный список всех имен, начинающихся с «Z», и проверяете, присутствует ли «Zakhir Н» в этом списке.

Теперь сравните эту гибридную структуру данных с массивами и связанными списками. Будет ли она быстрее или медленнее каждой исходной структуры при поиске и вставке? Приводить «О-большое» не нужно, просто выберите одно из двух: быстрее или медленнее.

Сортировка выбором

А теперь объединим все, что вы узнали, во втором алгоритме: сортировке выбором. Чтобы освоить этот алгоритм, вы должны понимать, как работают массивы и списки и «О-большое». Допустим, у вас на компьютере записана музыка и для каждого исполнителя хранится счетчик воспроизведений.



~ 🎵 ~	СЧЕТЧИК ВОСПРОИЗВЕДЕНИЙ
RADIOHEAD	156
KISHORE KUMAR	141
THE BLACK KEYS	35
NEUTRAL MILK HOTEL	94
BECK	88
THE STROKES	61
WILCO	111

Вы хотите отсортировать список по убыванию счетчика воспроизведений, чтобы самые любимые исполнители стояли на первых местах. Как это сделать?

Одно из возможных решений — пройти по списку и найти исполнителя с наибольшим количеством воспроизведений. Этот исполнитель добавляется в новый список.



Потом то же самое происходит со следующим по количеству воспроизведений исполнителем.

~♪~	СЧЕТЧИК ВОСПРОИЗВЕДЕНИЙ	→	♪ СПИСОК ♪	СЧЕТЧИК ВОСПРОИЗВЕДЕНИЙ
			RADIOHEAD	156
KISHORE KUMAR	141		KISHORE KUMAR	141
THE BLACK KEYS	35			
NEUTRAL MILK HOTEL	94			
BECK	88			
THE STROKES	61			
WILCO	111			

1. СЛЕДУЮЩИЙ ИСПОЛНИТЕЛЬ ПО КОЛИЧЕСТВУ ВОСПРОИЗВЕДЕНИЙ — KISHORE KUMAR

2. ПОЭТОМУ ОН СЛЕДУЮЩИМ ДОБАВЛЯЕТСЯ В НОВЫЙ СПИСОК

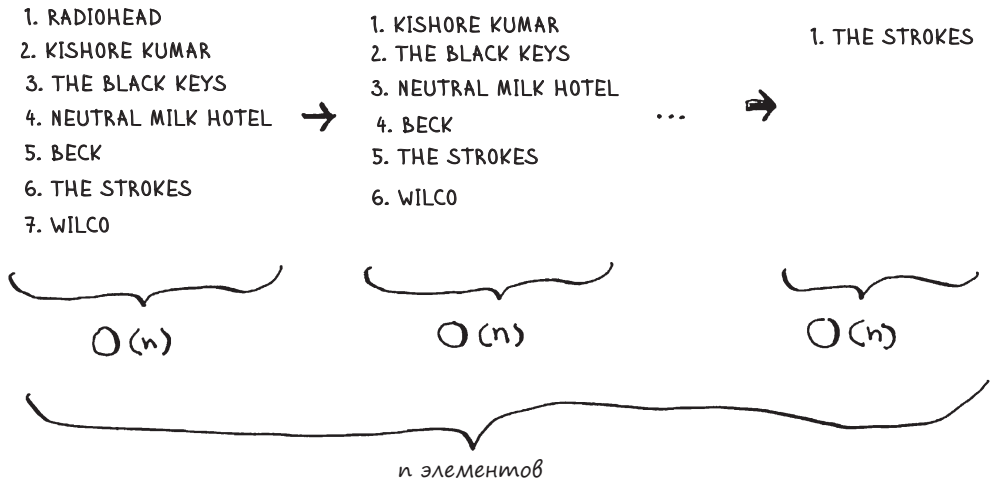
Продолжая действовать так, мы получаем отсортированный список.

~♪~	СЧЕТЧИК ВОСПРОИЗВЕДЕНИЙ
RADIOHEAD	156
KISHORE KUMAR	141
WILCO	111
NEUTRAL MILK HOTEL	94
BECK	88
THE STROKES	61
THE BLACK KEYS	35

А теперь попробуем оценить происходящее с точки зрения теории вычислений и посмотрим, сколько времени будут занимать операции. Напомним, что время $O(n)$ означает, что вы по одному разу обращаетесь к каждому элементу списка. Например, при простом поиске по списку исполнителей каждый исполнитель будет проверен один раз.

1. RADIOHEAD
 2. KISHORE KUMAR
 3. THE BLACK KEYS
 4. NEUTRAL MILK HOTEL
 5. BECK
 6. THE STROKES
 7. WILCO
- } n
элементов

Чтобы найти исполнителя с наибольшим значением счетчика воспроизведения, необходимо проверить каждый элемент в списке. Как вы уже видели, это делается за время $O(n)$. Итак, имеется операция, выполняемая за время $O(n)$, и ее необходимо выполнить n раз:



Все это требует времени $O(n \times n)$, или $O(n^2)$.

Алгоритмы сортировки очень полезны. Например, вы можете отсортировать:

- имена в телефонной книге;
- даты путешествий;
- сообщения электронной почты (от новых к старым).

УМЕНЬШЕНИЕ КОЛИЧЕСТВА ПРОВЕРЯЕМЫХ ЭЛЕМЕНТОВ

При каждом выполнении операций количество элементов, которые нужно проверить, сокращается. Со временем все сведется к проверке всего одного элемента. Почему же время выполнения все равно оценивается как $O(n^2)$? Это хороший вопрос, и ответ на него связан с ролью констант в «О-большом». Тема будет более подробно рассмотрена в главе 4, но я кратко объясню суть. Вы правы, вам действительно не нужно каждый раз проверять весь список из n элементов. Сначала проверяются n элементов, потом $n - 1, n - 2 \dots 2, 1$. В среднем проверяется список из $\frac{1}{2} \times n$ элементов. Его время выполнения составит $O(n \times \frac{1}{2} \times n)$. Однако константы (такие, как $\frac{1}{2}$) в «О-большом» игнорируются (еще раз: за полным обсуждением обращайтесь к главе 4), поэтому мы просто используем $O(n \times n)$, или $O(n^2)$.

Алгоритм сортировки выбором легко объясняется, но медленно работает. Быстрая сортировка — эффективный алгоритм сортировки, который выполняется за время $O(n \log n)$. Но мы займемся этой темой в главе 4!

Пример кода

Мы не будем приводить код сортировки музыкального списка, но написанный ниже код делает нечто очень похожее: он выполняет сортировку массива по возрастанию. Напишем функцию для поиска наименьшего элемента массива:

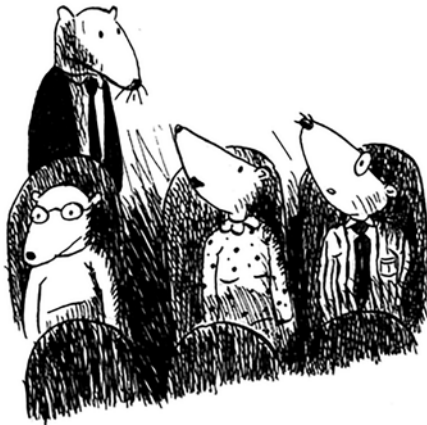
```
def findSmallest(arr):
    smallest = arr[0]
    smallest_index = 0
    for i in range(1, len(arr)):
        if arr[i] < smallest:
            smallest = arr[i]
            smallest_index = i
    return smallest_index
```

Для хранения наименьшего значения
 Для хранения индекса наименьшего значения

Теперь на основе этой функции можно написать функцию сортировки выбором:

```
def selectionSort(arr):  <..... Сортирует массив
    newArr = []
    for i in range(len(arr)):
        smallest = findSmallest(arr) <..... Находит наименьший элемент в массиве
        newArr.append(arr.pop(smallest))      и добавляет его в новый массив
    return newArr

print selectionSort([5, 3, 6, 2, 10])
```



Шпаргалка

- Память компьютера напоминает огромный шкаф с ящиками.
- Если вам потребуется сохранить набор элементов, воспользуйтесь массивом или списком.
- В массиве все элементы хранятся в памяти рядом друг с другом.
- В связанном списке элементы распределяются в произвольных местах памяти, при этом в каждом элементе хранится адрес следующего элемента.
- Массивы обеспечивают быстрое чтение.
- Списки обеспечивают быструю вставку и удаление.

3

Рекурсия



В этой главе

- ✓ Вы узнаете, что такое рекурсия — метод программирования, используемый во многих алгоритмах. Это важная концепция для понимания дальнейших глав книги.
- ✓ Вы научитесь разбивать задачи на базовый и рекурсивный случаи. В стратегии «разделяй и властвуй» (глава 4) эта простая концепция используется для решения более сложных задач.

Эта глава мне самому очень нравится, потому что в ней рассматривается *рекурсия* — элегантный метод решения задач. Рекурсия относится к числу моих любимых тем, но вызывает у людей противоречивые чувства. Они либо обожают ее, либо ненавидят, либо ненавидят, пока не полюбят через пару-тройку лет. Лично я отношусь к третьему лагерю. Чтобы вам было проще освоить эту тему, я дам несколько советов.

- Глава содержит множество примеров кода. Самостоятельно выполните этот код и посмотрите, как он работает.
- Мы будем рассматривать рекурсивные функции. Хотя бы один раз возьмите бумагу и карандаш и разберите, как работает рекурсивная функция: «Так, я передаю функции `factorial` значение 5, потом возвращаю

управление и передаю значение 4 функции `factorial`, которая...» и т. д. Такой разбор поможет вам понять, как работает рекурсивная функция.

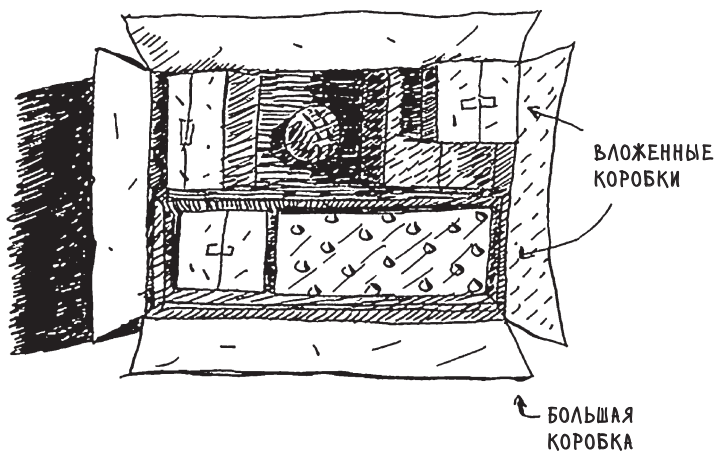
В этой главе также приводится большое количество псевдокода. *Псевдокод* представляет собой высокоуровневое описание решаемой задачи. Он записывается в форме, похожей на программный код, но в большей степени напоминает естественный язык.

Рекурсия

Допустим, вы разбираете чулан своей бабушки и натываетесь на загадочный запертый чемодан.



Бабушка говорит, что ключ к чемодану, скорее всего, лежит в коробке.



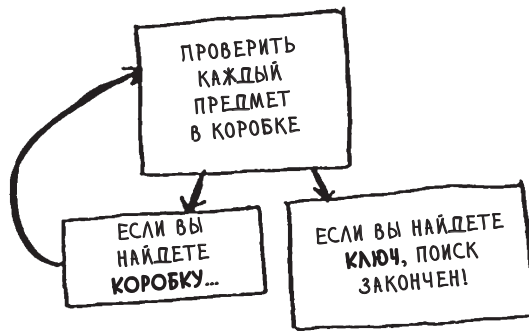
В коробке лежат другие коробки, а в них лежат маленькие коробочки. Ключ находится где-то там. Какой алгоритм поиска ключа предложите вы? Подумайте над алгоритмом, прежде чем продолжить чтение.

Одно из решений может выглядеть так:



1. Сложить все коробки в кучу.
2. Взять коробку и открыть.
3. Если внутри лежит коробка, добавить ее в кучу для последующего поиска.
4. Если внутри лежит ключ, поиск закончен!
5. Повторить.

Есть и альтернативное решение.



1. Просмотреть содержимое коробки.
2. Если вы найдете коробку, вернуться к шагу 1.
3. Если вы найдете ключ, поиск закончен!

Какое решение кажется вам более простым? Первое решение можно построить на цикле `while`. Пока куча коробок не пуста, взять очередную коробку и проверить ее содержимое:

```
def look_for_key(main_box):
    pile = main_box.make_a_pile_to_look_through()
    while pile is not empty:
        box = pile.grab_a_box()
        for item in box:
            if item.is_a_box():
                pile.append(item)
            elif item.is_a_key():
                print "found the key!"
```

Второй способ основан на рекурсии. *Рекурсией* называется вызов функцией самой себя. Второе решение на псевдокоде может выглядеть так:

```
def look_for_key(box):
    for item in box:
        if item.is_a_box():
            look_for_key(item)
        elif item.is_a_key():
            print "found the key!"
```

←..... Рекурсия!

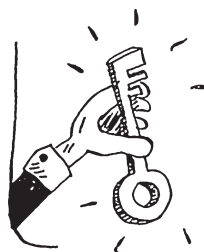
Оба решения делают одно и то же, но второе решение кажется мне более понятным. Применение рекурсии не ускоряет работу программы: более того, решение с циклами иногда работает быстрее. Мне нравится одна цитата Ли Колдуэлла с сайта Stack Overflow: «Циклы могут ускорить работу программы. Рекурсия может ускорить работу программиста. Выбирайте, что важнее в вашей ситуации!»¹

Рекурсия используется во многих нужных алгоритмах, поэтому важно понимать эту концепцию.

Базовый случай и рекурсивный случай

Так как рекурсивная функция вызывает сама себя, программисту легко ошибиться и написать функцию так, что возникнет бесконечный цикл. Предположим, вы хотите написать функцию для вывода обратного отсчета:

```
> 3...2...1
```

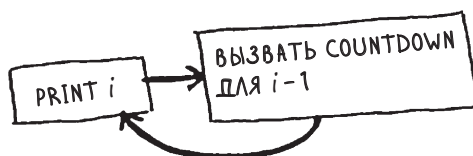


Ее можно записать в рекурсивном виде:

```
def countdown(i):
    print(i)
    countdown(i-1)

countdown(3)
```

Введите этот код и выполните его. И тут возникает проблема: эта функция выполняется бесконечно!



**Бесконечный
цикл**

¹ <http://stackoverflow.com/a/72694/139117>

> 3...2...1...0...-1...-2...

Чтобы прервать выполнение сценария, нажмите Ctrl+C.

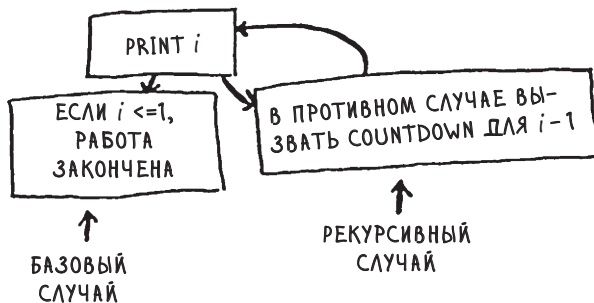
Когда вы пишете рекурсивную функцию, в ней необходимо указать, в какой момент следует прервать рекурсию. Вот почему *каждая рекурсивная функция состоит из двух частей: базового случая и рекурсивного случая*. В рекурсивном случае функция вызывает сама себя. В базовом случае функция себя не вызывает... чтобы предотвратить заикливание.

Добавим базовый случай в функцию countdown:

```
def countdown(i):
    print i
    if i <= 1:  ←..... Базовый случай
        return
    else:       ←..... Рекурсивный случай
        countdown(i-1)
```

countdown(3)

Теперь функция работает, как было задумано. Это выглядит примерно так:



Стек

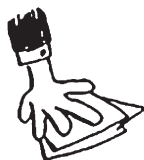
В этом разделе рассматривается *стек вызовов*. Концепция стека вызовов играет важную роль в программировании вообще; кроме того, ее важно понимать при использовании рекурсии.



Предположим, вы устраиваете вечеринку с барбекю. Вы составляете список задач и записываете дела на листках.



Помните, когда мы рассматривали массивы и списки, у вас тоже был список задач? Задачи, то есть элементы списка, можно было добавлять и удалять в произвольных позициях списка. Стопка листов работает куда проще. Новые (вставленные) элементы добавляются в начало списка, то есть на верх стопки. Читается только верхний элемент, и он исключается из списка. Таким образом, список задач поддерживает всего два действия: *занесение* (вставка) и *извлечение* (удаление из списка и чтение.)

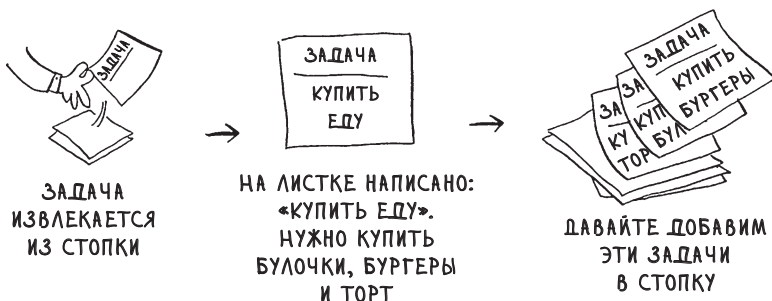


ЗАНЕСЕНИЕ
(НОВЫЙ ЭЛЕМЕНТ
ДОБАВЛЯЕТСЯ
НА ВЕРХ СТОПКИ)



ИЗВЛЕЧЕНИЕ
(ВЕРХНИЙ ЭЛЕМЕНТ
УДАЛЯЕТСЯ ИЗ СТОПКИ
И ЧИТАЕТСЯ)

Посмотрим, как работает список задач:



Такая структура данных называется *стеком*. Стек — простая структура данных. А теперь самое неожиданное: все это время вы пользовались стеком, не подозревая об этом!

Стек вызовов

Во внутренней работе вашего компьютера используется стек, называемый *стеком вызовов*. Давайте посмотрим, как он работает. Предположим, имеется простая функция:

```
def greet(name):  
    print "hello, " + name + "!"  
    greet2(name)  
    print "getting ready to say bye..."  
    bye()
```

Эта функция приветствует вас, после чего вызывает две другие функции. Вот эти две функции:

```
def greet2(name):  
    print "how are you, " + name + "?"  
def bye():  
    print "ok bye!"
```

Разберемся, что происходит при вызове функции.

ПРИМЕЧАНИЕ

Для простоты я привожу только вызовы функций `greet`, `greet2` и `bye` и опускаю вызовы функции `print`.

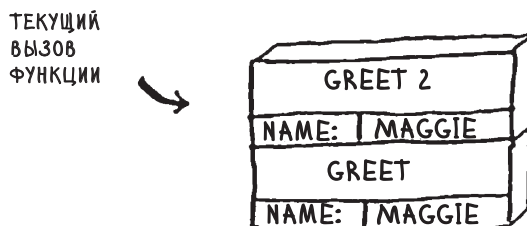
Предположим, в программе используется вызов `greet("maggie")`. Сначала ваш компьютер выделяет блок памяти для этого вызова функции.



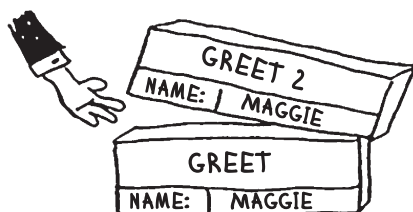
Затем эта память используется. Переменной `name` присваивается значение `"maggie"`; оно должно быть сохранено в памяти.



Каждый раз, когда вы вызываете функцию, компьютер сохраняет в памяти значения всех переменных для этого вызова. Далее выводится приветствие `hello, maggie!`, после чего следует второй вызов `greet2("maggie")`. И снова компьютер выделяет блок памяти для вызова функции.

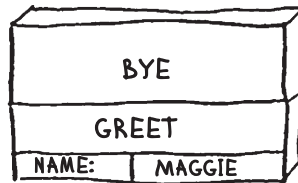


Ваш компьютер объединяет эти блоки в стек. Второй блок создается над первым. Вы выводите сообщение `how are you, maggie?`, после чего возвращаете управление из вызова функции. Когда это происходит, блок на вершине стека извлекается из него.

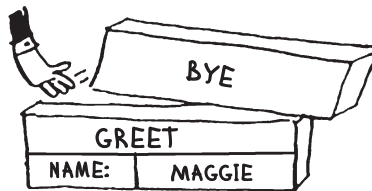


Теперь верхний блок в стеке относится к функции `greet`; это означает, что вы вернулись к функции `greet`. При вызове функции `greet2` функция `greet` еще не была завершена. Здесь-то и скрывается истинный смысл этого раздела: *когда вы вызываете функцию из другой функции, вызывающая функция*

приостанавливается в частично завершённом состоянии. Все значения переменных этой функции остаются в памяти. А когда выполнение функции `greet2` будет завершено, вы вернётесь к функции `greet` и продолжите её выполнение с того места, где оно прервалось. Сначала выводится сообщение `getting ready to say bye...`, после чего вызывается функция `bye`.



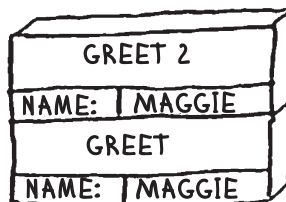
Блок для этой функции добавляется на вершину стека. Далее выводится сообщение `ok bye!` с выходом из вызова функции.



Управление снова возвращается функции `greet`. Делать больше нечего, так что управление возвращается и из функции `greet`. Этот стек, в котором сохранялись переменные разных функций, называется *стеком вызовов*.

УПРАЖНЕНИЕ

3.1 Предположим, имеется стек вызовов следующего вида:



Что можно сказать о текущем состоянии программы на основании этого стека вызовов?

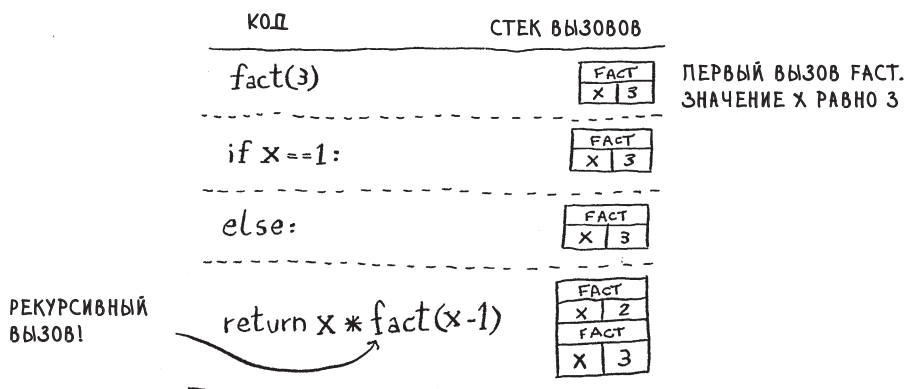
А теперь посмотрим, как работает стек вызовов с рекурсивными функциями.

Стек вызовов с рекурсией

Рекурсивные функции тоже используют стек вызовов! Посмотрим, как это делается, на примере функции вычисления факториала. Вызов `fact(5)` вычисляет факториал 5! и определяется следующим образом: $5! = 5 \times 4 \times 3 \times 2 \times 1$. По тому же принципу `fact(3)` соответствует $3 \times 2 \times 1$. Рекурсивная функция для вычисления факториала числа выглядит так:

```
def fact(x):
    if x == 1:
        return 1
    else:
        return x * fact(x-1)
```

В программу включается вызов `fact(3)`. Проанализируем этот вызов строку за строкой и посмотрим, как изменяется стек вызовов. Стоит напомнить, что верхний блок в стеке сообщает, какой вызов `fact` является текущим.



СЕЙЧАС ТЕКУЩИМ
СТАЛ ВТОРОЙ ВЫЗОВ
FACT. ЗНАЧЕНИЕ X
РАВНО 2

if $x == 1$:

FACT	
X	2
FACT	
X	3

ВЕРХНИЙ ВЫЗОВ ФУНКЦИИ - ТОТ, КОТОРЫЙ В ДАННЫЙ МОМЕНТ ЯВЛЯЕТСЯ ТЕКУЩИМ

else:

FACT	
X	2
FACT	
X	3

В ОБОИХ ВЫЗОВАХ СУЩЕСТВУЕТ ПЕРЕМЕННАЯ С ИМЕНЕМ X, КОТОРАЯ ИМЕЕТ В ЭТИХ ВЫЗОВАХ РАЗНЫЕ ЗНАЧЕНИЯ

return $x * \text{fact}(x-1)$

FACT	
X	1
FACT	
X	2
FACT	
X	3

ОБРАТИТЬСЯ К ЗНАЧЕНИЮ X ЭТОГО ВЫЗОВА ВНУТРИ ЭТОГО ВЫЗОВА НЕВОЗМОЖНО, И НАОБОРОТ

if $x == 1$:

FACT	
X	1
FACT	
X	2
FACT	
X	3

ОГО, ЭТО УЖЕ ТРЕТИЙ ВЫЗОВ, ПРИЧЕМ НИ ОДИН ВЫЗОВ ДО СИХ ПОР ТАК И НЕ ЗАВЕРШИЛСЯ!

ВОЗВРАЩАЕТ 1

FACT	
X	1
FACT	
X	2
FACT	
X	3

ПЕРВЫЙ БЛОК, КОТОРЫЙ БУДЕТ ИЗВЛЕЧЕН ИЗ СТЕКА; ЭТО ОЗНАЧАЕТ, ЧТО ИМЕННО ЭТОТ ВЫЗОВ ПЕРВЫМ ВЕРНЕТ УПРАВЛЕНИЕ

ВОЗВРАЩАЕТ 1

ВЫЗОВ ФУНКЦИИ, ТОЛЬКО ЧТО ВЕРНУВШИЙ УПРАВЛЕНИЕ

return $x * \text{fact}(x-1)$

ЗНАЧЕНИЕ X РАВНО 2

FACT	
X	2
FACT	
X	3

ВОЗВРАЩАЕТ 2

return $x * \text{fact}(x-1)$

x is 3

ЭТОТ ВЫЗОВ ВЕРНУЛ 2

FACT	
X	3

ВОЗВРАЩАЕТ 6

Здесь важно, что каждый вызов создает собственную копию x . Обратиться к переменной x , принадлежащей другой функции, невозможно.

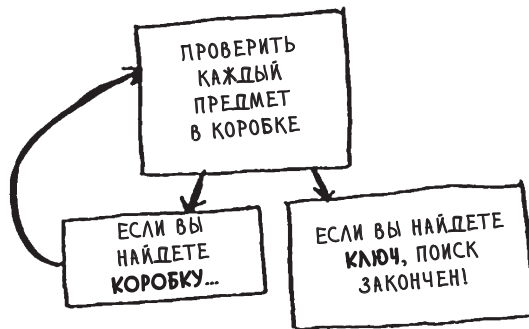
Стек играет важную роль в рекурсии. В начальном примере были представлены два решения поиска ключа. Вспомните, как выглядел первый:



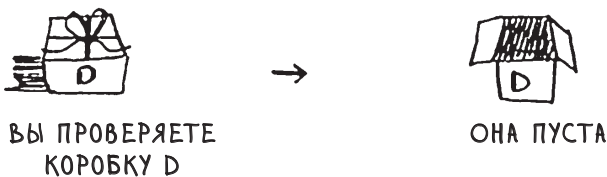
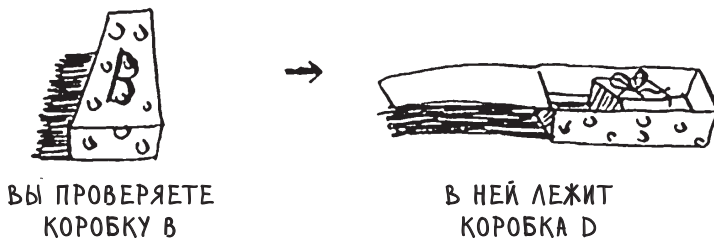
В этом случае все коробки лежат в одном месте и вы всегда знаете, в каких коробках еще нужно искать ключ.



Но в рекурсивном решении никакой кучи не существует.



Если кучи нет, то как ваш алгоритм узнает, в каких коробках еще нужно искать? Пример:



К этому моменту стек вызовов выглядит примерно так:



«Куча коробок» хранится в стеке! Это стек незавершенных вызовов функции, каждый из которых ведет собственный незаконченный список коробок для поиска. Стек в данном случае особенно удобен, потому что вам не нужно отслеживать коробки самостоятельно — стек делает это за вас.

Стек удобен, но у него есть своя цена: сохранение всей промежуточной информации может привести к значительным затратам памяти. Каждый вызов функции занимает немного памяти, но если стек станет слишком «высоким», это будет означать, что ваш компьютер сохраняет информацию по очень многим вызовам. На этой стадии есть два варианта.

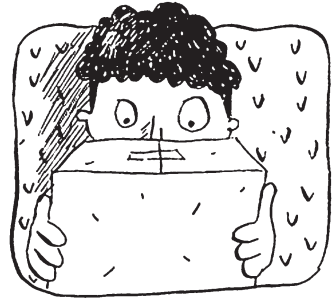
- Переписать код с использованием цикла.
- Иногда можно воспользоваться так называемой *хвостовой рекурсией*. Это непростая тема, которая выходит за рамки книги. Вдобавок она поддерживается далеко не во всех языках.

УПРАЖНЕНИЕ

- 3.2** Предположим, вы случайно написали рекурсивную функцию, которая бесконечно вызывает саму себя. Как вы уже видели, компьютер выделяет память в стеке при каждом вызове функции. А что произойдет со стеком при бесконечном выполнении рекурсии?

Шпаргалка

- Когда функция вызывает саму себя, это называется рекурсией.
- В каждой рекурсивной функции должно быть два случая: базовый и рекурсивный.
- Стек поддерживает две операции: занесение и извлечение элементов.
- Все вызовы функций сохраняются в стеке вызовов.
- Если стек вызовов станет очень большим, он займет слишком много памяти.



4

Быстрая сортировка



В этой главе

- ✓ Вы узнаете о стратегии «разделяй и властвуй». Случается так, что задача, над которой вы трудитесь, не решается ни одним из известных вам алгоритмов. Столкнувшись с такой задачей, хороший программист не сдается. У него существует целый арсенал приемов, которые он пытается использовать для получения решения. «Разделяй и властвуй» — первая общая стратегия, с которой вы познакомитесь.
- ✓ Далее рассматривается быстрая сортировка — элегантный алгоритм сортировки, часто применяемый на практике. Алгоритм быстрой сортировки использует стратегию «разделяй и властвуй».

Предыдущая глава была посвящена рекурсии. В этой главе вы воспользуетесь новыми знаниями для решения практических задач. Мы исследуем принцип «разделяй и властвуй», хорошо известный рекурсивный метод решения задач.

В этой главе мы постепенно добираемся до полноценных алгоритмов. В конце концов, алгоритм не особенно полезен, если он способен решать задачу только одного типа, — «разделяй и властвуй» помогает выработать новый подход к решению задач. Это всего лишь еще один инструмент в ва-

шем арсенале. Столкнувшись с новой задачей, не впадайте в ступор. Вместо этого спросите себя: «А нельзя ли решить эту задачу, применив стратегию “разделяй и властвуй”?»

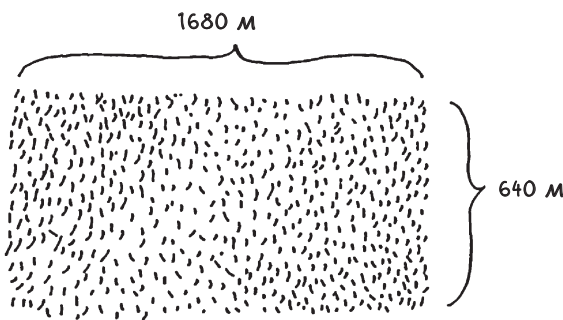
К концу этой главы вы освоите свой первый серьезный алгоритм «разделяй и властвуй»: *быструю сортировку*. Этот алгоритм сортировки работает намного быстрее сортировки выбором (о которой рассказывалось в главе 2). Он является хорошим примером элегантного кода.

«Разделяй и властвуй»

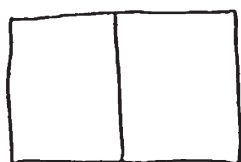
Возможно, вы не сразу поймете суть стратегии «разделяй и властвуй», поэтому мы рассмотрим три примера. Сначала я приведу наглядный пример. Потом мы разберем пример кода, который выглядит не так красиво, но, пожалуй, воспринимается проще. В завершение будет рассмотрена быстрая сортировка — алгоритм сортировки, использующий стратегию «разделяй и властвуй».



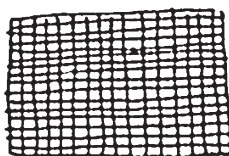
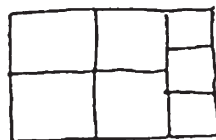
Представьте, что вы фермер, владеющий земельным участком.



Вы хотите равномерно разделить землю на одинаковые *квадратные* участки. Участки должны быть настолько большими, насколько это возможно, так что ни одно из следующих решений не подойдет.



НЕ КВАДРАТНЫЕ

СЛИШКОМ
МАЛЕНЬКИЕВСЕ УЧАСТКИ
ДОЛЖНЫ БЫТЬ
ОДИНАКОВЫМИ

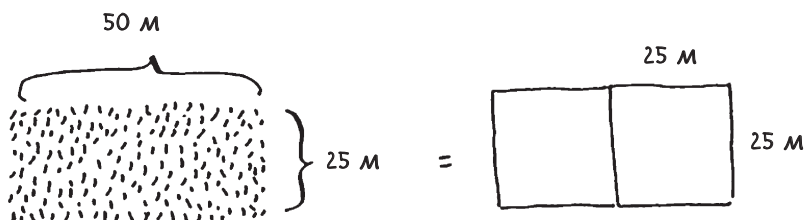
Как определить наибольший размер квадрата для участка? Воспользуйтесь стратегией «разделяй и властвуй»! Алгоритмы на базе этой стратегии являются рекурсивными.

Решение задачи методом «разделяй и властвуй» состоит из двух шагов.

1. Сначала определяется базовый случай. Это должен быть простейший случай из всех возможных.
2. Задача делится или сокращается до тех пор, пока не будет сведена к базовому случаю.

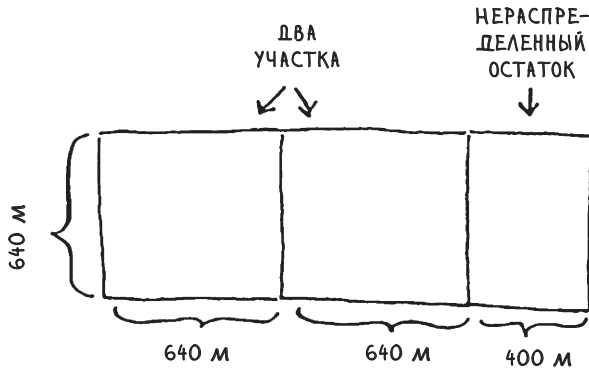
А теперь воспользуемся стратегией «разделяй и властвуй» для поиска решения этой задачи. Каков самый большой размер квадрата, который может использоваться?

Для начала нужно определить базовый случай. Самая простая ситуация — если длина одной стороны кратна длине другой стороны.

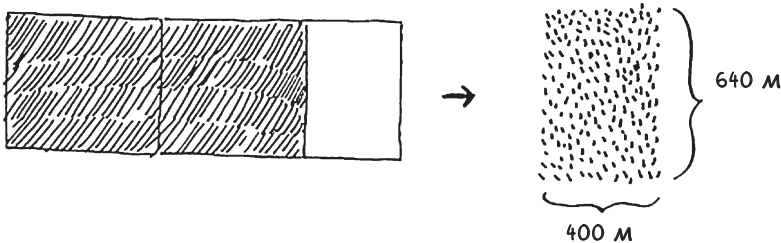


Предположим, длина одной стороны составляет 25 м, а длина другой — 50 м. В этом случае размер самого большого участка составляет $25\text{ м} \times 25\text{ м}$, и надел после деления будет состоять из двух участков.

Теперь нужно вычислить рекурсивный случай. Здесь-то вам на помощь и приходит стратегия «разделяй и властвуй». В соответствии с ней при каждом рекурсивном вызове задача должна сокращаться. Как сократить эту задачу? Для начала разметим самые большие участки, которые можно использовать.



В исходном наделе можно разместить два участка 640×640 м, и еще останется место. Тут-то и наступает момент истины. Нераспределенный остаток — это тоже надел земли, который нужно разделить. *Так почему бы не применить к нему тот же алгоритм?*



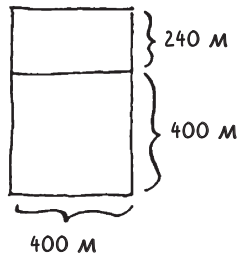
НОВЫЙ НАДЕЛ, КОТОРЫЙ ТОЖЕ
НУЖНО РАЗБИТЬ НА УЧАСТКИ

Итак, мы начали с надела 1680×640 м, который необходимо разделить на участки. Но теперь разделить нужно меньший сегмент — 640×400 м. Если

вы найдете самый большой участок, подходящий для такого размера, это будет самый большой участок, подходящий для всей фермы. Мы только что сократили задачу с размера 1680×640 м до 640×400 м!

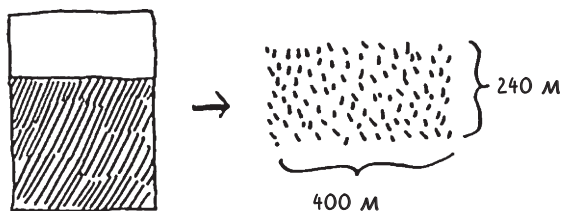
АЛГОРИТМ ЕВКЛИДА

«Если вы найдете самый большой участок, подходящий для такого размера, это будет самый большой участок, подходящий для всей фермы». Если истинность этого утверждения для вас неочевидна, не огорчайтесь. Она действительно неочевидна. К сожалению, доказательство получится слишком длинным, чтобы его можно было бы привести в книге, поэтому вам придется просто поверить мне на слово. Если вас интересует доказательство, поищите «алгоритм Евклида». Хорошее объяснение содержится на сайте Khan Academy: (<http://mng.bz/orm2>).

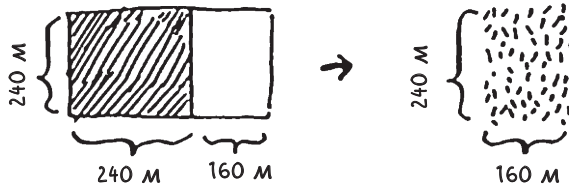


Применим тот же алгоритм снова. Если начать с участка 640×400 м, то размеры самого большого квадрата, который можно создать, составляют 400×400 м.

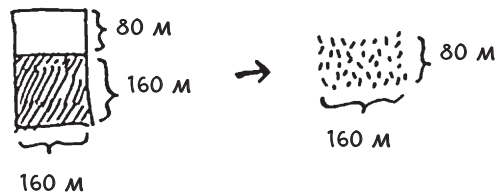
Остается меньший сегмент с размерами 400×240 м.



Отсекая поделенную часть, мы приходим к еще меньшему размеру сегмента, 240×160 м.

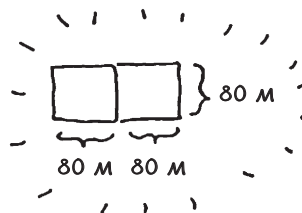


После очередного отсечения получается еще *меньший* сегмент.

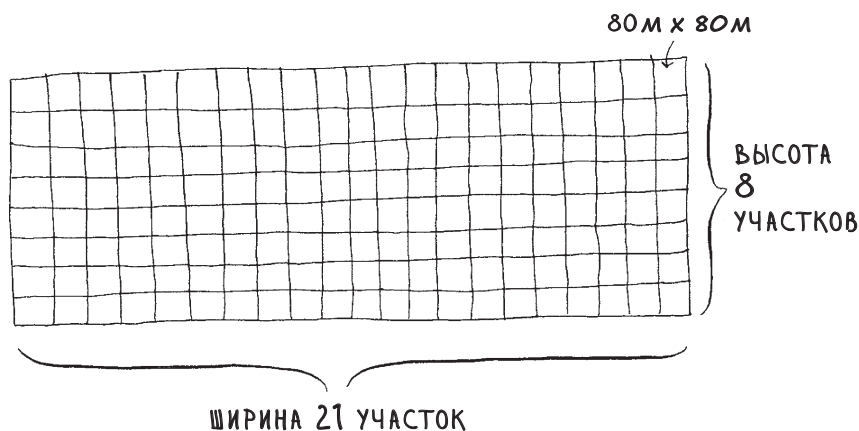


БАЗОВЫЙ СЛУЧАЙ!

Эге, да мы пришли к базовому случаю: 160 кратно 80. Если разбить этот сегмент на квадраты, ничего лишнего не останется!



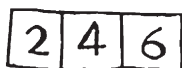
Итак, для исходного надела земли самый большой размер участка будет равен 80×80 м.



Вспомните, как работает стратегия «разделяй и властвуй».

1. Определите простейший случай как базовый.
2. Придумайте, как свести задачу к базовому случаю.

«Разделяй и властвуй» — не простой алгоритм, который можно применить для решения задачи. Скорее это подход к решению задачи. Рассмотрим еще один пример.



Имеется массив чисел.

Нужно просуммировать все числа и вернуть сумму. Сделать это в цикле совсем не сложно:

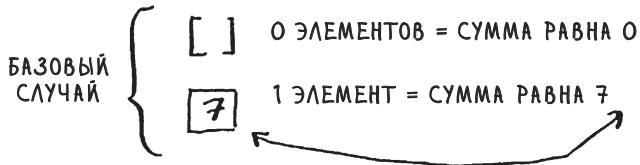
```
def sum(arr):
    total = 0
    for x in arr:
        total += x
    return total

print sum([1, 2, 3, 4])
```

Но как сделать то же самое с использованием рекурсивной функции?

Шаг 1: определить базовый случай. Как выглядит самый простой массив, который вы можете получить? Подумайте, как должен выглядеть про-

стейший случай, и продолжайте читать. Если у вас будет массив с 0 или 1 элементом, он суммируется достаточно просто.



Итак, с базовым случаем мы определились.

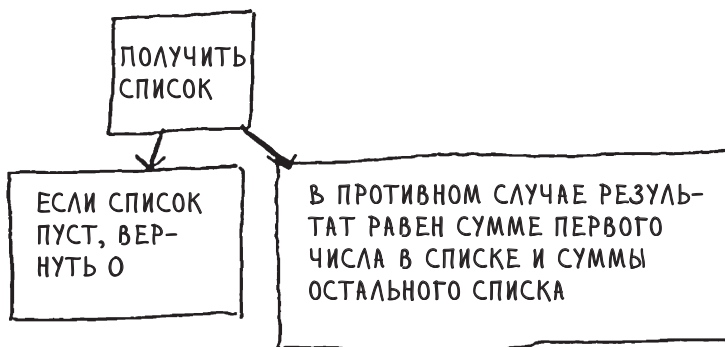
Шаг 2: каждый рекурсивный вызов должен приближать вас к пустому массиву. Как уменьшить размер задачи? Один из возможных способов:

$$\text{sum}([2, 4, 6]) = 12$$

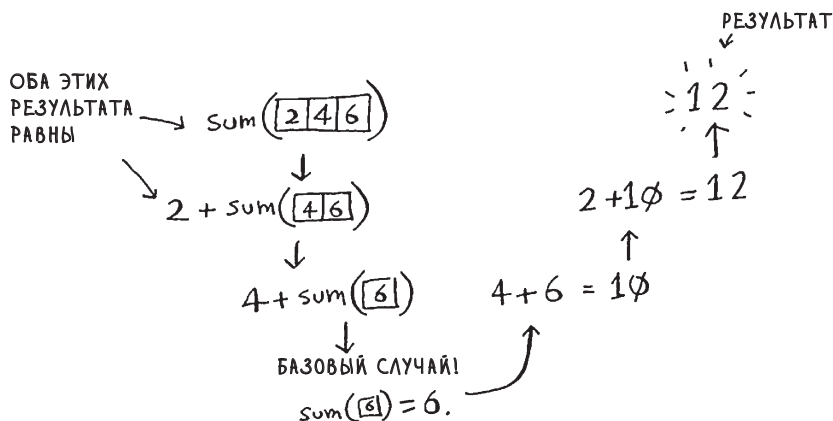
$$2 + \text{sum}([4, 6]) = 2 + 10 = 12$$

В любом случае результат равен 12. Но во второй версии функции `sum` передается меньший массив. А это означает, что вы сократили размер своей задачи!

Функция `sum` может работать по следующей схеме:



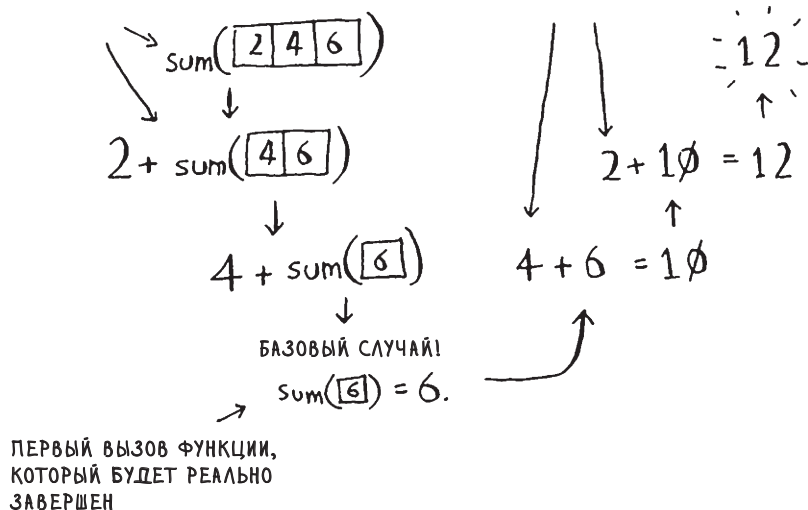
А вот как это выглядит в действии.



Вспомните, что при рекурсии сохраняется состояние.

НИ ОДИН ИЗ ЭТИХ
ВЫЗОВОВ ФУНКЦИИ НЕ
ЗАВЕРШИТСЯ ДО ТОГО,
КАК БУДЕТ ОБНАРУЖЕН
БАЗОВЫЙ СЛУЧАЙ!

ВСПОМНИТЕ, ЧТО РЕКУРСИЯ
СОХРАНЯЕТ СОСТОЯНИЕ ЭТИХ
ЧАСТИЧНО ЗАВЕРШЕННЫХ
ВЫЗОВОВ ФУНКЦИИ



СОВЕТ

Когда вы пишете рекурсивную функцию, в которой задействован массив, базовым случаем часто оказывается пустой массив или массив из одного элемента. Если вы не знаете, с чего начать, — начните с этого.

ПАРА СЛОВ О ФУНКЦИОНАЛЬНОМ ПРОГРАММИРОВАНИИ

Зачем применять рекурсию, если задача легко решается с циклом? Вполне резонный вопрос. Что ж, пора познакомиться с функциональным программированием!

В языках функционального программирования, таких как Haskell, циклов нет, поэтому для написания подобных функций приходится применять рекурсию. Если вы хорошо понимаете рекурсию, вам будет проще изучать функциональные языки. Например, вот как выглядит функция `sum` на языке Haskell:

<code>sum [] = 0</code>	◀.....	Базовый случай
<code>sum (x:xs) = x + (sum xs)</code>	◀.....	Рекурсивный случай

На первый взгляд кажется, что одна функция имеет два определения. Первое определение выполняется для базового случая, а второе — для рекурсивного. Функцию также можно записать на Haskell с использованием команды `if`:

```
sum arr = if arr == []
          then 0
          else (head arr) + (sum (tail arr))
```

Но первое определение проще читается. Так как рекурсия широко применяется в языке Haskell, в него включены всевозможные удобства для ее использования. Если вам нравится рекурсия или вы хотите изучить новый язык — присмотритесь к Haskell.

УПРАЖНЕНИЯ

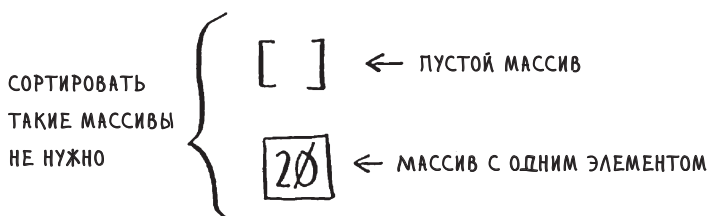
- 4.1 Напишите код для функции `sum` (см. выше).
- 4.2 Напишите рекурсивную функцию для подсчета элементов в списке.
- 4.3 Напишите рекурсивную функцию для нахождения наибольшего числа в списке.
- 4.4 Помните бинарный поиск из главы 1? Он тоже относится к классу алгоритмов «разделяй и властвуй». Сможете ли вы определить базовый и рекурсивный случаи для бинарного поиска?



Быстрая сортировка

Быстрая сортировка относится к алгоритмам сортировки. Она работает намного быстрее сортировки выбором и часто применяется в реальных программах. Быстрая сортировка также основана на стратегии «разделяй и властвуй».

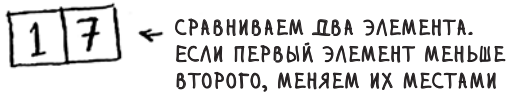
Воспользуемся быстрой сортировкой для упорядочения массива. Как выглядит самый простой массив, с которым может справиться алгоритм сортировки (помните подсказку из предыдущего раздела)? Некоторые массивы вообще не нуждаются в сортировке.



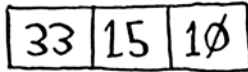
Пустые массивы и массивы, содержащие всего один элемент, станут базовым случаем. Такие массивы можно просто возвращать в исходном виде — сортировать ничего не нужно:

```
def quicksort(array):
    if len(array) < 2:
        return array
```

Теперь перейдем к массивам большего размера. Массив из двух элементов тоже сортируется без особых проблем.



А как насчет массива из трех элементов?

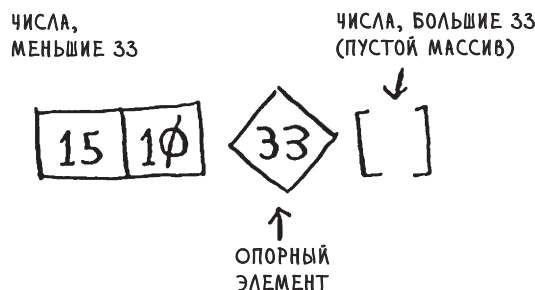


Помните: мы используем стратегию «разделяй и властвуй». Следовательно, массив должен разделяться до тех пор, пока мы не придем к базовому случаю. Алгоритм быстрой сортировки работает так: сначала в массиве выбирается элемент, который называется *опорным*.



О том, как выбрать хороший опорный элемент, будет рассказано далее. А пока предположим, что опорным становится первый элемент массива.

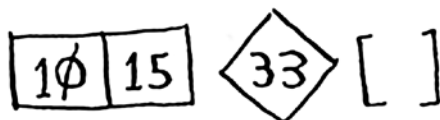
Теперь мы находим элементы, меньшие опорного, и элементы, бóльшие опорного.



Этот процесс называется *разделением*. Теперь у вас имеются:

- подмассив всех элементов, меньших опорного;
- опорный элемент;
- подмассив всех элементов, больших опорного.

Два подмассива не отсортированы — они просто выделены из исходного массива. Но если бы они *были* отсортированы, то провести сортировку всего массива было бы несложно.



Если бы подмассивы были отсортированы, то их можно было бы объединить в порядке «левый подмассив — опорный элемент — правый подмассив» и получить отсортированный массив. В нашем примере получается $[10, 15] + [33] + [] = [10, 15, 33]$, то есть отсортированный массив.

Как отсортировать подмассивы? Базовый случай быстрой сортировки уже знает, как сортировать массивы из двух элементов (левый подмассив) и пустые массивы (правый подмассив). Следовательно, если применить алгоритм быстрой сортировки к двум подмассивам, а затем объединить результаты, получится отсортированный массив!

```
quicksort([15, 10]) + [33] + quicksort([])
> [10, 15, 33]      <----- Отсортированный массив
```

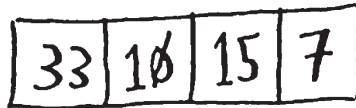
Этот метод работает при любом опорном элементе. Допустим, вместо 33 в качестве опорного был выбран элемент 15.



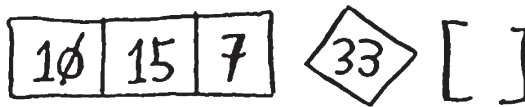
Оба подмассива состоят из одного элемента, а вы уже умеете сортировать такие подмассивы. Получается, что вы умеете сортировать массивы из трех элементов. Это делается так:

1. Выбрать опорный элемент.
2. Разделить массив на два подмассива: элементы, меньшие опорного, и элементы, большие опорного.
3. Рекурсивно применить быструю сортировку к двум подмассивам.

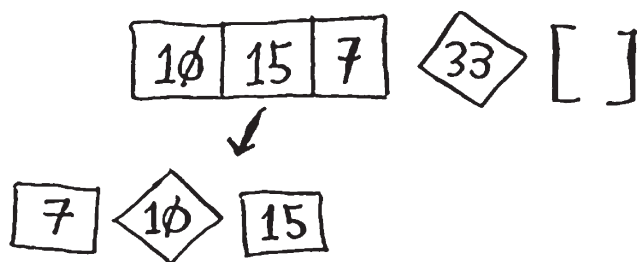
Как насчет массива из четырех элементов?



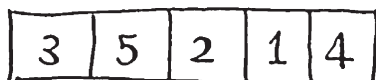
Предположим, опорным снова выбирается элемент 33.



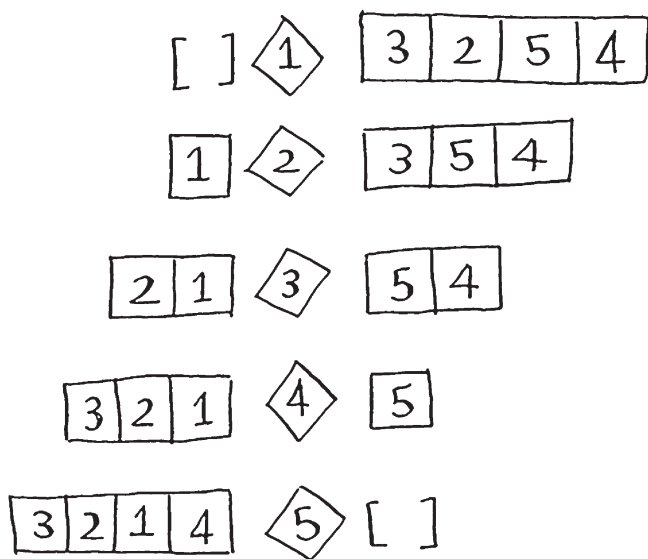
Левый подмассив состоит из трех элементов. Вы уже знаете, как сортируется массив из трех элементов: нужно рекурсивно применить к нему быструю сортировку.



Следовательно, вы можете отсортировать массив из 4 элементов. А если вы можете отсортировать массив из 4 элементов, то вы также можете отсортировать массив из 5 элементов. Почему? Допустим, имеется массив из 5 элементов.

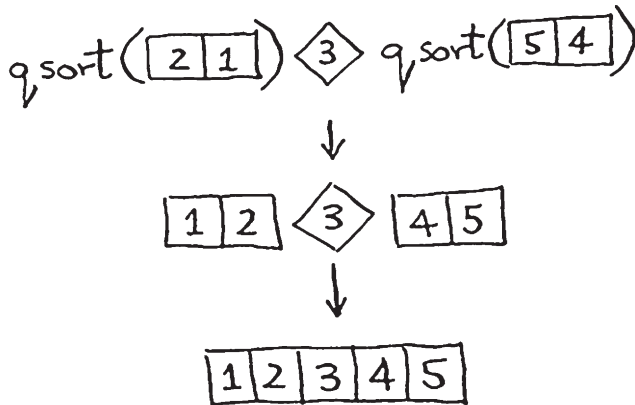


Вот как выглядят все варианты разделения этого массива в зависимости от выбранного опорного элемента:

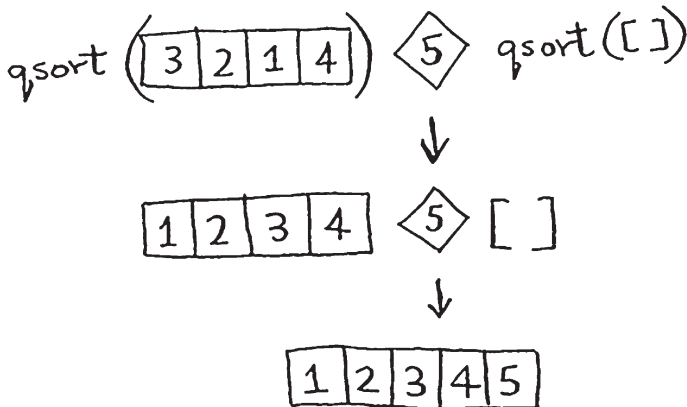


Все эти подмассивы содержат от 0 до 4 элементов. А вы уже знаете, как отсортировать массив, содержащий от 0 до 4 элементов, с использованием быстрой сортировки! Таким образом, вне зависимости от выбора опорного элемента вы можете рекурсивно вызывать быструю сортировку для двух подмассивов.

Например, предположим, что в качестве опорного выбирается элемент 3. Вы применяете быструю сортировку к подмассивам.



Подмассивы отсортированы, и теперь из них можно собрать отсортированный массив. Решение работает даже в том случае, если выбрать в качестве опорного элемент 5:



Итак, решение работает независимо от выбора опорного элемента. Следовательно, вы можете отсортировать массив из 5 элементов. По той же логике вы можете отсортировать массив из 6 элементов и т. д.

ДОКАЗАТЕЛЬСТВО ПО ИНДУКЦИИ

Вы только что познакомились с методом доказательства по индукции! Это один из способов, доказывающих, что ваш алгоритм работает. Каждое индуктивное доказательство состоит из двух частей: базы (базового случая) и индукционного перехода. Звучит знакомо? Допустим, я хочу доказать, что могу подняться на самый верх стремянки. Если мои ноги стоят на ступеньке, — это индукционный переход. Таким образом, если я стою на ступеньке 2, то могу подняться на ступеньку 3. Что касается базового случая, я сейчас стою на ступеньке 1. Из этого следует, что я могу подняться на самый верх стремянки, каждый раз поднимаясь на одну ступеньку.

Аналогичные рассуждения применимы к быстрой сортировке. Работоспособность алгоритма для базового случая — массивов с размером 0 и 1 — была продемонстрирована. В индукционном переходе я показал, что если быстрая сортировка работает для массива из 1 элемента, то она будет работать для массива из 2 элементов. А если она работает для массивов из 2 элементов, то она будет работать для массивов из 3 элементов и т. д. Из этого можно сделать вывод, что быстрая сортировка будет работать для всех массивов любого размера. Я не буду подробно рассматривать доказательства по индукции, но это интересный метод, который идет рука об руку со стратегией «разделяй и властвуй».



А вот как выглядит программный код быстрой сортировки:

```
def quicksort(array):
    if len(array) < 2:
        return array
    else:
        pivot = array[0]
        less = [i for i in array[1:] if i < pivot]
        greater = [i for i in array[1:] if i > pivot]
        return quicksort(less) + [pivot] + quicksort(greater)

print quicksort([10, 5, 2, 3])
```

Базовый случай: массивы с 0 и 1 элементом уже "отсортированы"

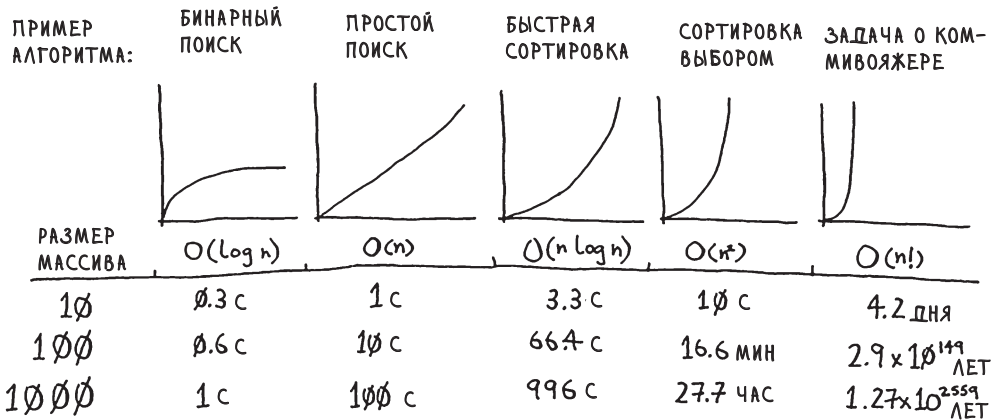
Рекурсивный случай

Подмассив всех элементов, меньших опорного

Подмассив всех элементов, больших опорного

Снова об «О-большом»

Алгоритм быстрой сортировки уникален тем, что его скорость зависит от выбора опорного элемента. Прежде чем рассматривать быструю сортировку, вспомним наиболее типичные варианты времени выполнения для «О-большое».



Оценки для медленного компьютера, выполняющего 10 операций в секунду

На графиках приведены примерные оценки времени при выполнении 10 операций в секунду. Они не претендуют на точность, а всего лишь дают представление о том, насколько различается время выполнения. Конечно,

на практике ваш компьютер способен выполнять гораздо больше 10 операций в секунду.

Для каждого времени выполнения также приведен пример алгоритма. Возьмем алгоритм сортировки выбором, о котором вы узнали в главе 2. Он обладает временем $O(n^2)$, и это довольно медленный алгоритм.

Другой алгоритм сортировки — так называемая *сортировка слиянием* — работает за время $O(n \log n)$. Намного быстрее! С быстрой сортировкой дело обстоит сложнее. В худшем случае быстрая сортировка работает за время $O(n^2)$.

Ничуть не лучше сортировки выбором! Но это худший случай, а в среднем быстрая сортировка выполняется за время $O(n \log n)$. Вероятно, вы спросите:

- что в данном случае понимается под «худшим» и «средним» случаем?
- если быстрая сортировка в среднем выполняется за время $O(n \log n)$, а сортировка слиянием выполняется за время $O(n \log n)$ всегда, то почему бы не использовать сортировку слиянием? Разве она не быстрее?

Сортировка слиянием и быстрая сортировка

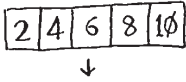
Допустим, у вас имеется простая функция для вывода каждого элемента в списке:

```
def print_items(list):
    for item in list:
        print item
```

Эта функция последовательно перебирает все элементы списка и выводит их. Так как функция перебирает весь список, она выполняется за время $O(n)$. Теперь предположим, что вы изменили эту функцию и она делает секундную паузу перед выводом:

```
from time import sleep
def print_items2(list):
    for item in list:
        sleep(1)
        print item
```

Перед выводом элемента функция делает паузу продолжительностью в 1 секунду. Предположим, вы выводите список из 5 элементов с использованием обеих функций:



`print.items: 2 4 6 8 10`

`print_items2: <ПАУЗА> 2 <ПАУЗА> 4 <ПАУЗА> 6 <ПАУЗА> 8 <ПАУЗА> 10`

Обе функции проходят по списку один раз, и обе выполняются за время $O(n)$. Как вы думаете, какая из них работает быстрее? Я думаю, `print_items` работает намного быстрее, потому что она не делает паузу перед выводом каждого элемента. Следовательно, даже притом что обе функции имеют одинаковую скорость «О-большое», реально `print_items` работает быстрее. Когда вы используете «О-большое» (например, $O(n)$), в действительности это означает следующее:

$$c * n$$

ФИКСИРОВАННЫЙ
ПРОМЕЖУТОК
ВРЕМЕНИ

Здесь c — некоторый фиксированный промежуток времени для вашего алгоритма. Он называется *константой*. Например, время выполнения может составлять *10 миллисекунд * n* для `print_items` против *1 секунды * n* для `print_items2`.

Обычно константа игнорируется, потому что если два алгоритма имеют разное время «О-большое», она роли не играет. Для примера возьмем бинарный и простой поиск. Допустим, такие константы присутствуют в обоих алгоритмах.

$\frac{10 \text{ мс} * n}{\text{ПРОСТОЙ ПОИСК}}$	$\frac{1 \text{ с} * \log n}{\text{БИНАРНЫЙ ПОИСК}}$
--	--

Первая реакция: «Ого! У простого поиска константа равна 10 миллисекундам, а у бинарного поиска — 1 секунда. Простой поиск намного быстрее!» Теперь предположим, что поиск ведется по списку из 4 миллиардов элементов. Время будет таким:

$$\begin{array}{l|l} \text{ПРОСТОЙ ПОИСК} & 10_{\text{мс}} \times 4 \text{ МИЛЛИАРДА} = 463 \text{ ДНЯ} \\ \text{БИНАРНЫЙ ПОИСК} & 1 \text{ с} \times 32 = 32 \text{ СЕКУНДЫ} \end{array}$$

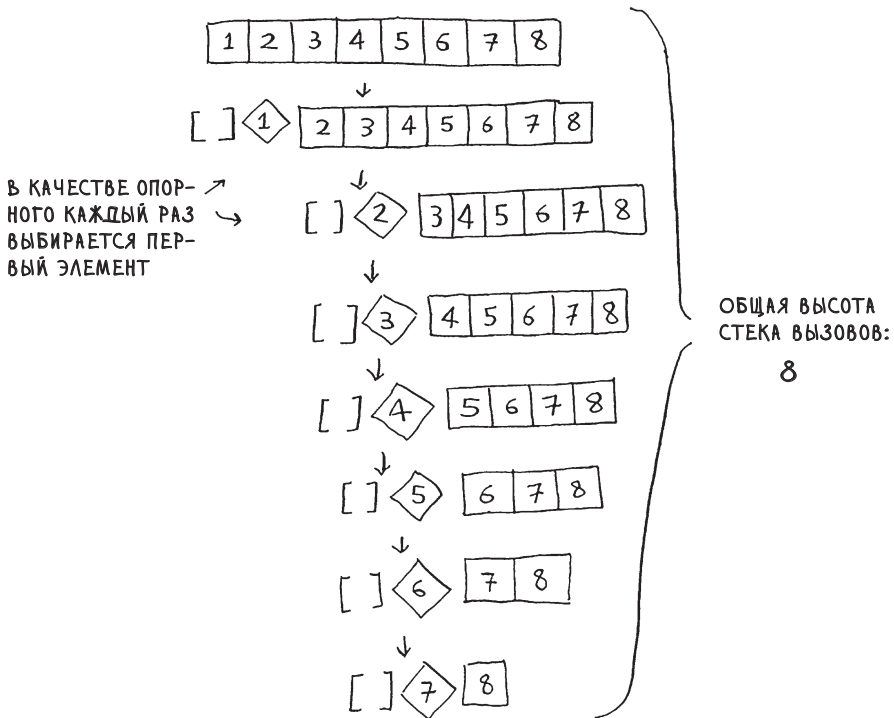
При расчете времени бинарного поиска используется значение 32, поскольку $\log(4\,000\,000\,000)$ равен 32. Как видите, бинарный поиск все равно работает намного быстрее. Константа ни на что не повлияла.

Однако в некоторых случаях константа *может* иметь значение. Один из примеров такого рода — быстрая сортировка и сортировка слиянием. У быстрой сортировки константа меньше, чем у сортировки слиянием, поэтому, несмотря на то что оба алгоритма характеризуются временем $O(n \log n)$, быстрая сортировка работает быстрее. А на практике быстрая сортировка работает быстрее, потому что средний случай встречается намного чаще худшего.

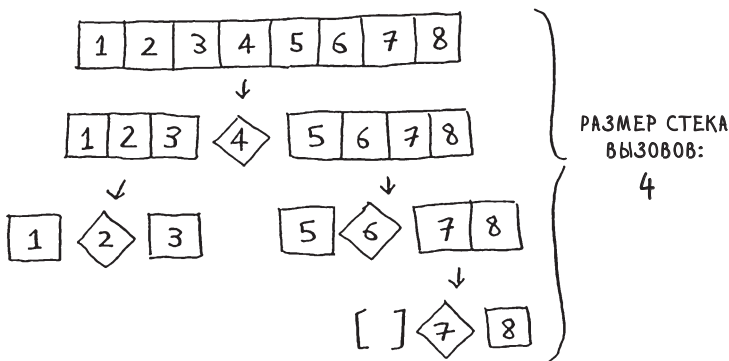
А теперь ответим на первый вопрос: как выглядит средний случай по сравнению с худшим?

Средний и худший случаи

Быстродействие быстрой сортировки сильно зависит от выбора опорного элемента. Предположим, опорным всегда выбирается первый элемент, а быстрая сортировка применяется к *уже отсортированному* массиву. Быстрая сортировка не проверяет, отсортирован входной массив или нет, и все равно пытается его отсортировать.



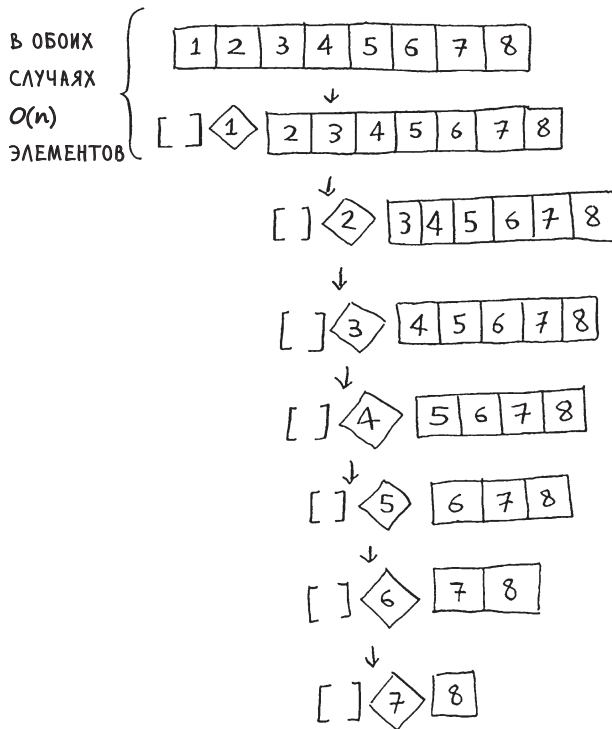
Обратите внимание: на этот раз массив не разделяется на две половины. Вместо этого один из двух подмассивов всегда пуст, так что стек вызовов получается очень длинным. Теперь предположим, что в качестве опорного всегда выбирается средний элемент. Посмотрим, как выглядит стек вызовов в этом случае.



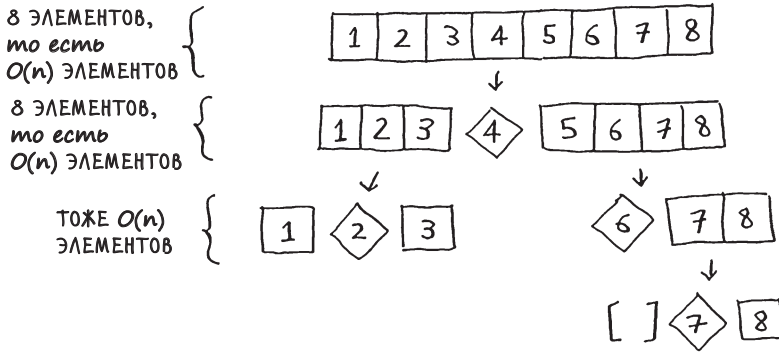
Стек намного короче! Массив каждый раз делится надвое, поэтому такое количество рекурсивных вызовов излишне. Вы быстрее добираетесь до базового случая, и стек вызовов получается более коротким.

Первый из рассмотренных примеров описывает худший сценарий, а второй — лучший. В худшем случае размер стека описывается как $O(n)$. В лучшем случае он составит $O(\log n)$.

Теперь рассмотрим первый уровень стека. Один элемент выбирается опорным, а остальные элементы делятся на подмассивы. Вы перебираете все 8 элементов массива, поэтому первая операция выполняется за время $O(n)$. На этом уровне стека вызовов вы обратились ко всем 8 элементам. Но на самом деле вы обращаетесь к $O(n)$ элементам на каждом уровне стека вызовов!



Даже если массив будет разделен другим способом, вы все равно каждый раз обращаетесь к $O(n)$ элементам.



В этом примере существуют $O(\log n)$ (с технической точки зрения правильнее сказать «высота стека вызовов равна $O(\log n)$ ») уровней. А так как каждый уровень занимает время $O(n)$, то весь алгоритм займет время $O(n) * O(\log n) = O(n \log n)$. Это сценарий лучшего случая.

В худшем случае существуют $O(n)$ уровней, поэтому алгоритм займет время $O(n) * O(n) = O(n^2)$.

А теперь сюрприз: лучший случай также является средним. Если вы всегда будете выбирать опорным элементом случайный элемент в массиве, быстрая сортировка в среднем завершится за время $O(n \log n)$. (За одним исключением: если все элементы массива одинаковы, время выполнения всегда будет худшим без дополнительной логики.) Это один из самых быстрых существующих алгоритмов сортировки, который заодно является хорошим примером стратегии «разделяй и властвуй».

УПРАЖНЕНИЯ

Запишите «О-большое» для каждой из следующих операций.

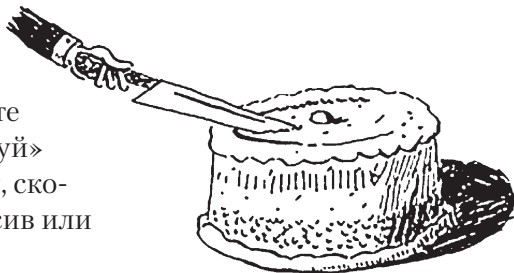
- 4.5 Вывод значения каждого элемента массива.
- 4.6 Удвоение значения каждого элемента массива.
- 4.7 Удвоение значения только первого элемента массива.

- 4.8 Создание таблицы умножения для всех элементов массива. Например, если массив состоит из элементов [2, 3, 7, 8, 10], сначала каждый элемент умножается на 2, затем каждый элемент умножается на 3, затем на 7 и т. д.

	2	3	7	8	10
2	4	6	14	16	20
3	6	9	21	24	30
7	14	21	49	56	70
8	16	24	56	64	80
10	20	30	70	80	100

Шпаргалка

- Стратегия «разделяй и властвуй» основана на разбиении задачи на уменьшающиеся фрагменты. Если вы используете стратегию «разделяй и властвуй» со списком, то базовым случаем, скорее всего, является пустой массив или массив из одного элемента.
- Если вы реализуете алгоритм быстрой сортировки, выберите в качестве опорного случайный элемент. Среднее время выполнения быстрой сортировки составляет $O(n \log n)$!
- Из двух алгоритмов с одинаковой скоростью «О-большое» один может быть быстрее другого. Именно по этой причине быстрая сортировка быстрее сортировки слиянием.
- При сравнении простой сортировки с бинарной константа почти никогда роли не играет, потому что $O(\log n)$ слишком сильно превосходит $O(n)$ по скорости при большом размере списка.



5

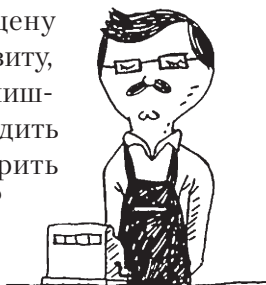
Хеш-таблицы

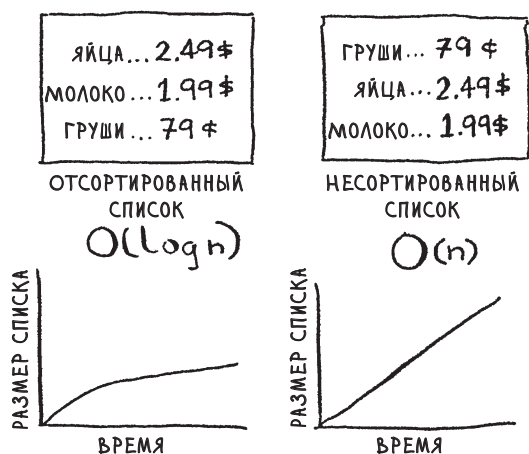


В этой главе

- ✓ Вы узнаете о хеш-таблицах — одной из самых полезных базовых структур данных. Хеш-таблицы находят множество применений; в этой главе рассматриваются распространенные варианты использования.
- ✓ Вы изучите внутреннее устройство хеш-таблиц: реализацию, коллизии и хеш-функции. Это поможет вам понять, как анализируется производительность хеш-таблицы.

Представьте, что вы продавец в маленьком магазинчике. Когда клиент покупает товары, вы проверяете их цену по книге. Если записи в книге не упорядочены по алфавиту, то поиск слова «апельсины» в каждой строке займет слишком много времени. Фактически вам придется проводить простой поиск из главы 1, а для этого нужно проверить каждую запись. Помните, сколько времени это займет? $O(n)$. Если же книга упорядочена по алфавиту, вы сможете воспользоваться бинарным поиском, время которого составляет всего $O(\log n)$.





На всякий случай напомним, что время $O(n)$ и $O(\log n)$ — далеко не одно и то же! Предположим, вы можете просмотреть 10 записей в книге за секунду. В следующей таблице показано, сколько времени займет простой и бинарный поиск.

КОЛИЧЕСТВО ЗАПИСЕЙ В КНИГЕ	$O(n)$	$O(\log n)$
100	10с	1с ← НЕОБХОДИМО ПРОВЕРИТЬ $\log_2 100 = 7$ СТРОК
1000	1.66 мин	1с ← $\log_2 1000 = 10$ СТРОК
10000	16.6 мин	2с ← $\log_2 10000 = 14$ СТРОК = 2 с

Вы уже знаете, что бинарный поиск работает очень быстро. Но поиск данных в книге — головная боль для кассира, даже если ее содержимое отсортировано. Пока вы листаете страницы, клиент потихоньку начинает выходить из себя. Гораздо удобнее было бы завести помощницу, которая помнит все названия товаров и цены. Тогда ничего искать вообще не придется: вы спрашиваете помощницу, а она мгновенно отвечает.



Ваша помощница Мэгги может за время $O(1)$ сообщить цену любого товара независимо от размера книги. Она работает еще быстрее, чем бинарный поиск.

КОЛИЧЕСТВО ЭЛЕМЕНТОВ В КНИГЕ	ПРОСТОЙ ПОИСК	БИНАРНЫЙ ПОИСК	МЭГГИ
	$O(n)$	$O(\log n)$	$O(1)$
1000	10 с	1 с	МГНОВЕННО
10000	1.6 мин	1 с	МГНОВЕННО
100000	16.6 мин	2 с	МГНОВЕННО

Просто чудо, а не девушка! И где взять такую Мэгги?

Обратимся к структурам данных. Пока вам известны две структуры данных: массивы и списки. (О стеках я не говорю, потому что нормальный поиск в стеке невозможен.) Книгу можно реализовать в виде массива.

(яйца, 2.49)	(молоко, 1.49)	(груши, 0.79)
--------------	----------------	---------------

Каждый элемент массива на самом деле состоит из двух элементов: названия товара и его цены. Если отсортировать массив по имени, вы сможете провести по нему бинарный поиск для определения цены товара. Это означает, что поиск будет выполняться за время $O(\log n)$. Но нам нужно, чтобы поиск выполнялся за время $O(1)$ (другими словами, вы хотите создать «Мэгги»). В этом вам помогут хеш-функции.

Хеш-функции

Хеш-функция представляет собой функцию, которая получает строку¹ и возвращает число:



В научной терминологии говорят, что хеш-функция «отображает строки на числа». Можно подумать, что найти закономерности получения чисел для подаваемых на вход строк невозможно. Однако хеш-функция должна соответствовать некоторым требованиям.

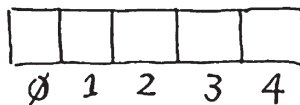
- Она должна быть последовательной. Допустим, вы передали ей строку «апельсины» и получили 4. Это значит, что каждый раз в будущем, передавая ей строку «апельсины», вы будете получать 4. Без этого хеш-таблица бесполезна.
- Разным словам должны соответствовать разные числа. Например, хеш-функция, которая возвращает 1 для каждого полученного слова, никуда

¹ Под «строкой» в данном случае следует понимать любые данные — последовательность байтов.

не годится. В идеале каждое входное слово должно отображаться на свое число.

Итак, хеш-функция связывает строки с числами. Зачем это нужно, спросите вы? Так ведь это позволит нам реализовать «Мэгги»!

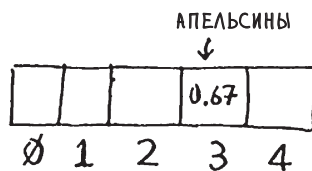
Начнем с пустого массива:



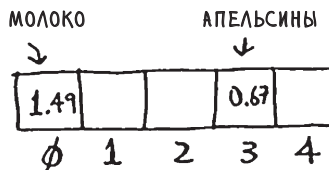
Все цены будут храниться в этом массиве; передадим хеш-функции строку «апельсины».



Хеш-функция выдает значение «3». Сохраним цену апельсинов в элементе массива с индексом 3.



Добавим молоко. Передадим хеш-функции строку «молоко».



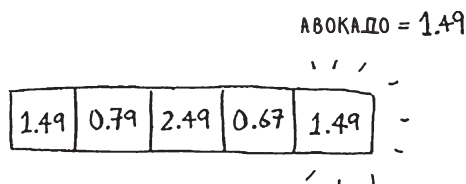
Продолжайте действовать так, и со временем весь массив будет заполнен ценами на товары.

1.49	0.79	2.49	0.67	1.49
------	------	------	------	------

А теперь вы спрашиваете: сколько стоит авокадо? Искать в массиве ничего не нужно, просто передайте строку «авокадо» хеш-функции.



Результат показывает, что значение хранится в элементе с индексом 4. И оно, конечно, там и находится!



Хеш-функция сообщает, где хранится цена, и вам вообще не нужно ничего искать! Такое решение работает, потому что:

- Хеш-функция неизменно связывает название с одним индексом. Каждый раз, когда она вызывается для строки «авокадо», вы получаете обратно одно и то же число. При первом вызове этой функции вы узнаете, где следует сохранить цену авокадо, а при последующих вызовах она сообщает, где взять эту цену.
- Хеш-функция связывает разные строки с разными индексами. «Авокадо» связывается с индексом 4, а «молоко» — с индексом 0. Для каждой строки находится отдельная позиция массива, в которой сохраняется цена этого товара.

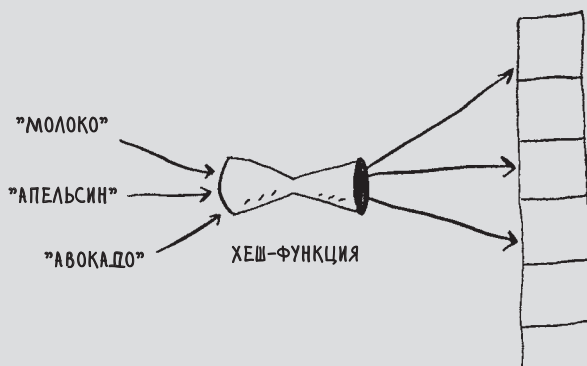
- Хеш-функция знает размер массива и возвращает только действительные индексы. Таким образом, если длина массива равна 5 элементам, хеш-функция не вернет 100, потому что это значение не является действительным индексом в массиве.

Поздравляю: вы создали «Мэгги»! Свяжите воедино хеш-функцию и массив, и вы получите структуру данных, которая называется *хеш-таблицей*. Хеш-таблица станет первой изученной вами структурой данных, с которой связана дополнительная логика. Массивы и списки напрямую отображаются на адреса памяти, но хеш-таблицы устроены более умно. Они определяют место хранения элементов при помощи хеш-функций.

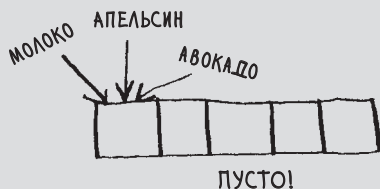
Вероятно, хеш-таблицы станут самой полезной из сложных структур данных, с которыми вы познакомитесь. Они также известны под другими названиями: «ассоциативные массивы», «словари», «отображения», «хеш-карты» или просто «хеши». Хеш-таблицы исключительно быстро работают! Помните описание массивов и связанных списков из главы 2? Обращение к элементу массива происходит мгновенно. А хеш-таблицы используют массивы для хранения данных, поэтому при обращении к элементам они не уступают массивам.

В ЧЕМ ПОДВОХ?

Мы только что рассмотрели идеальную хеш-функцию. Она связывает каждый товар точно с отдельной позицией в массиве:



Посмотрите: все товары находятся в своих позициях. В реальности такого идеального однозначного сопоставления, скорее всего, не получится. Товарам придется соседствовать. В одной позиции будут находиться несколько товаров, а другие позиции останутся пустыми.



Мы обсудим эту ситуацию в разделе о коллизиях. А пока просто знайте, что, хотя хеш-таблицы очень полезны, они редко точно связывают элементы с позициями.

Кстати, такое взаимно-однозначное связывание называется инъективной функцией.

Запомните это словечко: пригодится, чтобы произвести впечатление на друзей!

Скорее всего, вам никогда не придется заниматься реализацией хеш-таблиц самостоятельно. В любом приличном языке существует реализация хеш-таблиц. В Python тоже есть хеш-таблицы; они называются *словарями*. Новая хеш-таблица задается пустыми фигурными скобками:

```
>>> book = {}
```



`book` — новая хеш-таблица. Добавим в `book` несколько цен:

```
>>> book["апельсин"] = 0.67
>>> book["молоко"] = 1.49
>>> book["авокадо"] = 1.49
>>> print book
{'avocado': 1.49, 'apple': 0.67, 'milk': 1.49}
```

←..... Апельсины стоят 67 центов
 ←..... Молоко стоит 1 доллар 49 центов

Пока все просто! А теперь запросим цену авокадо:

```
>>> print(book["авокадо"])
1.49  <----- Цена авокадо
```

Хеш-таблица состоит из ключей и значений. В хеше `book` имена продуктов являются ключами, а цены — значениями. Хеш-таблица связывает ключи со значениями.

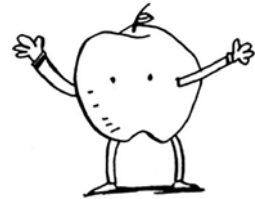
В следующем разделе приведены примеры, в которых хеш-таблицы приносят большую пользу.

АПЕЛЬ-СИНЫ	0.67
МОЛОКО	1.49
АВО-КАДО	1.49

ХЕШ-ТАБЛИЦА С ЦЕНАМИ
НА ПРОДУКТЫ

УПРАЖНЕНИЯ

Очень важно, чтобы хеш-функции были последовательными, то есть неизменно возвращали один и тот же результат для одинаковых входных данных. Если это условие будет нарушено, вы не сможете найти свой элемент после того, как он будет помещен в хеш-таблицу!



Какие из следующих функций являются последовательными?

- 5.1 `f(x) = 1` <----- Возвращает "1" для любых входных значений
- 5.2 `f(x) = rand()` <----- Возвращает случайное число
- 5.3 `f(x) = next_empty_slot()` <----- Возвращает индекс следующего пустого элемента в хеш-таблице
- 5.4 `f(x) = len(x)` <----- Возвращает длину полученной строки

Примеры использования

Хеш-таблицы повсеместно применяются на практике. В этом разделе представлены некоторые примеры.

Использование хеш-таблиц для поиска

В вашем телефоне есть удобная встроенная телефонная книга.

С каждым именем связывается номер телефона.



BADE MAMA → 581 660 9820
 ALEX MANNING → 484 234 4680
 JANE MARIN → 415 567 3579

Предположим, вы хотите построить такую телефонную книгу. Имена людей в этой книге связываются с номерами. Телефонная книга должна поддерживать следующие функции:

- добавление имени человека и номера телефона, связанного с этим именем;
- получение номера телефона, связанного с введенным именем.

Такая задача идеально подходит для хеш-таблиц! Хеш-таблицы отлично работают, когда вы хотите:

- создать связь, отображающую один объект на другой;
- найти значение в списке.

Построить телефонную книгу, в общем-то, несложно. Начните с создания новой хеш-таблицы:

```
>>> phone_book = {}
```

Добавим в телефонную книгу несколько номеров:

```
>>> phone_book["Дженни"] = 8675309
>>> phone_book["служба спасения"] = 911
```

Вот и все! Теперь предположим, что вы хотите найти номер телефона Дженни (Jenny). Просто передайте ключ хешу:

```
>>> print(phone_book["jenny"])
8675309  ←..... Номер Дженни
```

А теперь представьте, что то же самое вам пришлось бы делать с массивом.

Как бы вы это сделали? Хеш-таблицы упрощают моделирование отношений между объектами.



ДЖЕННИ	8675309
СЛУЖБА СПАСЕНИЯ	911

ХЕШ-ТАБЛИЦА
КАК ТЕЛЕФОННАЯ КНИГА

Хеш-таблицы используются для поиска соответствий в гораздо большем масштабе. Например, представьте, что вы хотите перейти на веб-сайт — допустим, <http://adit.io>. Ваш компьютер должен преобразовать символическое имя *adit.io* в IP-адрес.

ADIT.IO → 173.255.248.55

Для любого посещаемого веб-сайта его имя преобразуется в IP-адрес:

GOOGLE.COM → 74.125.239.133
FACEBOOK.COM → 173.252.128.6
SCRIBD.COM → 23.235.47.175

Связать символическое имя с IP-адресом? Идеальная задача для хеш-таблиц! Этот процесс называется *преобразованием DNS*. Хеш-таблицы — всего лишь один из способов реализации этой функциональности. На компьютерах существует кэш DNS, в котором хранятся подобные сопоставления для недавно посещенных сайтов и который удобно создавать с помощью хеш-таблицы.

Исключение дубликатов

Предположим, вы руководите избирательным участком. Естественно, каждый избиратель может проголосовать всего один раз. Как проверить, что он не голосовал ранее? Когда человек приходит голосовать, вы узнаете его полное имя, а затем проверяете по списку уже проголосовавших избирателей.



Если имя входит в список, значит, этот человек уже проголосовал — гоните наглеца! В противном случае вы добавляете имя в список и разрешаете ему проголосовать. Теперь предположим, что желающих проголосовать много и список уже проголосовавших достаточно велик.



Каждый раз, когда кто-то приходит голосовать, вы вынуждены просматривать этот гигантский список и проверять, голосовал он или нет. Однако существует более эффективное решение: воспользоваться хешем!

Сначала создадим хеш для хранения информации об уже проголосовавших людях:

```
>>> voted = {}
```

Когда кто-то приходит голосовать, проверьте, присутствует ли его имя в хеше:

```
>>> value = voted.get("Том")
```

Функция `value` принимает значение `True`, если ключ "Том" присутствует в хеш-таблице. В противном случае она принимает значение `False`. С помощью этой функции можно проверить, голосовал избиратель ранее или нет!



Код выглядит так:

```
voted = {}

def check_voter(name):
    if voted.get(name):
        print("Выгнать его!")
    else:
        voted[name] = True
        print("Допустить к голосованию!")
```

Давайте протестируем его на нескольких примерах:

```
>>> check_voter("Том")
Допустить к голосованию!
>>> check_voter("Майк")
Допустить к голосованию!
>>> check_voter("Майк")
Выгнать его!
```

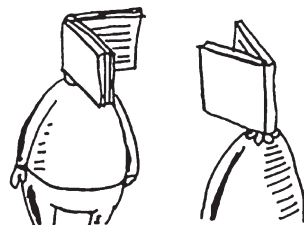
Когда Том приходит на участок в первый раз, программа разрешает ему проголосовать. Потом приходит Майк, который тоже допускается к голосованию. Но потом Майк делает вторую попытку, и на этот раз у него ничего не получается.

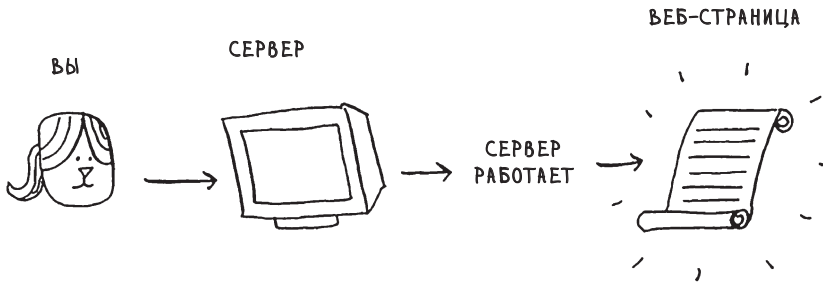
Если бы имена проголосовавших хранились в списке, то выполнение функции со временем замедлилось бы, потому что функции пришлось бы проводить простой поиск по всему списку. Но имена хранятся в хеш-таблице, а хеш-таблица мгновенно сообщает, присутствует имя избирателя в списке или нет. Проверка дубликатов в хеш-таблице выполняется очень быстро.

Использование хеш-таблицы как кэша

Последний пример: кэширование. Если вы работаете над созданием веб-сайтов, то, вероятно, уже слышали о пользе кэширования. Общая идея кэширования такова: допустим, вы заходите на сайт <https://facebook.com>.

1. Вы обращаетесь с запросом к серверу Facebook.
2. Сервер ненадолго задумывается, генерирует веб-страницу и отправляет ее вам.
3. Вы получаете веб-страницу.



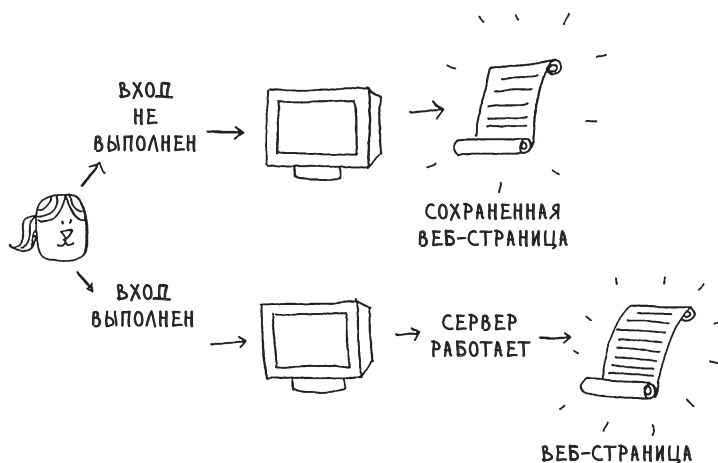


Например, на Facebook сервер может собирать информацию о действиях всех ваших друзей, чтобы представить ее вам. На то, чтобы собрать всю информацию и передать ее вам, требуется пара секунд. С точки зрения пользователя, пара секунд — это очень долго. Он начинает думать: «Почему Facebook работает так медленно?» Кроме того серверам Facebook приходится обслуживать миллионы людей, и эти пары секунд для них суммируются. Серверы Facebook трудятся в полную силу, чтобы сгенерировать все эти страницы. Нельзя ли сделать так, чтобы серверы Facebook выполняли меньше работы?

Представьте, что у вас есть племянница, которая пристает к вам с вопросами о планетах: «Сколько километров от Земли до Марса?», «А сколько километров до Луны?», «А до Юпитера?» Каждый раз вы вводите запрос в Google и сообщаете ей ответ. На это уходит пара минут. А теперь представьте, что она всегда спрашивает: «Сколько километров от Земли до Луны?» Довольно быстро вы запоминаете, что Луна находится на расстоянии 384 400 километров от Земли. Искать информацию в Google не нужно... вы просто запоминаете и выдаете ответ. Вот так работает механизм кэширования: сайт просто запоминает данные, вместо того чтобы пересчитывать их заново.

Если вы вошли на Facebook, то весь контент, который вы видите, адаптирован специально для вас. Каждый раз, когда вы заходите на *facebook.com*, серверам приходится думать, какой контент вас интересует. Если же вы не ввели учетные данные на Facebook, то вы видите страницу входа. Все пользователи видят одну и ту же страницу входа. Facebook постоянно получает одинаковые запросы: «Я еще не вошел на сайт, выдайте мне домашнюю

страницу». Сервер перестает выполнять лишнюю работу и генерировать домашнюю страницу снова и снова. Вместо этого он запоминает, как выглядит домашняя страница, и отправляет ее вам.



Такой механизм хранения называется *кэшированием*. Он обладает двумя преимуществами:

- вы получаете веб-страницу намного быстрее, как и в том случае, когда запомнили расстояние от Земли до Луны. Когда племянница в следующий раз задаст вопрос, вам не придется гуглить. Вы можете выдать ответ мгновенно;
- Facebook приходится выполнять меньше работы.

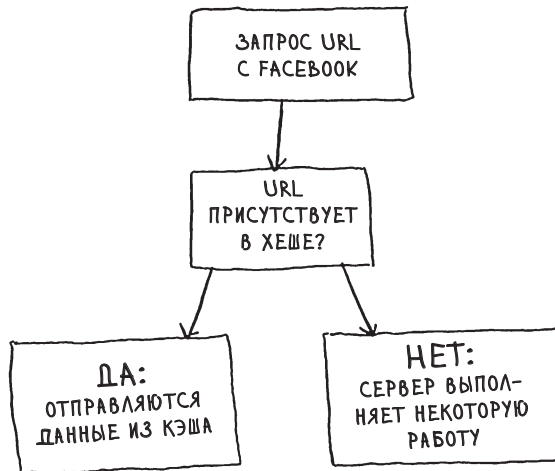
Кэширование — стандартный способ ускорения работы. Все крупные веб-сайты применяют кэширование. А кэшируемые данные хранятся в хеше!

Facebook не просто кэширует домашнюю страницу. Также кэшируются страницы «О нас», «Условия использования» и многие другие. Следовательно, необходимо создать связь URL-адреса страницы и данных страницы.

`facebook.com/about` → ДАННЫЕ СТРАНИЦЫ С ИНФОРМАЦИЕЙ О FACEBOOK

`facebook.com` → ДАННЫЕ ДОМАШНЕЙ СТРАНИЦЫ

Когда вы посещаете страницу на сайте Facebook, сайт сначала проверяет, хранится ли страница в кэше.



Вот как это выглядит в коде:

```

cache = {}

def get_page(url):
    if cache.get(url):
        return cache[url]  ← Возвращаются кэшированные данные
    else:
        data = get_data_from_server(url)
        cache[url] = data  ← Данные сначала сохраняются в кэше
        return data
  
```

Здесь сервер выполняет работу только в том случае, если URL не хранится в кэше. Однако перед тем, как возвращать данные, вы сохраняете их в кэше. Когда пользователь в следующий раз запросит тот же URL-адрес, данные

можно отправить из кэша (вместо того, чтобы заставлять сервер выполнять работу).

Шпаргалка

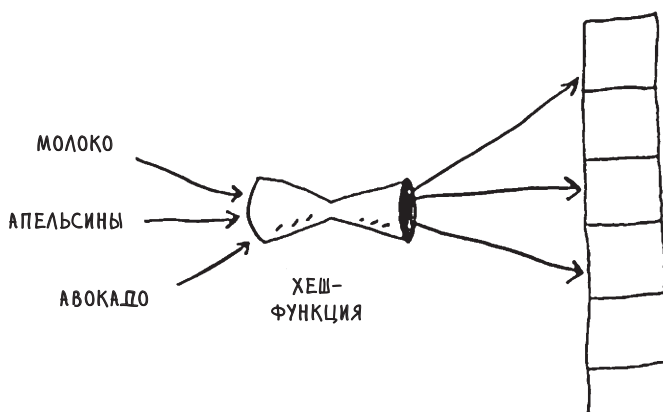
Хеши хорошо подходят для решения следующих задач:

- моделирование отношений между объектами;
- устранение дубликатов;
- кэширование/запоминание данных вместо выполнения работы на сервере.

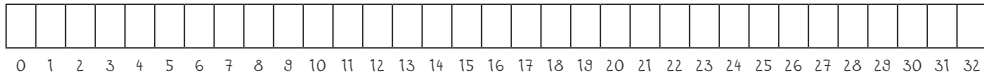
Коллизии

Как я уже сказал, в большинстве языков существуют свои хеш-таблицы. Вам не нужно знать, как написать собственную реализацию, поэтому я не буду надолго останавливаться на внутреннем строении хеш-таблиц. Но быстродействие-то важно всегда! Чтобы понять быстродействие хеш-таблиц, необходимо сначала понять, что такое коллизии. В следующих двух разделах рассматриваются коллизии и быстродействие хеш-таблиц.

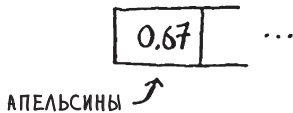
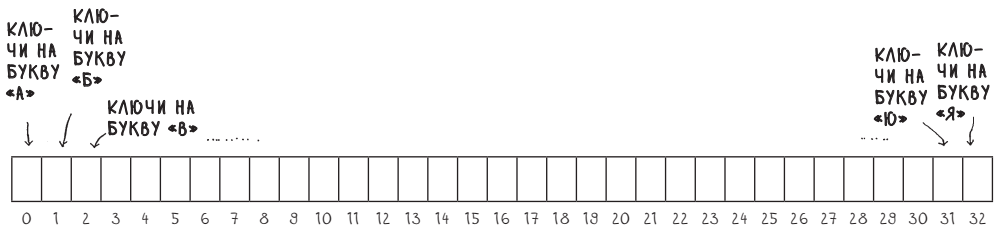
Прежде всего я немножко приукрасил действительность. Я сказал, что хеш-функция всегда отображает разные ключи на разные позиции в массиве.



На самом деле написать такую хеш-функцию почти невозможно. Рассмотрим простой пример: допустим, массив состоит всего из 33 ячеек.

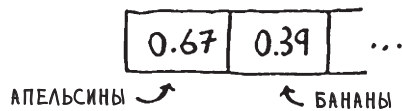


И хеш-функция очень простая: элемент массива просто назначается по алфавитному признаку.

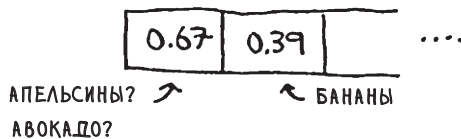


Может быть, вы уже поняли суть проблемы. Вы хотите поместить цену апельсинов в хеш. Для этого выделяется первая ячейка.

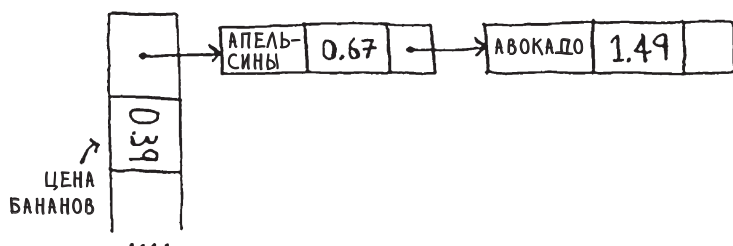
После апельсинов в хеш заносится цена бананов. Для бананов выделяется вторая ячейка.



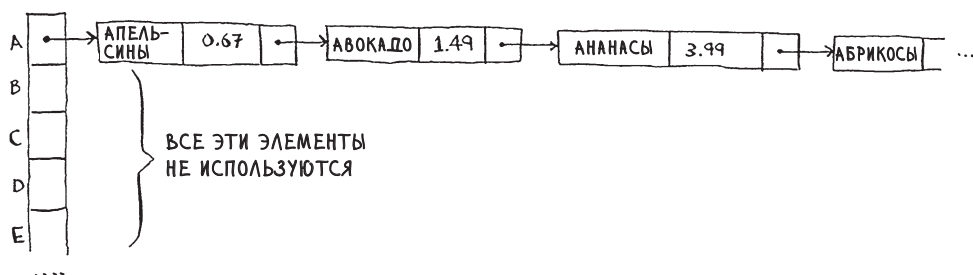
Пока все прекрасно! Но теперь в хеш нужно включить цену авокадо. И для авокадо снова выделяется первая ячейка.



О нет! Элемент уже занят апельсинами! Что же делать? Такая ситуация называется *коллизией*: двум ключам назначается один элемент массива. Возникает проблема: если сохранить в этом элементе цену авокадо, то она запишется на место цены апельсинов. И когда кто-нибудь спросит, сколько стоят апельсины, вы вместо этого сообщите цену авокадо! Коллизии — неприятная штука, и вам придется как-то разбираться с ними. Существует много разных стратегий обработки коллизий. Простейшая из них выглядит так: если несколько ключей отображаются на один элемент, в этом элементе создается связанный список.



В этом примере и «апельсины», и «авокадо» отображаются на один элемент массива, поэтому в элементе создается связанный список. Если вам потребуется узнать цену бананов, эта операция по-прежнему выполнится быстро. Если потребуется узнать цену апельсинов, работа пойдет чуть медленнее. Вам придется провести поиск по связанному списку, чтобы найти в нем «апельсины». Если связанный список мал, это не так страшно — поиск будет ограничен тремя или четырьмя элементами. Но предположим, что вы работаете в специализированной лавке, в которой продаются только продукты на букву «а».



Одну минуту! Вся хеш-таблица полностью пуста, кроме одной ячейки. И эта ячейка содержит огромный связанный список! Каждый элемент этой хеш-таблицы хранится в связанном списке. Ситуация ничуть не лучше той, когда все данные сразу хранятся в связанном списке. Работа с данными замедляется.

Из этого примера следуют два важных урока:

- *выбор хеш-функции действительно важен.* Хеш-функция, отображающая все ключи на один элемент массива, никуда не годится. В идеале хеш-функция должна распределять ключи равномерно по всему хешу;
- если связанные списки становятся слишком длинными, работа с хеш-таблицей сильно замедляется. Но они не станут слишком длинными *при использовании хорошей хеш-функции!*

Хеш-функции играют важную роль. Хорошая хеш-функция создает минимальное число коллизий. Как же выбрать хорошую хеш-функцию? Об этом в следующем разделе!

Быстродействие

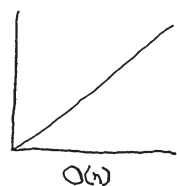
Глава началась с примера магазинчика. Вы хотели построить механизм, который мгновенно выдает цены на продукты. Что ж, хеш-таблицы работают очень быстро.

	СРЕДНИЙ СЛУЧАЙ	ХУДШИЙ СЛУЧАЙ
ПОИСК	$O(1)$	$O(n)$
ВСТАВКА	$O(1)$	$O(n)$
УДАЛЕНИЕ	$O(1)$	$O(n)$

БЫСТРОДЕЙСТВИЕ
ХЕШ-ТАБЛИЦ

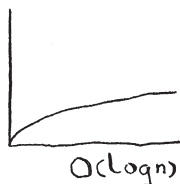
В среднем хеш-таблицы выполняют любые операции за время $O(1)$. Время $O(1)$ называется *постоянным*. Ранее примеры постоянного времени вам еще

не встречались. Оно не означает, что операции выполняются мгновенно; просто время остается постоянным независимо от размера хеш-таблицы. Например, вы знаете, что простой поиск выполняется за линейное время.



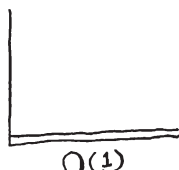
$O(n)$
ЛИНЕЙНОЕ ВРЕМЯ
(ПРОСТОЙ ПОИСК)

Бинарный поиск работает быстрее — за логарифмическое время:



$O(\log n)$
ЛОГАРИФМИЧЕСКОЕ ВРЕМЯ
(БИНАРНЫЙ ПОИСК)

Поиск данных в хеш-таблице выполняется за постоянное время.



$O(1)$
ПОСТОЯННОЕ ВРЕМЯ
(ХЕШ-ТАБЛИЦЫ)

Видите горизонтальную линию? Она означает, что при любом размере хеш-таблицы — 1 элемент или 1 миллиард элементов — выборка данных займет одинаковое время. На самом деле вы уже сталкивались с постоянным временем: получение элемента из массива выполняется за постоянное

время. От размера массива оно не зависит. В среднем случае хеш-таблицы работают действительно быстро.

В худшем случае все операции с хеш-таблицей выполняются за время $O(n)$ (линейное время), а это очень медленно. Сравним хеш-таблицы с массивами и списками.

	ХЕШ- ТАБЛИЦЫ (СРЕДНИЙ СЛУЧАЙ)	ХЕШ- ТАБЛИЦЫ (ХУДШИЙ СЛУЧАЙ)	МАС- СИВЫ	СВЯ- ЗАННЫЕ СПИСКИ
ПОИСК	$O(1)$	$O(n)$	$O(1)$	$O(n)$
ВСТАВКА	$O(1)$	$O(n)$	$O(n)$	$O(1)$
УДАЛЕ- НИЕ	$O(1)$	$O(n)$	$O(n)$	$O(1)$

Взгляните на средний случай для хеш-таблиц. При поиске хеш-таблицы не уступают в скорости массивам (получение значения по индексу). А при вставке и удалении они так же быстры, как и связанные списки. Получается, что они взяли лучшее от обеих структур! Но в худшем случае хеш-таблицы медленно выполняют все эти операции, поэтому очень важно избегать худшего случая быстродействия при работе с хеш-таблицами. А для этого следует избегать коллизий. Для предотвращения коллизий необходимы:

- низкий коэффициент заполнения;
- хорошая хеш-функция.

ПРИМЕЧАНИЕ

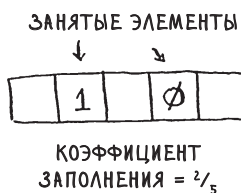
Материал следующего раздела не является обязательным. Речь пойдет о том, как реализовать хеш-таблицу, но вам никогда не придется делать это самостоятельно. Какой бы язык программирования вы ни выбрали, в нем найдется готовая реализация хеш-таблиц. Вы можете воспользоваться встроенной реализацией хеш-таблицы, не сомневаясь в том, что она имеет хорошую эффективность. А в следующем разделе мы заглянем во внутреннее устройство хеш-таблиц.

Коэффициент заполнения

Коэффициент заполнения хеш-таблицы вычисляется по простой формуле.

$$\frac{\text{КОЛИЧЕСТВО ЭЛЕМЕНТОВ
В ХЕШ-ТАБЛИЦЕ}}{\text{ОБЩЕЕ КОЛИЧЕСТВО
ЭЛЕМЕНТОВ}}$$

Хеш-таблицы используют массив для хранения данных, поэтому для вычисления коэффициента заполнения можно подсчитать количество заполненных элементов в массиве. Например, в следующей хеш-таблице коэффициент заполнения равен $\frac{2}{5}$, или 0,4.

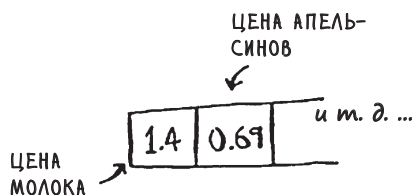


Скажите, каков коэффициент заполнения этой таблицы?



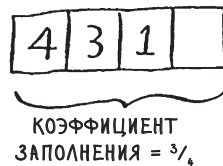
Если вы ответили « $\frac{1}{3}$ » — все правильно. По коэффициенту заполнения можно оценить, насколько заполнена хеш-таблица.

Предположим, в хеш-таблице нужно сохранить цены 100 товаров и хеш-таблица состоит из 100 элементов. В лучшем случае каждому товару будет выделен отдельный элемент.

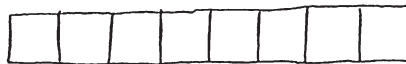


Коэффициент заполнения этой хеш-таблицы равен 1. А если хеш-таблица состоит всего из 50 элементов? Тогда ее коэффициент заполнения будет равен 2. Выделить под каждый товар отдельный элемент ни при каких условиях не удастся, потому что элементов попросту не хватит! Коэффициент заполнения больше 1 означает, что количество товаров превышает количество элементов в массиве.

С ростом коэффициента заполнения в хеш-таблицу приходится добавлять новые элементы, то есть изменять ее размер. Представим, что эта хеш-таблица приближается к заполнению.



Хеш-таблицу необходимо расширить. Расширение начинается с создания нового массива большего размера. Обычно в таком случае создается массив вдвое большего размера.



Теперь все эти элементы необходимо заново вставить в новую хеш-таблицу функцией hash:

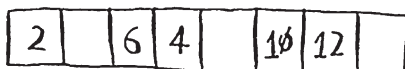


Новая таблица имеет коэффициент заполнения $\frac{3}{8}$. Гораздо лучше! С меньшим коэффициентом загрузки число коллизий уменьшается, и ваша таблица начинает работать более эффективно. Хорошее приближенное правило:

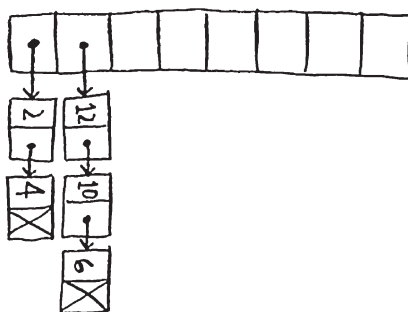
изменяйте размер хеш-таблицы, когда коэффициент заполнения превышает 0,7. Но ведь на изменение размеров уходит много времени, скажете вы, и будете абсолютно правы! Да, изменение размеров требует значительных затрат ресурсов, поэтому оно не должно происходить слишком часто. В среднем хеш-таблицы работают за время $O(1)$ даже с изменением размеров.

Хорошая хеш-функция

Хорошая хеш-функция должна обеспечивать равномерное распределение значений в массиве.



Плохая хеш-функция создает скопления и порождает множество коллизий.



Какую хеш-функцию считать хорошей? К счастью, вам об этом никогда не придется беспокоиться — пусть об этом думают умные люди в темных кабинетах. Если вам действительно интересно, рассмотрите функцию CityHash. Именно ее использует Google Abseil. Это библиотека C++ с открытым исходным кодом, основанная на внутреннем коде Google. Она содержит все общие функции C++. Abseil лежит в основе кода Google, поэтому если в ней используется CityHash, можно быть уверенным, что эта функция неплоха. Смело применяйте ее.

УПРАЖНЕНИЯ

Очень важно, чтобы хеш-функции обеспечивали хорошее распределение. Они должны распределять значения как можно шире. Худший случай — хеш-функция, которая отображает все значения на одну позицию в хеш-таблице.

Предположим, имеются четыре хеш-функции, которые получают строки.

1. Первая функция возвращает «1» для любого входного значения.
2. Вторая функция возвращает длину строки в качестве индекса.
3. Третья функция возвращает первый символ строки в качестве индекса. Таким образом, все строки, начинающиеся с «a», хешируются в одну позицию, все строки, начинающиеся с «b», — в другую и т. д.
4. Четвертая функция ставит в соответствие каждой букве простое число: $a = 2$, $b = 3$, $c = 5$, $d = 7$, $e = 11$ и т. д. Для строки хеш-функцией становится остаток от деления суммы всех значений на размер хеша. Например, если размер хеша равен 10, то для строки «bag» будет вычислен индекс $3 + 2 + 17 \% 10 = 22 \% 10 = 2$.

В каком из этих примеров хеш-функции будут обеспечивать хорошее распределение? Считайте, что хеш-таблица содержит 10 элементов.

- 5.5** Телефонная книга, в которой ключами являются имена, а значениями — номера телефонов. Задан следующий список имен: Esther, Ben, Bob, Dan.
- 5.6** Связь размера батарейки с напряжением. Размеры батареек: А, АА, ААА, АААА.
- 5.7** Связь названий книг с именами авторов. Названия книг: «Maus», «Fun Home», «Watchmen».

Шпаргалка

Вам почти никогда не придется реализовывать хеш-таблицу самостоятельно. Язык программирования, который вы используете, должен предоставить необходимую реализацию. Вы можете пользоваться хеш-таблицами Python, и при этом вам будет обеспечена производительность среднего случая: постоянное время.

Хеш-таблицы чрезвычайно полезны, потому что они обеспечивают высокую скорость операций и позволяют по-разному моделировать данные. Возможно, вскоре выяснится, что вы постоянно используете их в своей работе.

- Хеш-таблица создается объединением хеш-функции с массивом.
- Коллизии нежелательны. Хеш-функция должна свести количество коллизий к минимуму.
- Хеш-таблицы обеспечивают очень быстрое выполнение поиска, вставки и удаления.
- Хеш-таблицы хорошо подходят для моделирования отношений между объектами.
- Как только коэффициент заполнения превышает 0,7, пора изменять размер хеш-таблицы.
- Хеш-таблицы используются для кэширования данных (например, на веб-серверах).
- Хеш-таблицы хорошо подходят для обнаружения дубликатов.

6

Поиск в ширину



В этой главе

- ✓ Вы научитесь моделировать сети при помощи новой абстрактной структуры данных — графов.
- ✓ Вы освоите поиск в ширину — алгоритм, который применяется к графам для получения ответов на вопросы вида «Какой кратчайший путь ведет к X?»
- ✓ Вы узнаете, чем направленные графы отличаются от ненаправленных.
- ✓ Вы освоите топологическую сортировку — другой алгоритм сортировки, раскрывающий связи между узлами.

Эта глава посвящена графам. Сначала вы узнаете, что такое граф. Затем я покажу первый алгоритм, работающий с графами. Он называется *поиском в ширину* (BFS, Breadth-First Search).

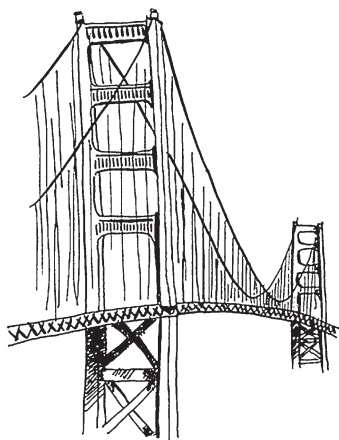
Поиск в ширину позволяет найти кратчайшее расстояние между двумя объектами. Однако сам термин «кратчайшее расстояние» может иметь много разных значений! Например, с помощью поиска в ширину можно:

- реализовать проверку правописания (минимальное количество изменений, преобразующих ошибочно написанное слово в правильное, например АЛГОРИФМ -> АЛГОРИТМ — одно изменение);

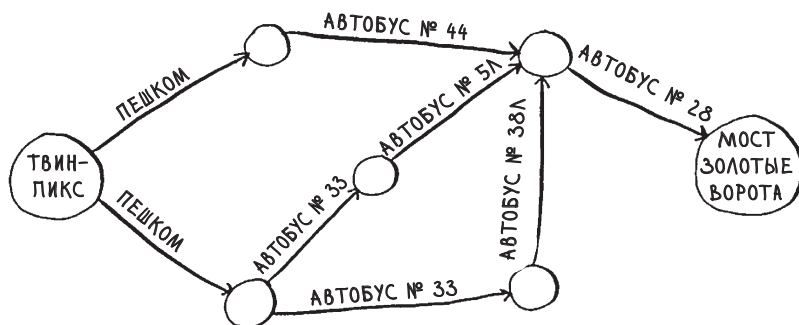
- найти ближайшего к вам врача;
- создать поискового робота.

Одни из самых полезных алгоритмов, известных мне, работают с графами. Внимательно прочитайте несколько следующих глав — этот материал неоднократно пригодится вам в работе.

Знакомство с графами

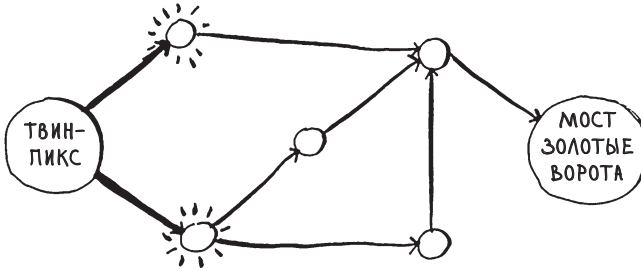


Предположим, вы находитесь в Сан-Франциско и хотите добраться из Твин-Пикс к мосту Золотые Ворота. Вы намереваетесь доехать на автобусе с минимальным количеством пересадок. Возможные варианты:

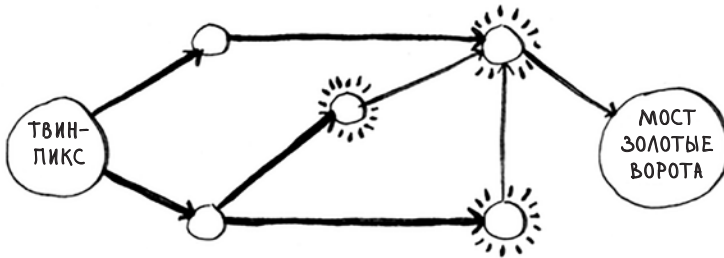


Какой алгоритм вы использовали бы для поиска пути с наименьшим количеством шагов?

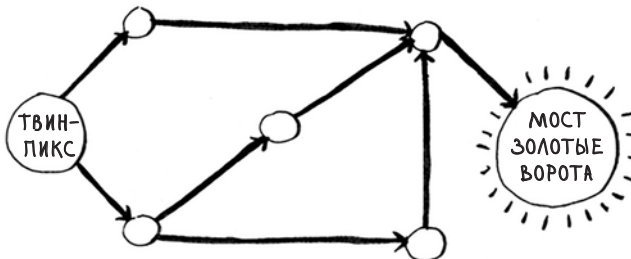
Можно ли сделать это за один шаг? На следующем рисунке выделены все места, в которые можно добраться за один шаг.



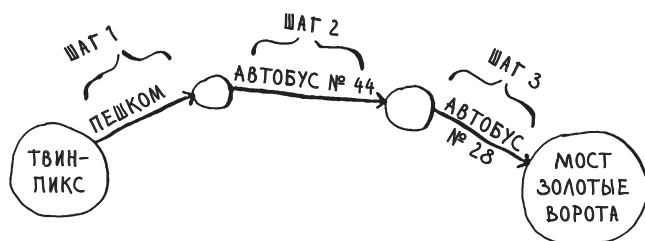
Мост на этой схеме не выделен; до него невозможно добраться за один шаг. А можно ли добраться до него за два шага?



И снова мост не выделен, а значит, до него невозможно добраться за два шага. Как насчет трех шагов?



Ага! На этот раз мост Золотые Ворота выделен. Следовательно, чтобы добраться из Твин-Пикс к мосту по этому маршруту, необходимо сделать три шага.



Есть и другие маршруты, которые приведут вас к мосту, но они длиннее (четыре шага). Алгоритм обнаружил, что кратчайший путь к мосту состоит из трех шагов. Задача такого типа называется *задачей поиска кратчайшего пути*. Часто требуется найти некий кратчайший путь: путь к дому вашего друга, путь к победе в шахматной партии (за наименьшее количество ходов) и т. д. Алгоритм для решения задачи поиска кратчайшего пути называется *поиском в ширину*.

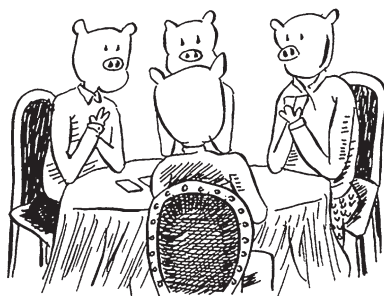
Чтобы найти кратчайший путь из Твин-Пикс к мосту Золотые Ворота, нам пришлось выполнить два шага:

- 1) смоделировать задачу в виде графа;
- 2) решить задачу методом поиска в ширину.

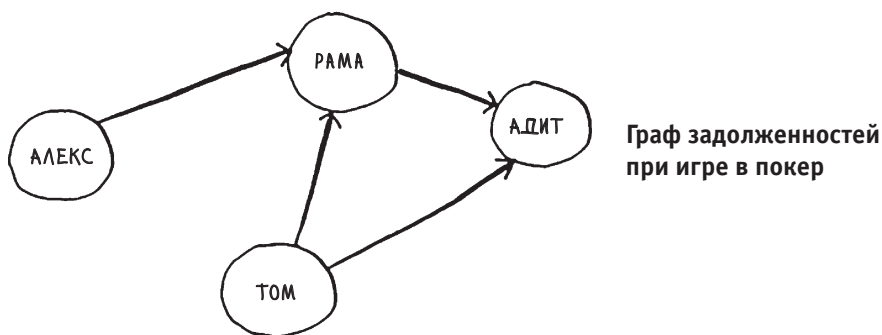
В следующем разделе я расскажу, что такое графы. Затем будет рассмотрен более подробно поиск в ширину.

Что такое граф?

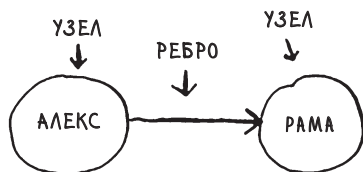
Граф моделирует набор связей. Представьте, что вы с друзьями играете в покер и хотите смоделировать, кто кому сейчас должен. Например, условие «Алекс должен Раме» можно смоделировать так:



А полный граф может выглядеть так:

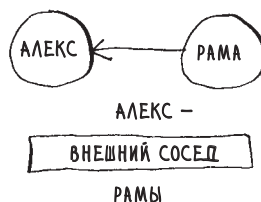
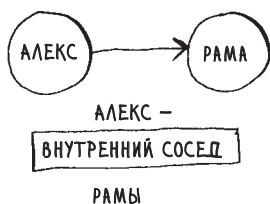


Алекс должен Раме, Том должен Адиту и т. д. Каждый граф состоит из *узлов* и *ребер*.



Вот и все! Графы состоят из узлов и ребер. Узел может быть напрямую соединен с несколькими другими узлами. Эти узлы называются *внутренними* или *внешними соседями*.

Поскольку связь идет в направлении от Алекса к Раме, то Алекс — внутренний сосед Рамы, А Рама — внешний сосед Алекса. Терминология непривычная, поэтому иллюстрация поможет разобраться.



На ней видно, что Адит не является внешним или внутренним соседом Алекса, потому что они не связаны напрямую. Но он является внешним соседом Рамы и Тома.

Графы используются для моделирования связей между разными объектами. А теперь посмотрим, как работает поиск в ширину.

Поиск в ширину

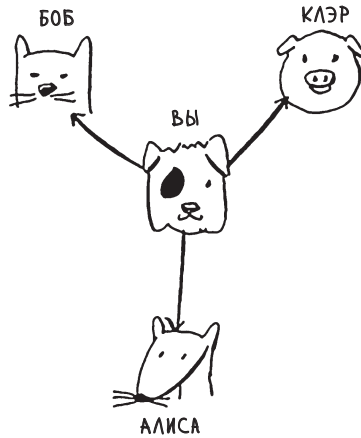
В главе 1 уже рассматривался пример алгоритма поиска: бинарный поиск. Поиск в ширину также относится к категории алгоритмов поиска, но этот алгоритм работает с графами. Он помогает ответить на вопросы двух типов:

- тип 1: существует ли путь от узла А к узлу В?
- тип 2: как выглядит кратчайший путь от узла А к узлу В?

Вы уже видели пример поиска в ширину, когда мы просчитывали кратчайший путь из Твин-Пикс к мосту Золотые Ворота. Это был вопрос типа 2: как выглядит кратчайший путь? Теперь разберем работу алгоритма более подробно с вопросом типа 1: существует ли путь?



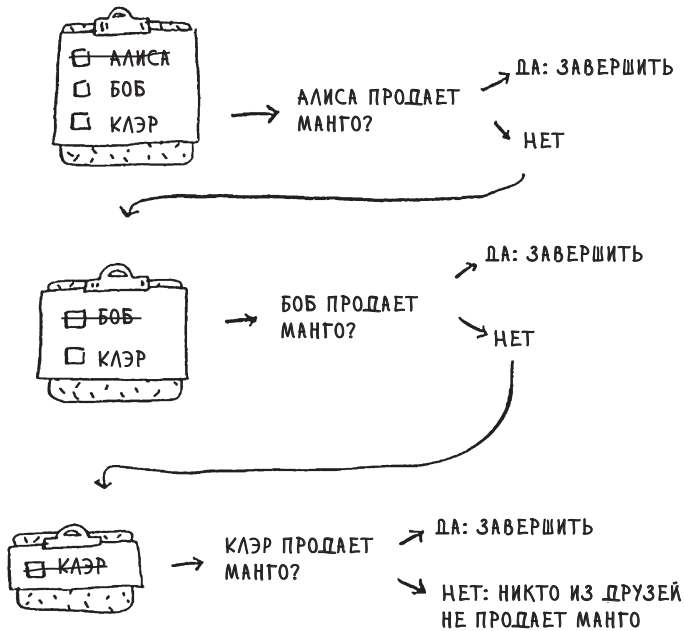
Представьте, что вы выращиваете манго. Вы ищете продавца, который будет продавать ваши замечательные манго. А может, продавец найдется среди ваших контактов на Facebook? Для начала стоит поискать среди друзей.



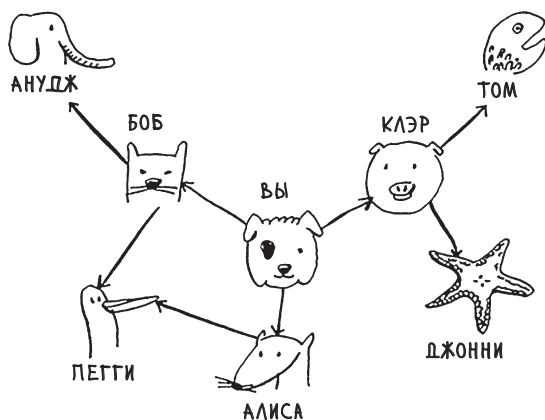
Поиск происходит вполне тривиально.

Сначала нужно построить список друзей для поиска.

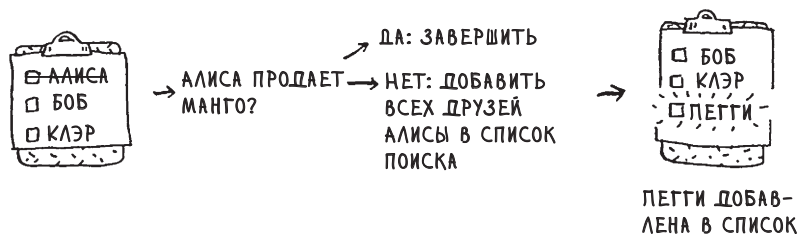
Теперь нужно обратиться к каждому человеку в списке и проверить, продает ли этот человек манго.



Предположим, ни один из ваших друзей не продает манго. Теперь поиск продолжается среди друзей ваших друзей.



Каждый раз, когда вы проверяете кого-то из списка, вы добавляете в список всех его друзей.



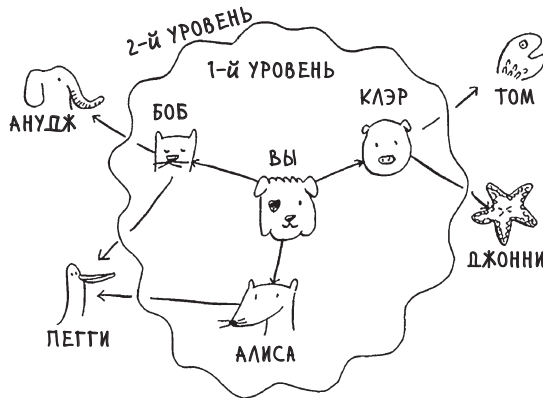
В таком случае поиск ведется не только среди друзей, но и среди друзей друзей тоже. Напомним: нужно найти в сети хотя бы одного продавца манго. Если Алиса не продает манго, то в список добавляются ее друзья. Это означает, что со временем вы проверите всех ее друзей, а потом их друзей и т. д. С этим алгоритмом поиск рано или поздно пройдет по всей сети, пока вы все-таки не наткнетесь на продавца манго. Такой алгоритм и называется поиском в ширину.

Поиск кратчайшего пути

На всякий случай напомним два вопроса, на которые может ответить алгоритм поиска в ширину:

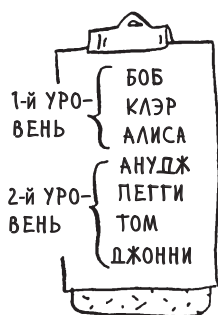
- тип 1: существует ли путь от узла А к узлу В? (Есть ли продавец манго в вашей сети?)
- тип 2: как выглядит кратчайший путь от узла А к узлу В? (Кто из продавцов манго находится ближе всего к вам?)

Вы уже знаете, как ответить на вопрос 1; теперь попробуем ответить на вопрос 2. Удается ли вам найти ближайшего продавца манго? Будем считать, что ваши друзья — это связи первого уровня, а друзья друзей — связи второго уровня.



Связи первого уровня предпочтительнее связей второго уровня, связи второго уровня предпочтительнее связей третьего уровня и т. д. Отсюда следует, что поиск по контактам второго уровня не должен производиться, пока вы не будете полностью уверены в том, что среди связей первого уровня нет ни одного продавца манго. Но ведь поиск в ширину именно это и делает! Поиск в ширину распространяется от начальной точки. А это означает, что связи первого уровня будут проверены до связей второго уровня. Контрольный вопрос: кто будет проверен первым, Клар или Анудж? Ответ: Клар является связью первого уровня, а Анудж — связью второго уровня. Следовательно, Клар будет проверена первой.

Также можно объяснить это иначе: связи первого уровня добавляются в список поиска раньше связей второго уровня.



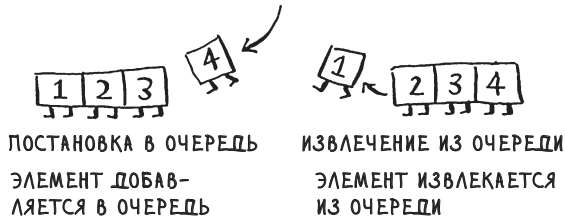
Вы двигаетесь вниз по списку и проверяете каждого человека (является ли он продавцом манго). Связи первого уровня будут проверены до связей второго уровня, так что вы найдете продавца манго, ближайшего к вам. Поиск в ширину находит не только путь из А в В, но и кратчайший путь.

Обратите внимание: это условие выполняется только в том случае, если поиск осуществляется в порядке добавления людей. Другими словами, если Клэр была добавлена в список до Ануджа, то проверка Клэр должна быть выполнена до проверки Ануджа. А что произойдет, если вы проверите Ануджа раньше, чем Клэр, и оба они окажутся продавцами манго? Анудж является связью второго уровня, а Клэр — связью первого уровня. В результате будет найден продавец манго, не ближайший к вам в сети. Следовательно, проверять связи нужно в порядке их добавления. Для операций такого рода существует специальная структура данных, которая называется *очередью*.

Очереди

Очередь работает точно так же, как и в реальной жизни. Предположим, вы с другом стоите в очереди на автобусной остановке. Если вы стоите ближе к началу очереди, то первым сядете в автобус. Структура данных очереди работает аналогично. Очереди чем-то похожи на стеки: вы не можете обращаться к произвольным элементам очереди. Вместо этого поддерживаются всего две операции: *постановка в очередь* и *извлечение из очереди*.





Если вы поставите в очередь два элемента, то элемент, добавленный первым, будет извлечен из очереди раньше второго. А ведь это свойство можно использовать для реализации списка поиска! Люди, добавленные в список первыми, будут извлечены из очереди и проверены первыми.

Очередь относится к категории структур данных FIFO: First In, First Out («первым вошел, первым вышел»). А стек принадлежит к числу структур данных LIFO: Last In, First Out («последним пришел, первым вышел»).

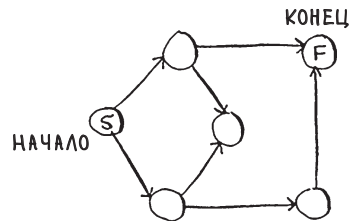


Теперь, когда вы знаете, как работает очередь, можно переходить к реализации поиска в ширину!

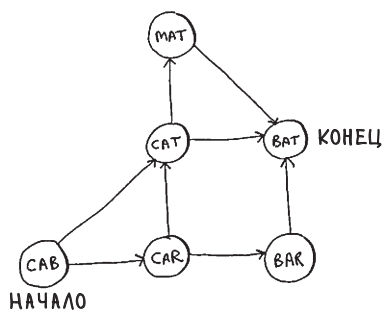
УПРАЖНЕНИЯ

Примените алгоритм поиска в ширину к каждому из этих графов, чтобы найти решение.

- 6.1** Найдите длину кратчайшего пути от начального до конечного узла.



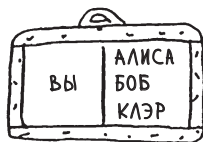
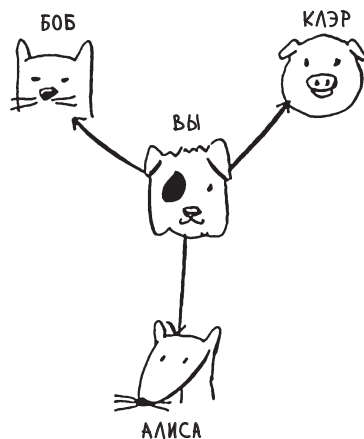
6.2 Найдите длину кратчайшего пути от «cab» к «bat».



Реализация графа

Для начала необходимо реализовать граф на программном уровне. Граф состоит из нескольких узлов. И каждый узел соединяется с соседними узлами. Как выразить отношение типа «вы \rightarrow Боб»? К счастью, вам уже известна структура данных, способная выражать отношения: *хеш-таблица*!

Вспомните: хеш-таблица связывает ключ со значением. В данном случае узел должен быть связан со всеми его внешними соседями.

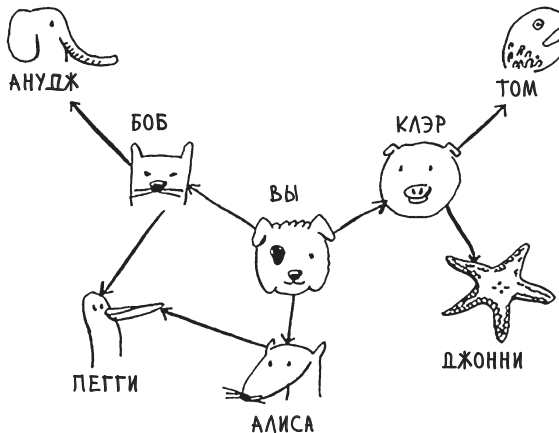


А вот как это записывается на Python:

```
graph = {}
graph["вы"] = ["Алиса", "Боб", "Клэр"]
```

Обратите внимание: элемент «вы» (you) отображается на массив. Следовательно, результатом выражения `graph["вы"]` является массив всех ваших внешних соседей. Помните, что внешние соседи — это узлы, к которым обращается узел «вы».

Граф — всего лишь набор узлов и ребер, поэтому для представления графа на Python ничего больше не потребуется. А как насчет большего графа, например такого?



Код на языке Python выглядит так:

```
graph = {}
graph["вы"] = ["Алиса", "Боб", "Клэр"]
graph["Боб"] = ["Анудж", "Пегги"]
graph["Алиса"] = ["Пегги"]
graph["Клэр"] = ["Том", "Джонни"]
graph["Анудж"] = []
graph["Пегги"] = []
graph["Том"] = []
graph["Джонни"] = []
```

Контрольный вопрос: важен ли порядок добавления пар «ключ — значение»?

Важно ли, какую запись вы будете использовать, — такую:

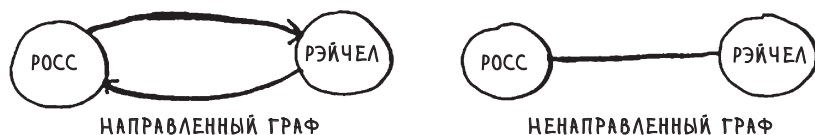
```
graph["Клэр"] = ["Том", "Джонни"]
graph["Анудж"] = []
```

или такую:

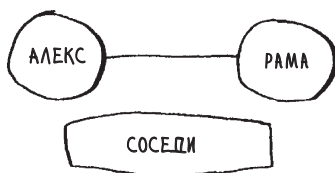
```
graph["Анудж"] = []
graph["Клэр"] = ["Том", "Джонни"]
```

Вспомните предыдущую главу. Ответ: нет, неважно. В хеш-таблицах элементы не упорядочены, поэтому добавлять пары «ключ — значение» можно в любом порядке.

У Ануджа, Пегги, Тома и Джонни внешних соседей нет. У них есть внутренние соседи, поскольку линии со стрелками указывают на них, но не существует стрелок от них к другим узлам. Такой граф называется *направленным* — отношения действуют только в одну сторону. В ненаправленном графе стрелок нет. Например, оба следующих графа эквивалентны.

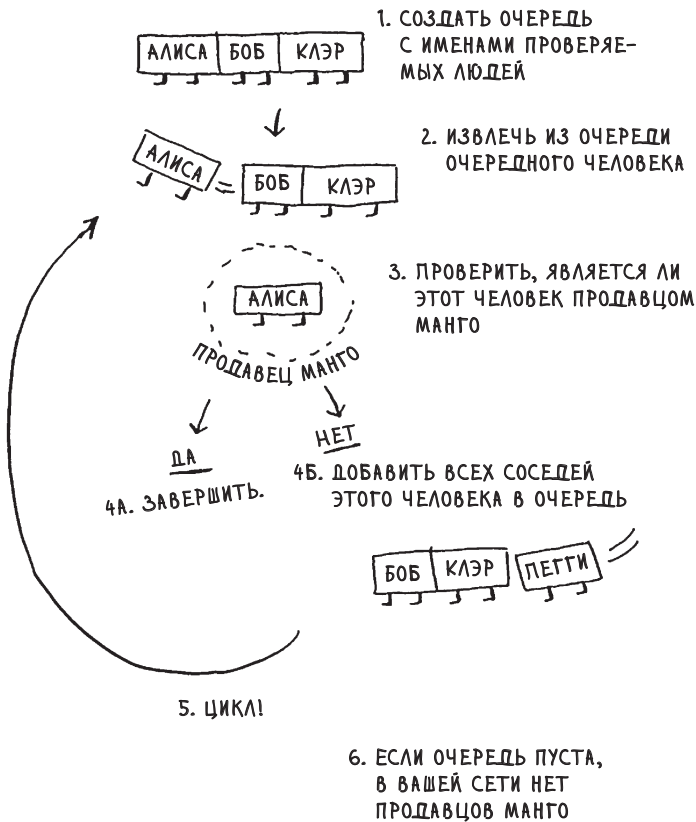


В ненаправленных графах не существует деления соседей на «внутренних» и «внешних» — это просто «соседи».



Реализация алгоритма

Напомню, как работает реализация.



Все начинается с создания очереди. В Python для создания *двусторонней очереди* (дека) используется функция `deque`:

```
from collections import deque
search_queue = deque()
search_queue += graph["вы"]
```

Создание новой очереди
Все соседи добавляются в очередь поиска

Напомню, что выражение `graph["вы"]` вернет список всех ваших соседей, например `["Алиса", "Боб", "Клэр"]`. Все они добавляются в очередь поиска.



А теперь рассмотрим остальное:

```
while search_queue:  <----- Пока очередь не пуста...
    person = search_queue.popleft()  <----- из очереди извлекается первый человек
    if person_is_seller(person):  <----- Провераем, является ли этот человек
                                    продавцом манго
        print person + " продавец манго!"  <----- Да, это продавец манго
        return True
    else:
        search_queue += graph[person]  <----- Нет, не является. Все друзья этого че-
                                                ловека добавляются в очередь поиска
return False  <----- Если выполнение дошло
                                   до этой строки, значит,
                                   в очереди нет продавца манго
```

И последнее: нужно определить функцию `person_is_seller`, которая сообщает, является ли человек продавцом манго. Например, функция может выглядеть так:

```
def person_is_seller(name):
    return name[-1] == 'm'
```

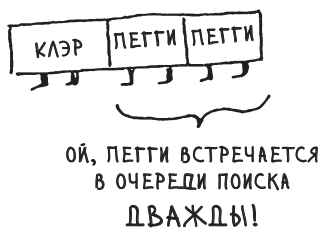
Эта функция проверяет, заканчивается ли имя на букву «m», и если заканчивается, этот человек считается продавцом манго. Проверка довольно глупая, но для нашего примера сойдет. А теперь посмотрим, как работает поиск в ширину.



И так далее. Алгоритм продолжает работать до тех пор, пока:

- не будет найден продавец манго,
- или
- очередь не опустеет (в этом случае продавца манго нет).

У Алисы и Боба есть один общий друг: Пегги. Следовательно, Пегги будет добавлена в очередь дважды: при добавлении друзей Алисы и при добавлении друзей Боба. В результате Пегги появится в очереди поиска в двух экземплярах.



Но проверить, является ли Пегги продавцом манго, достаточно всего один раз. Проверая ее дважды, вы выполняете лишнюю, ненужную работу. Следовательно, после проверки человека нужно пометить как уже проверенного.

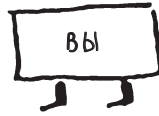
Если этого не сделать, может возникнуть бесконечный цикл. Предположим, граф выглядит так:



В начале очередь поиска содержит всех ваших соседей.



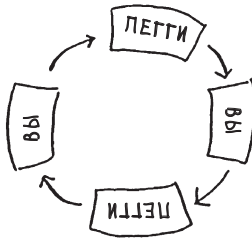
Теперь вы проверяете Пегги. Она не является продавцом манго, поэтому все ее соседи добавляются в очередь поиска.



Вы проверяете себя. Вы не являетесь продавцом манго, поэтому все ваши соседи добавляются в очередь поиска.



И так далее. Возникает бесконечный цикл, потому что очередь поиска будет поочередно переходить от вас к Пегги.



Прежде чем проверять человека, следует убедиться в том, что он не был проверен ранее. Для этого мы будем вести список уже проверенных людей.



А вот окончательная версия кода поиска в ширину, в которой учтено это обстоятельство:

```
def search(name):
    search_queue = deque()
    search_queue += graph[name]
    searched = []  # ←..... Этот массив используется для отслеживания
                   # уже проверенных людей
    while search_queue:
        person = search_queue.popleft()
        if not person in searched:  # ←..... Человек проверяется только в том случае,
            if person_is_seller(person):  # если он не проверялся ранее
                print person + " продавец манго!"
                return True
            else:
                search_queue += graph[person]
                searched.append(person)  # ←..... Человек помечается как
                                         # уже проверенный
    return False

search("вы")
```

Попробуйте выполнить этот код самостоятельно. Замените функцию `person_is_seller` чем-то более содержательным и посмотрите, выведет ли она то, что вы ожидали.

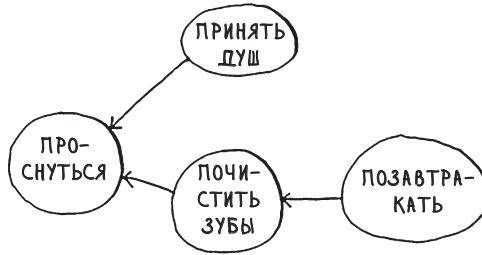
Время выполнения

Если поиск продавца манго был выполнен по всей сети, значит, вы прошли по каждому ребру (напомню: ребром называется соединительная линия или линия со стрелкой, ведущая от одного человека к другому). Таким образом, время выполнения составляет как минимум $O(\text{количество ребер})$.

Также в программе должна храниться очередь поиска. Добавление одного человека в очередь выполняется за постоянное время: $O(1)$. Выполнение операции для каждого человека потребует суммарного времени $O(\text{количество людей})$. Поиск в ширину выполняется за время $O(\text{количество людей} + \text{количество ребер})$, что обычно записывается в форме $O(V+E)$ (V — количество вершин, E — количество ребер).

УПРАЖНЕНИЯ

Перед вами небольшой граф моего утреннего распорядка.



Из графа видно, что я завтракаю только после того, как почищу зубы. Таким образом, узел «Позавтракать» зависит от узла «Почистить зубы».

А вот душ не зависит от чистки зубов, потому что я могу сначала принять душ, а потом почистить зубы. На основании графа можно сформулировать порядок, в котором я действую утром.

1. Проснуться.
2. Принять душ.
3. Почистить зубы.
4. Позавтракать.

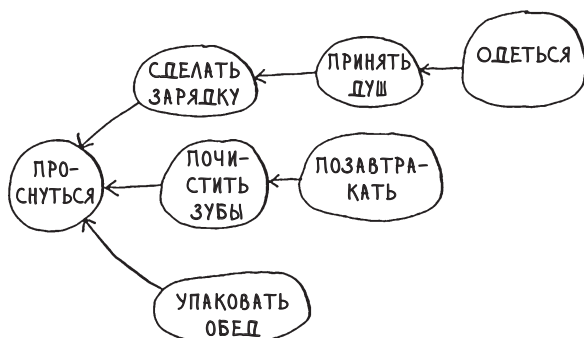
Следует заметить, что действие «Принять душ» может перемещаться в списке, поэтому следующий список тоже действителен.

1. Проснуться.
2. Почистить зубы.
3. Принять душ.
4. Позавтракать.

6.3 Для каждого из следующих трех списков укажите, действителен он или недействителен.

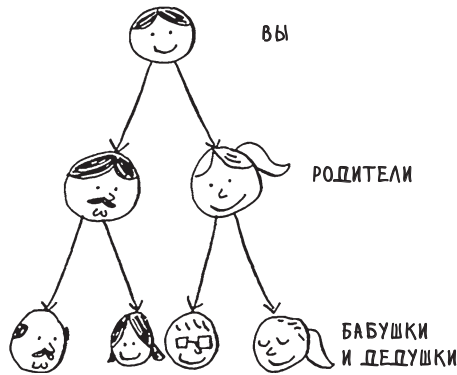
А	Б	В
1. ПРОСНУТЬСЯ	1. ПРОСНУТЬСЯ	1. ПРИНЯТЬ ДУШ
2. ПРИНЯТЬ ДУШ	2. ПОЧИСТИТЬ ЗУБЫ	2. ПРОСНУТЬСЯ
3. ПОЗАВТРАКАТЬ	3. ПОЗАВТРАКАТЬ	3. ПОЧИСТИТЬ ЗУБЫ
4. ПОЧИСТИТЬ ЗУБЫ	4. ПРИНЯТЬ ДУШ	4. ПОЗАВТРАКАТЬ

6.4 Немного увеличим исходный граф. Постройте действительный список для этого графа.

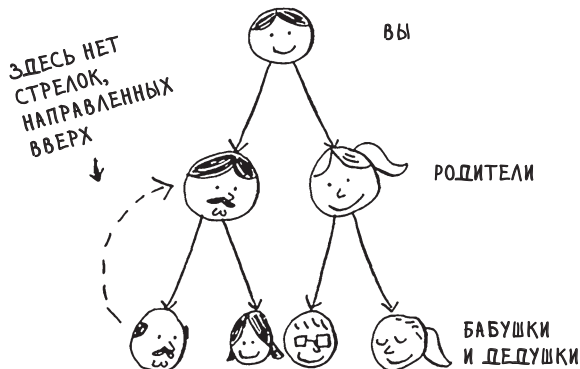


Можно сказать, что этот список в некотором смысле отсортирован. Если задача А зависит от задачи В, то задача А находится в более поздней позиции списка. Такая сортировка называется *топологической*; фактически она предоставляет способ построения упорядоченного списка на основе графа. Предположим, вы планируете свадьбу и у вас составлен большой граф со множеством задач, но вы не знаете, с чего начать. Проведите *топологическую сортировку* графа — и получите список задач, которые можно выполнять одну за другой.

Допустим, имеется генеалогическое древо.

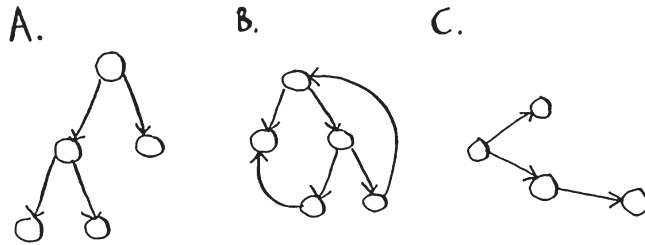


Генеалогическое древо — тоже граф, потому что в нем есть узлы (люди) и ребра. Ребра указывают на родителей человека. Естественно, все ребра направлены вниз — в генеалогическом древе ребро, указывающее вверх, не имеет смысла. Ведь ваш отец никак не может быть дедушкой вашего дедушки!



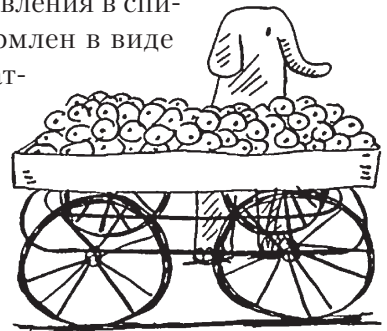
Такая особая разновидность графа, в которой нет ребер, указывающих в обратном направлении, называется *деревом*.

6.5 Какие из следующих графов также являются деревьями?



Шпаргалка

- Поиск в ширину позволяет определить, существует ли путь из А в В.
- Если путь существует, то поиск в ширину находит кратчайший путь.
- Если в вашей задаче требуется найти «кратчайшее X», попробуйте смоделировать свою задачу графом и воспользуйтесь поиском в ширину для ее решения.
- В направленном графе есть стрелки, а отношения действуют в направлении стрелки (Рама → Адит означает «Рама должен Адиту»).
- В ненаправленных графах стрелок нет, а отношение идет в обе стороны (Росс — Рэйчел означает «Росс встречается с Рэйчел, а Рэйчел встречается с Россом».)
- Очереди относятся к категории FIFO («первым вошел, первым вышел»).
- Стек относится к категории LIFO («последним пришел, первым вышел»).
- Людей следует проверять в порядке их добавления в список поиска, поэтому он должен быть оформлен в виде очереди, иначе найденный путь не будет кратчайшим.
- Позаботьтесь о том, чтобы уже проверенный человек не проверялся заново, иначе может возникнуть бесконечный цикл.



7

Деревья

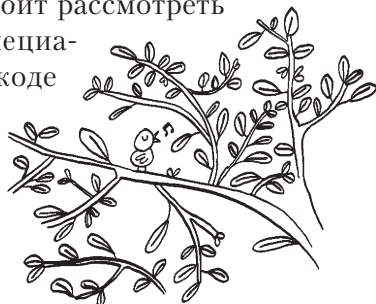


В этой главе

- ✓ Вы узнаете, что такое дерево и чем деревья отличаются от графов.
- ✓ Вы освоите выполнение алгоритмов с деревьями.
- ✓ Вы изучите поиск в глубину и узнаете, чем он отличается от поиска в ширину.
- ✓ Вы познакомитесь с кодом Хаффмана — алгоритмом сжатия, использующим деревья.

Что общего у алгоритмов сжатия и хранения информации в базе данных? В их основе часто лежит дерево, выполняющее всю черную работу. Деревья составляют подмножество графов. Деревья стоит рассмотреть отдельно, так как у них существует много специализированных разновидностей. Например, в коде Хаффмана — алгоритме сжатия, с которым вы познакомитесь в этой главе, — применяются бинарные деревья.

Многие базы данных используют сбалансированное дерево (например, B-дерево, о котором

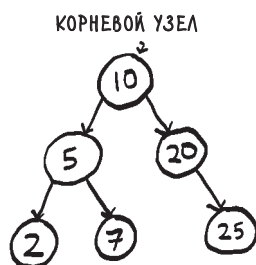
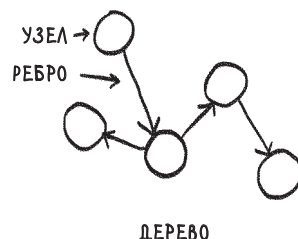


вы узнаете в следующей главе). Видов деревьев довольно много. Чтобы вам было легче разобраться в них, в главах 7 и 8 будут представлены необходимая терминология и концепции.

Ваше первое дерево

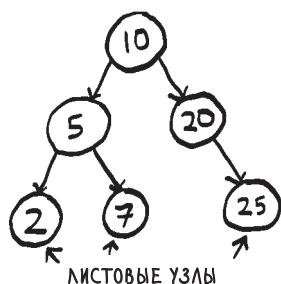
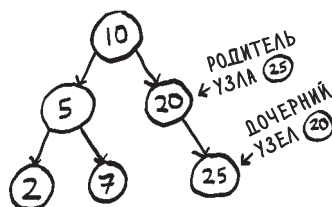
Дерево — это разновидность графа. Более подробное определение дадим ниже, а пока узнаем несколько терминов и рассмотрим пример.

Как и графы, деревья состоят из узлов и ребер.



В этой книге мы будем работать с корневыми деревьями. У корневого дерева имеется один узел, от которого можно перейти к любому другому узлу.

Мы будем работать исключительно с корневыми деревьями, так что, говоря о *дереве* в этой главе, я всегда имею в виду корневое дерево. У узлов могут быть дочерние узлы, а у дочерних узлов может быть родительский узел.



В дереве узлы имеют по крайней мере одного родителя.

Существует только один узел без родителя — это корневой узел.

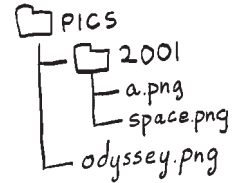
Узлы, не имеющие дочерних узлов, называются листовыми узлами (листьями).

Если вы поняли, что такое корень, лист, родительский и дочерний узел, значит, вы готовы двигаться дальше!

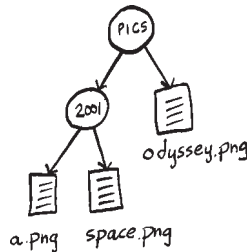
Каталоги файлов

Так как дерево является разновидностью графа, к нему можно применить алгоритм, работающий с графом. В главе 6 был рассмотрен поиск в ширину — алгоритм для нахождения кратчайшего пути в графе. Поиск в ширину можно использовать и с деревом. Если вы еще не освоились с поиском в ширину, обращайтесь к главе 6.

Каталог файлов представляет собой дерево, с которым каждый, кто работает с компьютером, сталкивается ежедневно. Представьте, что у нас есть такой каталог.



Требуется вывести имена всех файлов в каталоге `pics`, включая все его подкаталоги. Однако в данном случае существует только один подкаталог — `2001`. Задачу можно решить поиском в ширину! Для начала я покажу, как этот каталог выглядит в виде дерева.



Ранее мы использовали алгоритм поиска в ширину как инструмент поиска, однако этим его возможности не ограничиваются. Поиск в ширину является алгоритмом обхода; это значит, что он посещает каждый узел дерева («обходит» его). А это именно то, что нужно! Нам нужен алгоритм, который перейдет к каждому файлу в дереве и выведет его имя. Для перебора файлов в каталоге можно воспользоваться поиском в ширину. Алгоритм также заходит в подкаталоги, ищет в них файлы и выводит имена обнаруженных файлов. Логика выглядит так:

1. Посетить каждый узел в дереве.
2. Если узел является файлом, вывести его имя.
3. Если узел является папкой, добавить его в очередь папок, чтобы найти находящиеся в нем файлы.

Ниже приведен код реализации. Он очень похож на код поиска продавца манго из главы 6:

```
from os import listdir
from os.path import isfile, join
from collections import deque

def printnames(start_dir):
    search_queue = deque()
    search_queue.append(start_dir)
    while search_queue:
        dir = search_queue.popleft()
        for file in sorted(listdir(dir)):
            fullpath = join(dir, file)
            if isfile(fullpath):
                print(file)
            else:
                search_queue.append(fullpath)

printnames("pics")
```

Очередь используется для отслеживания папок, в которых должен проводиться поиск

Пока очередь не пуста, извлечь из нее папку для проверки

Перебрать все файлы и папки в этой папке

Если это файл, вывести его имя

Если это папка, добавить ее в конец очереди папок для поиска

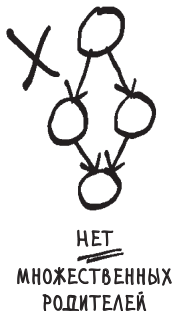
Здесь используется очередь, как и в примере с продавцом манго. В очереди хранится информация о том, в каких еще папках нужно провести поиск файлов. Конечно, в том примере мы останавливались сразу же после нахождения продавца манго, но здесь пройдем по всему дереву.

Впрочем, в этом примере присутствует одно важное отличие от кода поиска продавца манго. Удастся ли вам его найти?

Как вы помните, в примере с продавцом манго нам приходилось следить, проводился ли уже поиск этого человека:

```
...
if person not in searched:
    if person_is_seller(person):
        ...
    ...
```

Искать этого человека, только если его поиск еще не проводился



Здесь этого делать не нужно! В деревьях нет циклов, и у каждого узла только один родитель. Мы никогда не сможем случайно провести поиск в одной папке несколько раз либо ввести программу в бесконечный цикл, поэтому нет необходимости отслеживать, в каких папках мы уже искали. Зайти снова в одну и ту же папку при обходе попросту невозможно.

Благодаря этому свойству деревьев наш код стал проще. Это важный вывод, который стоит запомнить: в деревьях не бывает циклов.

О СИМВОЛИЧЕСКИХ ССЫЛКАХ

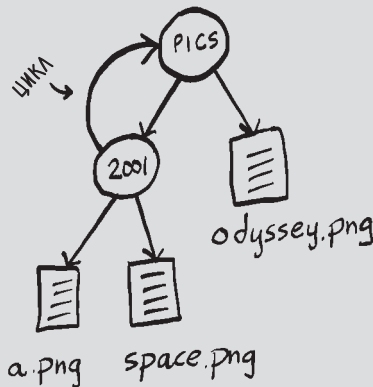
Вероятно, вы знаете, что такое символические ссылки. Для тех, кто не знает: символические ссылки открывают возможность создания циклов в каталоге файлов. Символическую ссылку в macOS или Linux можно создать командой

```
ln -s pics/ pics/2001/pics
```

а в Windows — командой

```
mklink /d pics/ pics/2001/pics
```

Если бы я добавил символическую ссылку, дерево выглядело бы примерно так:



Теперь каталог уже не является деревом! Чтобы не переусложнять, в этом примере мы не будем рассматривать символические ссылки. Если у вас они есть, Python достаточно сообразителен, чтобы не войти в бесконечный цикл. Он выдаст такую ошибку:

```
OSError: [Errno 62] Too many levels of symbolic links':
'pics/2001/pics'
```

¹ Слишком много уровней символических ссылок. — Примеч. пер.

Космическая одиссея: поиск в глубину

Еще раз обойдем каталог файлов, но на этот раз рекурсивно:

```
from os import listdir
from os.path import isfile, join

def printnames(dir):
    for file in sorted(listdir(dir)):
        fullpath = join(dir, file)
        if isfile(fullpath):
            print(file)
        else:
            printnames(fullpath)
printnames("pics")
```

Перебрать все файлы и все папки в текущей папке
 Если это файл, вывести его имя
 Если это папка, вызвать для нее функцию рекурсивно, чтобы провести поиск файлов и папок

Обратите внимание: на этот раз очередь не используется. Вместо этого при обнаружении папки мы сразу заходим в нее, чтобы найти новые файлы и папки. В нашем распоряжении появились два способа вывода имен файлов. Удивительно, но *они выводят имена файлов в разном порядке!*

Одно решение выводит имена так:

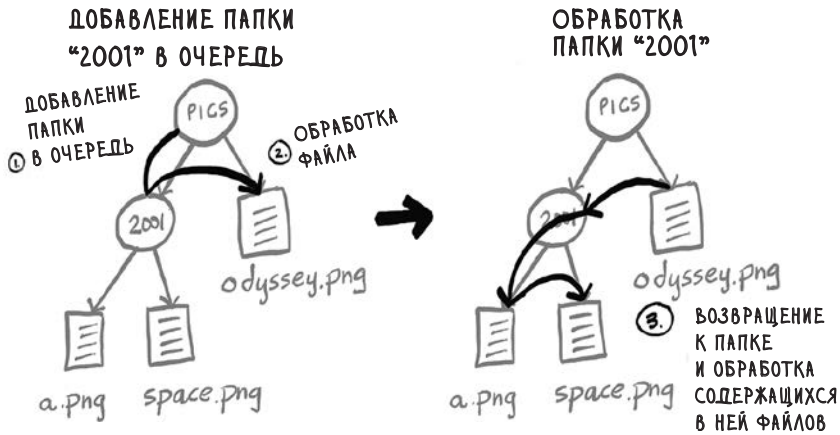
```
a.png
space.png
odyssey.png
```

А другое — так:

```
odyssey.png
a.png
space.png
```

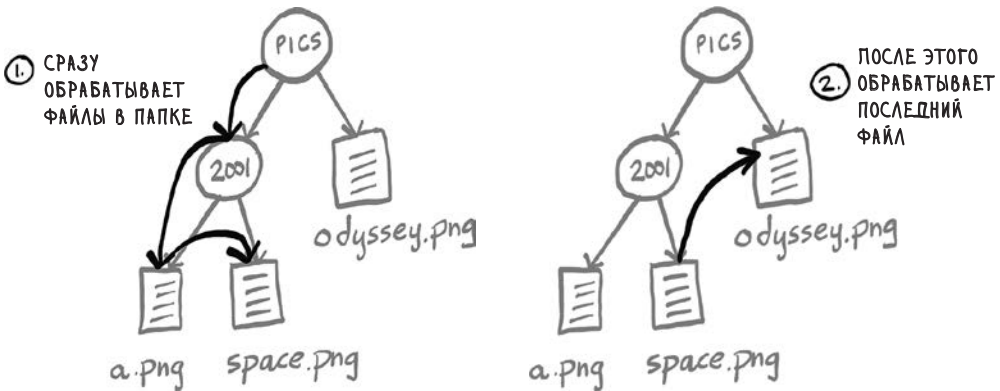
Можете ли вы определить, какое решение использует тот или иной порядок и почему? Проверьте себя, прежде чем двигаться дальше.

В первом решении используется поиск в ширину. Когда он находит папку, эта папка добавляется в очередь для проверки в будущем. Таким образом, алгоритм переходит к папке 2001, не заходит в нее, но добавляет в очередь для последующей проверки. Он выводит все имена файлов в папке pics/, после чего переходит к папке 2001/ и выводит имена содержащихся в ней файлов.



Как видите, алгоритм сначала посещает папку 2001, но не заходит в нее. Эта папка просто добавляется в очередь, а поиск в ширину переходит к `odyssey.png`.

Во втором решении используется алгоритм, называемый поиском в глубину. Поиск в глубину тоже работает с графом и представляет собой алгоритм обхода дерева. При обнаружении папки он сразу заходит в нее, вместо того чтобы добавлять ее в очередь.

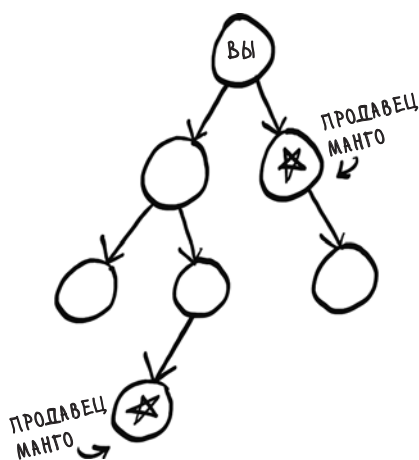


Второе решение выводит следующий результат:

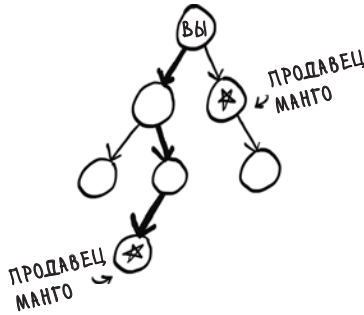
a.png
space.png
odyssey.png

Поиск в ширину и поиск в глубину тесно связаны друг с другом, и там, где речь идет об одном из них, будет упоминаться и другой. Оба алгоритма выводят все имена файлов, так что они оба подойдут для нашего примера. Однако между ними существует важное различие. Поиск в глубину не может использоваться для нахождения кратчайшего пути!

В примере с продавцом манго мы не могли применить поиск в глубину. Мы опирались на тот факт, что все друзья первого уровня проверяются раньше друзей второго уровня и т. д. Именно так работает поиск в ширину. В отличие от него, поиск в глубину сразу заходит на максимально возможную глубину. Он может сначала обнаружить продавца манго в трех уровнях от вас, хотя есть контакт ближе. Представим, что ваше дерево друзей выглядит так:

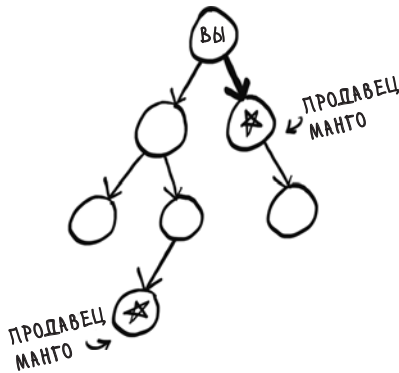


Предположим, что узлы обрабатываются слева направо. Поиск в глубину переходит к крайнему левому дочернему узлу, после чего продолжает работу.



Так как поиск зашел в глубину от левого узла, он не смог понять, что правый узел — продавец манго, находящийся намного ближе к вам.

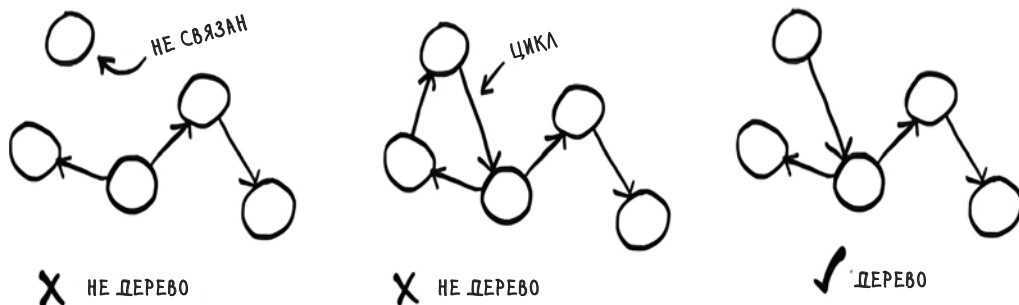
Поиск в ширину верно находит ближайшего продавца манго.



Таким образом, хотя оба алгоритма выводят списки файлов, для нахождения кратчайшего пути подходит только поиск в ширину. У поиска в глубину находятся другие применения. Его можно использовать для топологической сортировки — мы кратко рассмотрели ее в главе 6.

Правильное определение дерева

После рассмотрения примера можно привести более точное определение дерева. Дерево представляет собой *связный ациклический граф*.



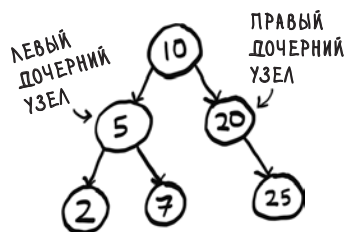
Как я уже сказал, мы работаем исключительно с корневыми деревьями, так что у всех наших деревьев есть корень. А еще мы работаем только со связными графами. Таким образом, самое важное, что следует запомнить, — *в деревьях не может быть циклов*.

Итак, вы увидели, как работает дерево. Рассмотрим подробнее один из его типов.

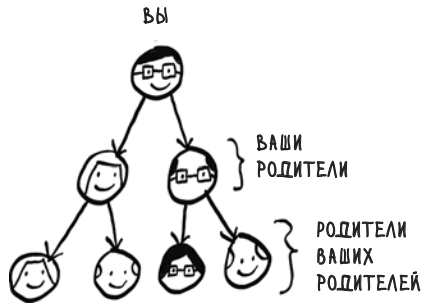
Бинарные деревья

В computer science полно разных деревьев. Бинарные деревья используются очень часто. Далее в этой и следующей главе мы будем в основном работать с бинарными деревьями.

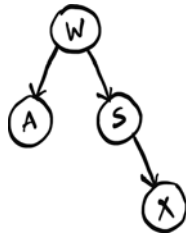
Бинарное дерево представляет собой особую разновидность дерева, узлы которого могут иметь не более двух дочерних узлов (отсюда и название). Дочерние узлы традиционно называются *левым* и *правым* узлами.



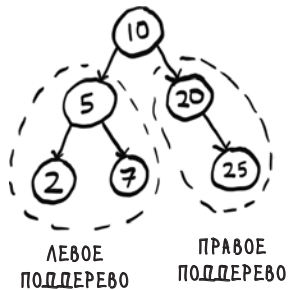
Пример бинарного дерева — генеалогическое древо, так как у каждого узла имеются два биологических родителя.



В этом примере между узлами существует четкая связь — все они составляют одну семью. Однако данные могут быть абсолютно произвольными.



Важно то, что ни у одного узла не может быть больше двух дочерних узлов. Иногда встречаются термины «левое поддерево» и «правое поддерево».



Бинарные деревья очень распространены в computer science. Далее мы разберем пример, в котором используется бинарное дерево.

Код Хаффмана

Код Хаффмана — хороший пример использования бинарных деревьев. Он также лежит в основе алгоритмов сжатия текста. Мы не будем формально описывать алгоритм, а расскажем о том, как он работает и как в нем применяются деревья.

Небольшое введение в тему. Чтобы понять, как работает сжатие, необходимо знать, сколько места занимает текстовый файл. Представим файл, содержащий всего одно слово: `tilt`. Сколько места он занимает? Чтобы узнать это, можно воспользоваться командой `stat` (в Unix). Сохраните слово в файле с именем `test.txt`, а затем воспользуйтесь командой `stat`:

```
$ cat test.txt  
tilt  
  
$ stat -f%z test.txt  
4
```

Итак, файл занимает 4 байта: по 1 байту на символ.

Выглядит логично. Если предположить, что мы используем кодировку ISO-8859-1 (в следующей врезке объясняется, что это такое), каждая буква занимает ровно 1 байт. Например, букве *a* в ISO-8859-1 соответствует код 97, который в двоичной системе записывается в виде 01100001, — итого 8 бит. Битом называется двоичная цифра с возможными значениями 0 или 1. 8 бит образуют 1 байт. Таким образом, буква *a* представлена одним байтом. Коды ISO-8859-1 содержатся в интервале от 00000000, что соответствует нуль-символу, до 11111111, что соответствует *j* (латинская буква *y* в нижнем регистре с тремой). Всего существуют 256 возможных 8-битовых комбинаций из 0 и 1, так что кодировка ISO-8859-1 способна представить до 256 букв.

КОДИРОВКИ

Как показывает этот пример, существует много способов кодирования символов. Другими словами, букву *a* можно записать в двоичной форме разными способами.

Все началось с кодировки ASCII, которая была создана в 1960-х годах. ASCII является 7-битной кодировкой. К сожалению, в ASCII не вошли многие символы, в том числе символы с диакритическими знаками (например, *ï* или *ö*) и распространенные обозначения валют, например британского фунта или японской иены.

Тогда была создана кодировка ISO-8859-1. Это 8-битная кодировка, поэтому количество символов в ней вдвое больше количества символов в ASCII. Вместо 128 символов стало доступно 256. Тем не менее этого было недостаточно, и многие страны стали создавать собственные кодировки. Например, в Японии существует несколько кодировок для японского языка, так как ISO-8859-1 и ASCII были ориентированы на европейские языки. Словом, хаос нарастал, пока не появился Юникод.

Юникод — стандарт кодирования символов, который поддерживает символы всех языков. Юникод версии 15 включает 149 186 символов — большой шаг вперед по сравнению с 256! Более 1000 из них составляют эмодзи.

Юникод является стандартом, а в программах должна использоваться кодировка, соответствующая стандарту. В настоящее время самой популярной кодировкой такого рода считается UTF-8. Это кодировка с переменной длиной кодирования символов; это означает, что представление одного символа может занимать от 1 до 4 байт (8–32 бит).

Впрочем, для вас это не так важно. Я намеренно упростил пример, используя ISO-8859-1 с 8-битным представлением символов, — это удобное представление с постоянной длиной кодирования.

Просто запомните основное:

- Алгоритмы сжатия сокращают количество битов, необходимых для представления каждого символа.
- UTF-8 — хороший вариант кодировки по умолчанию для программных проектов.

Декодируем двоичное представление в ISO-8859-1:

011100100110000101100100. Вы можете воспользоваться таблицей ISO-8859-1 или преобразователем двоичной записи в ISO-8859-1 (их можно найти в интернете), чтобы упростить задачу.

Мы знаем, что каждая буква занимает 8 бит, поэтому для начала разделим последовательность на 8-битные блоки, чтобы ее было проще читать:

```
01110010 01100001 01100100
```

Отлично, теперь мы видим, что запись состоит из трех букв. По таблице ISO-8859-1 можно найти, что это буквы *rad*: 01110010 соответствует *r*, и т. д. Так текстовый редактор берет двоичные данные из текстового файла и отображает их в виде текста ISO-8859-1. Для просмотра двоичной информации можно воспользоваться `xxd` — утилитой для Unix. Вот как текст *tilt* выглядит в двоичной форме:

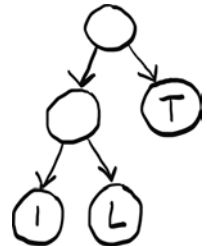
```
$ xxd -b test.txt
00000000: 01110100 01101001 01101100 01110100
tilt
```

Здесь в игру вступает сжатие. Для слова *tilt* не нужны все 256 возможных букв, достаточно трех. Таким образом, для представления буквы не требуется 8 бит, можно обойтись всего двумя. Можно было бы определить собственную 2-битную кодировку для этих трех букв:

```
t = 00
i = 01
l = 10
```

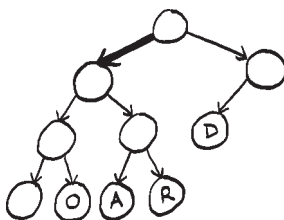
В таком случае можно записать *tilt* в новой кодировке: 00011000 (для удобства можно разделить последовательность, добавив в нее пробелы: 00 01 10 00). Сравнив ее с таблицей, вы увидите, что она соответствует исходному слову *tilt*.

Именно так работает алгоритм Хаффмана: он ищет повторяющиеся символы и использует для их представления менее 8 бит. В результате происходит сжатие данных. Код Хаффмана генерирует дерево.

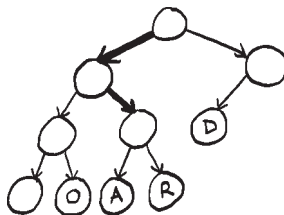


Приходится перебирать цифры по одной, словно просматривая пленку с цифрами.

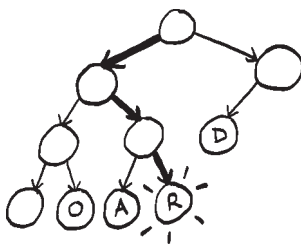
Вот как это делается: первое число равно 0, поэтому идем налево (здесь показана только часть дерева).



Далее следует 1, поэтому идем направо.



Снова 1, снова идем направо.



Буква найдена! Остались двоичные данные 01010. Вы можете снова начать от корневого узла и найти остальные буквы. Попробуйте декодировать остаток текста, а затем продолжайте читать. Нашли закодированное слово? Да, это слово *rad*. Здесь проявляется принципиальное отличие кода Хаффмана

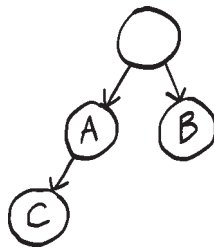
от ISO-8859-1. Длины кодов могут изменяться, поэтому декодирование должно выполняться по-другому.

Такая схема требует больших усилий по сравнению с разделением на блоки. С другой стороны, у нее есть одно большое преимущество. Обратите внимание: у букв, встречающихся чаще, коды более короткие. *D* встречается в тексте три раза, поэтому ее код состоит всего из двух цифр — в отличие от буквы *I*, встречающейся два раза, и буквы *P*, которая встречается всего один раз. Вместо того чтобы кодировать все 4 битами, мы применяем повышенное сжатие для часто используемых букв. В длинном тексте это обеспечит большой выигрыш!

Теперь, когда вы в целом понимаете, как работает алгоритм Хаффмана, посмотрим, какие свойства деревьев в нем используются. Прежде всего возможно ли наложение кодов? Для примера возьмем следующий код:

a = 0
b = 1
c = 00

Имеется двоичная последовательность 001. Что она представляет — *AA**B* или *CB*? Код для *c* и *a* содержит общую часть, поэтому ответ неясен. Дерево для этого кода будет выглядеть так:

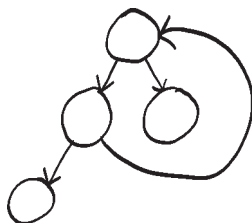


Буква *A* встречается на пути к *C*, и это создает проблему.

С кодами Хаффмана такой проблемы нет, потому что буквы помещаются только в листовых узлах. И от корня к каждому листовому узлу существует уникальный путь — это одно из свойств деревьев. Таким образом, можно гарантировать, что наложения не будет.

Это свойство также гарантирует, что каждой букве соответствует только один код. Существование нескольких путей к каждой букве означало бы, что за каждой буквой закреплено сразу несколько кодов, что излишне.

Когда вы читаете код по одной цифре, предполагается, что в конечном итоге вы получите букву. Если бы в графе присутствовал цикл, то такое предположение не работало бы из-за риска попасть в бесконечный цикл.



Но так как это дерево, циклы в нем заведомо отсутствуют, так что вы гарантированно придете к какой-нибудь букве.

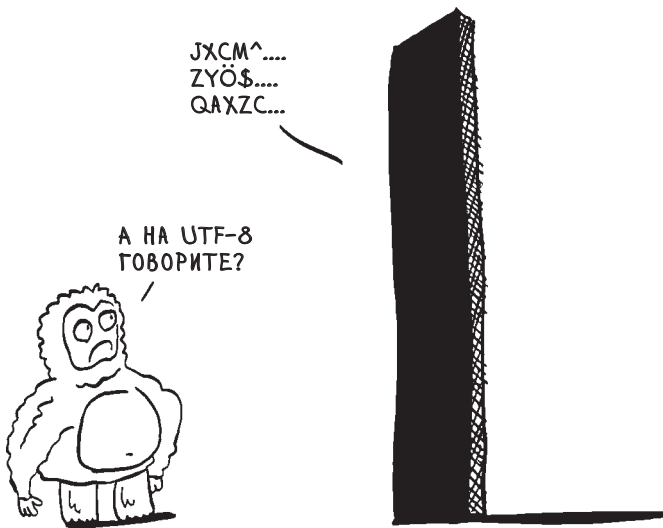
Мы используем корневое дерево. У корневых деревьев имеется корневой узел, и это важно, ведь вы должны знать, с чего начинать декодирование! С другой стороны, не каждый граф имеет корневой узел.

Наконец, дерево, используемое здесь, называется *бинарным*. Узлы в бинарных деревьях имеют не более двух дочерних узлов — левый и правый. Это логично, так как в двоичной записи всего две цифры. При наличии третьего дочернего узла было бы неясно, какую цифру он должен представлять.

Эта глава познакомила вас с деревьями. В следующей главе рассматриваются некоторые типы деревьев и их возможные применения.

Шпаргалка

- Деревья — это разновидность графов, но в деревьях не может быть циклов.
- Поиск в глубину — еще один алгоритм обхода графа. Он не подходит для поиска кратчайших путей.
- Бинарное дерево — особая разновидность деревьев, у которой каждый узел может иметь не более двух дочерних узлов.
- Существует много типов кодирования символов. Юникод является международным стандартом, а UTF-8 — самая популярная кодировка в Юникоде.



8

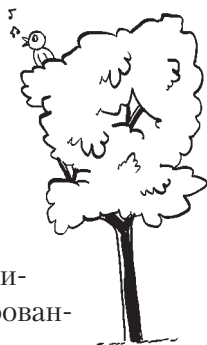
Сбалансированные деревья



В этой главе

- ✓ Вы узнаете о новой структуре данных, называемой бинарным деревом поиска (BST, Binary Search Tree).
- ✓ Вы узнаете, что такое сбалансированные деревья и почему они часто эффективнее массивов или связанных списков.
- ✓ Вы также узнаете об АВЛ-деревьях — разновидности сбалансированного дерева BST. В самом неблагоприятном сценарии бинарные деревья могут работать медленно. Сбалансированное дерево поможет им стать эффективнее.

В предыдущей главе вы узнали о новой структуре данных — дереве. Теперь, когда вы подружились с деревьями, пришло время узнать, для чего они используются. Если массивы и связанные списки не дают нужной производительности, целесообразно обратиться к структуре дерева. В этой главе рассматривается, какой производительности можно добиться с помощью деревьев. Мы изучим особую разновидность дерева, которая может обеспечить исключительную производительность, — так называемые сбалансированные деревья.

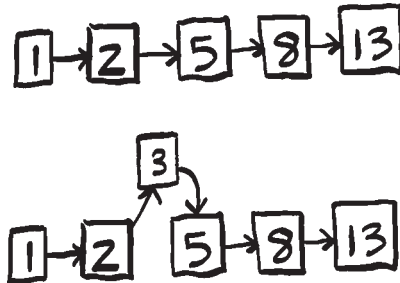


Балансировка

Помните бинарный поиск, о котором речь шла еще в главе 1? С его помощью можно найти информацию намного быстрее, чем простым поиском — со сложностью $O(\log n)$ вместо $O(n)$. Впрочем, есть одна проблема: вставка. Конечно, поиск занимает время $O(\log n)$, но массив должен быть отсортирован. Если требуется вставить новое число в отсортированный массив, это займет время $O(n)$. Проблема в том, чтобы найти место для нового значения. Придется передвинуть несколько элементов, чтобы освободить для него место.



Вот если бы вставку можно было выполнять как в связанном списке, где достаточно поменять пару указателей...



Но поиск в связанных списках выполняется за линейное время. Как использовать лучшие стороны обоих решений?

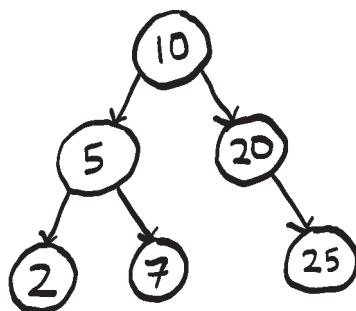
Деревья повышают скорость вставки

Итак, по сути, нам нужны скорость поиска в отсортированном массиве и более быстрая вставка. Как известно, вставка выполняется быстрее в связанных списках. А значит, нам понадобится структура данных, объединяющая эти идеи.

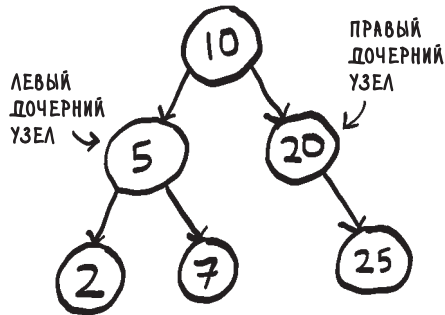
	ПОИСК	ВСТАВКА
ОТСОРТИРОВАННЫЙ МАССИВ	$O(\log n)$	$O(n)$
СВЯЗАННЫЙ СПИСОК	$O(n)$	$O(1)$
???	$O(\log n)$	БЫСТРЕЕ, ЧЕМ $O(n)$

И такая структура существует — это дерево! Деревья бывают десятков разных видов, поэтому я специально отмечаю, что нам нужно сбалансированное бинарное дерево поиска (BST). В этой главе я покажу, как работает BST, а затем вы научитесь его балансировать.

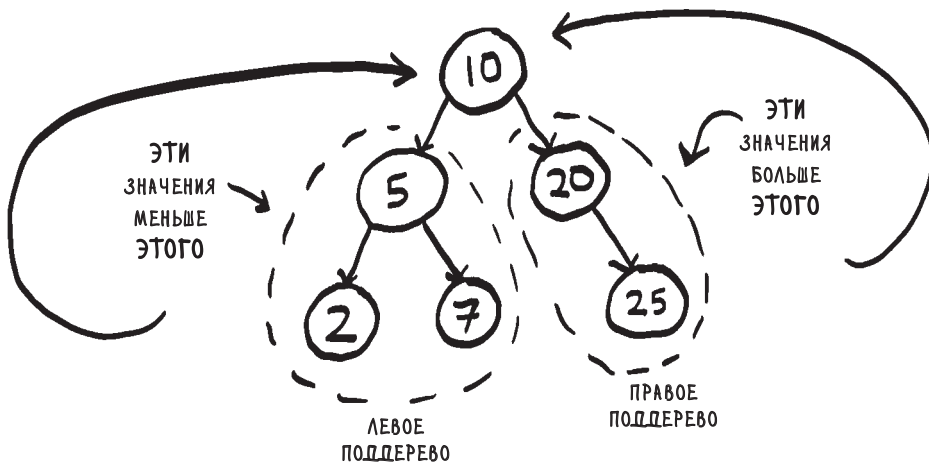
BST — это разновидность бинарного дерева. Пример BST:



Как и в бинарном дереве, каждый узел имеет до двух дочерних узлов: левый и правый. Но у этого дерева есть одно свойство, относящее его к BST: значение левого дочернего узла *всегда меньше*, чем значение узла, а значение правого дочернего узла *всегда больше*. Таким образом, для узла 10 левый дочерний узел имеет меньшее значение (5), а правый дочерний узел — большее значение (20).

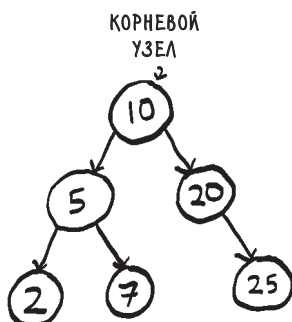


Более того, *все* числа в поддереве левого дочернего узла меньше самого узла!

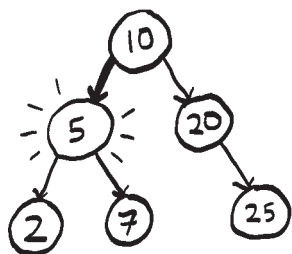


Это особое свойство означает, что поиск будет выполняться очень быстро.

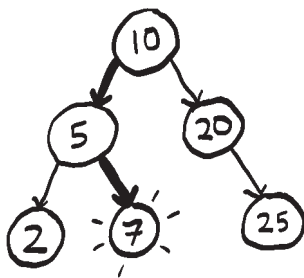
Давайте посмотрим, содержится ли число 7 в этом дереве. Начнем с корневого узла.



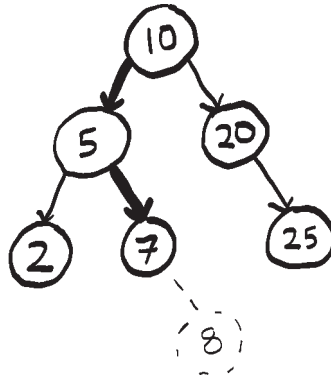
Число 7 меньше 10, поэтому проверяется левое поддерево. Вспомните: все узлы с меньшими значениями находятся в левом поддереве, а все узлы с большими значениями — в правом. Следовательно, мы сразу понимаем, что проверять узлы справа не нужно, потому что 7 там не будет. Опускаясь по левому поддереву от 10, мы переходим к узлу 5.



Число 7 больше 5, поэтому на этот раз идем направо.



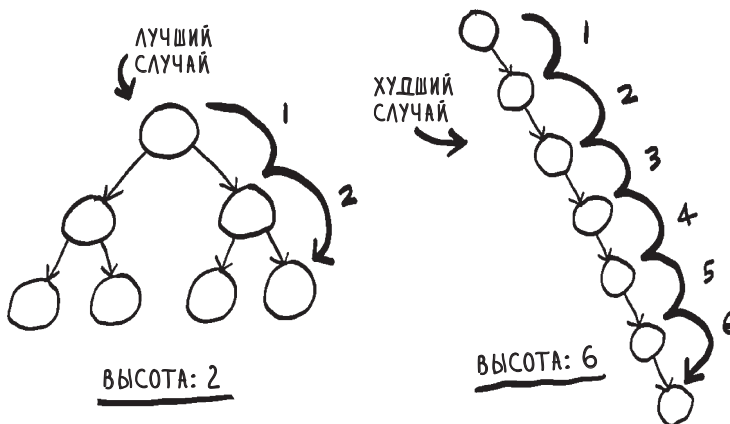
Нашли! Теперь поищем другое число — 8. Поиск проходит точно по такому же пути.



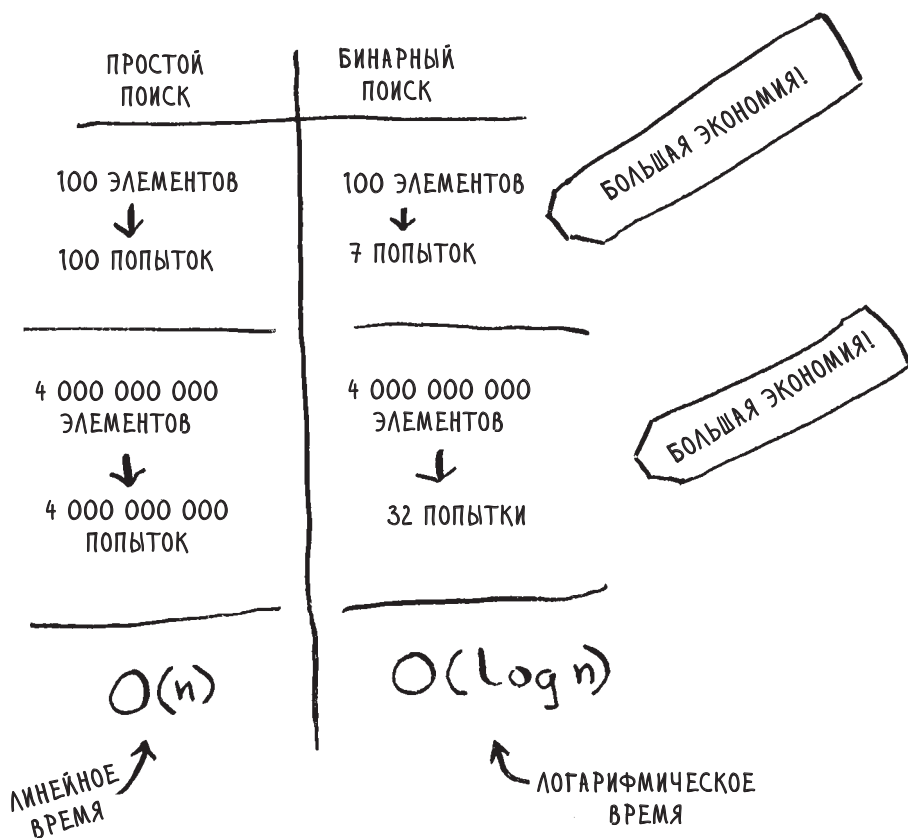
Только на этот раз его там нет! Если бы искомый узел присутствовал в дереве, он находился бы в месте, обозначенном пунктиром. Собственно, все наши рассуждения о деревьях обусловлены одной необходимостью: знать, работают ли они быстрее массивов и связанных списков. Итак, рассмотрим производительность деревьев, но для этого необходимо учитывать их высоту.

Короткие деревья работают быстрее

Рассмотрим два дерева. Оба дерева состоят из 7 узлов, но сильно различаются по производительности.



Высота дерева для лучшего случая равна 2. Это означает, что к любому узлу можно перейти от корневого узла максимум за 2 шага. Высота дерева для худшего случая равна 6. Это означает, что к любому узлу можно перейти от корневого узла максимум за 6 шагов. Сравним эти показатели с производительностью бинарного поиска в сравнении с простым поиском. На всякий случай напомним производительность бинарного и простого поиска.



Помните игру с отгадыванием чисел? Чтобы угадать число из набора 100 чисел, бинарный поиск потребует 7 попыток, а простой — 100. Нечто похожее происходит и с деревьями.

ДЕРЕВО ДЛЯ ХУДШЕГО СЛУЧАЯ	ДЕРЕВО ДЛЯ ЛУЧШЕГО СЛУЧАЯ
7 УЗЛОВ ↓ 6 ШАГОВ (ТО ЕСТЬ ВЫСОТА РАВНА 6)	7 УЗЛОВ ↓ 2 ШАГА (ТО ЕСТЬ ВЫСОТА РАВНА 2)

Дерево для худшего случая выше, и у него хуже производительность. В нем все узлы выстроены в одну линию. Дерево имеет высоту $O(n)$, так что поиск будет выполняться за время $O(n)$. Это можно представить так: дерево в действительности напоминает связанный список, так как один узел содержит ссылку на другой и т. д. А поиск по связанному списку выполняется за время $O(n)$.

Дерево для лучшего случая имеет высоту $O(\log n)$, а поиск по нему займет время $O(\log n)$.

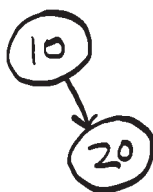
	ДЕРЕВО ДЛЯ ХУДШЕГО СЛУЧАЯ	ДЕРЕВО ДЛЯ ЛУЧШЕГО СЛУЧАЯ
	7 УЗЛОВ ↓ 6 ШАГОВ	7 УЗЛОВ ↓ 2 ШАГА
ВРЕМЯ ПОИСКА:	$O(n)$	$O(\log n)$

Таким образом, ситуация очень похожа на сравнение бинарного поиска с простым! Если можно обеспечить высоту дерева в $O(\log n)$, то поиск по дереву будет выполняться за время $O(\log n)$ — именно то, что требовалось.

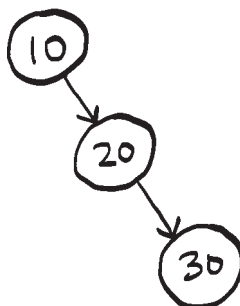
	ПОИСК	ВСТАВКА
ОТСОРТИРОВАННЫЙ МАССИВ	$O(\log n)$	$O(n)$
СВЯЗАННЫЙ СПИСОК	$O(n)$	$O(1)$
???	$O(\log n)$	БЫСТРЕЕ, ЧЕМ $O(n)$

Но как добиться, чтобы высота составляла $O(\log n)$? В следующем примере мы построим дерево, которое оказывается деревом для худшего случая (а этого нам желательно избежать). Начнем с одного узла, затем добавим другой.

10

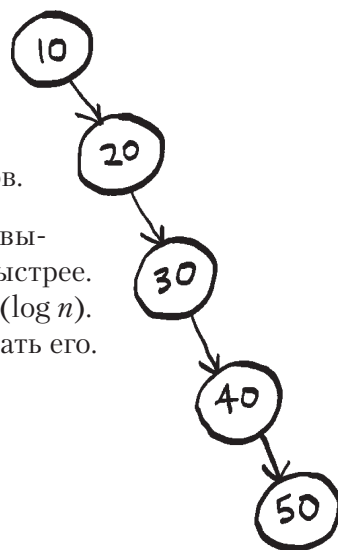


Пока неплохо. Добавим еще несколько узлов.



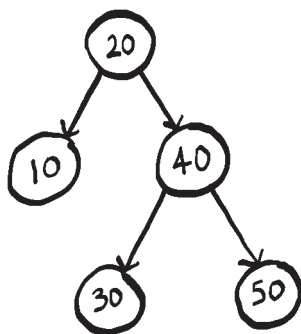
Узлы приходится добавлять справа, потому что все они больше всех своих предшественников.

Получается дерево для худшего случая — его высота равна $O(n)$! Короткие деревья работают быстрее. Кратчайшее время для BST может составлять $O(\log n)$. Чтобы укоротить BST, необходимо сбалансировать его. Рассмотрим сбалансированное дерево BST.



АВЛ-деревья: разновидность сбалансированных деревьев

АВЛ¹-деревья составляют разновидность самобалансируемых BST. Это означает, что АВЛ-деревья сохраняют высоту $O(\log n)$. Каждый раз, когда дерево разбалансируется, то есть его высота становится отличной от $O(\log n)$, оно корректирует себя. В последнем примере дерево может перебалансироваться и принять следующий вид:

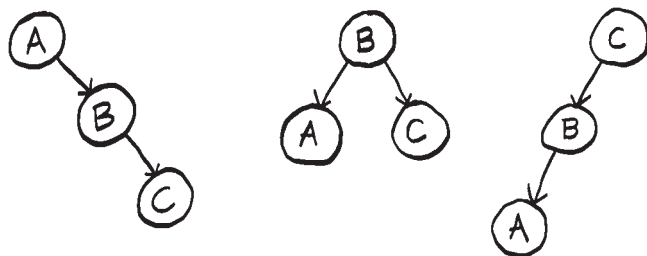


АВЛ-дерево обеспечивает нужную высоту $O(\log n)$, которая достигается самобалансировкой и поворотами.

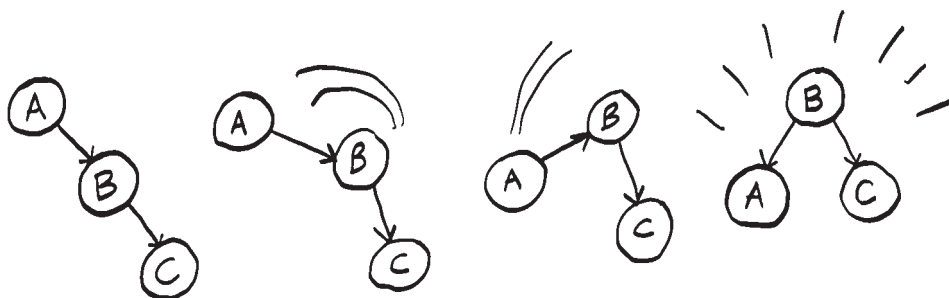
¹ Аббревиатура АВЛ образована от фамилий ученых, придумавших эту структуру: Адельсон-Вельский и Ландис. — *Примеч. ред.*

Повороты

Представьте, что имеется дерево с тремя узлами. Любой из них может быть корневым.



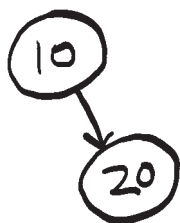
В результате поворота набор узлов смещается, образуя новую конфигурацию. Рассмотрим поворот в замедлении.



Мы выполняем поворот влево, начиная с несбалансированного дерева с корневым узлом A и заканчивая сбалансированным деревом с корневым узлом B.

Повороты — популярный метод балансировки деревьев. Снова начнем с одного узла.

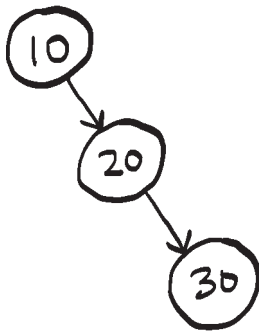
10



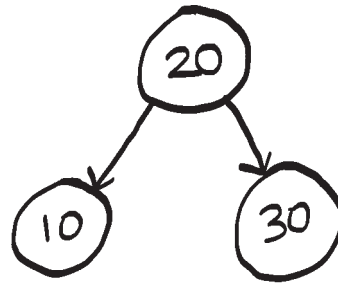
Добавим один узел.

Пока все нормально. Высоты дочерних узлов неодинаковы; разность между ними не превышает 1. Однако разность 1 приемлема для АВЛ-деревьев. Теперь добавим еще один узел.

Стоп! Дерево разбалансировалось. Пришло время поворотов!



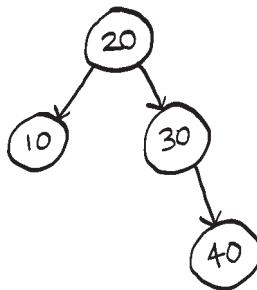
ПОВОРОТ!



СБАЛАНСИРОВАНО!

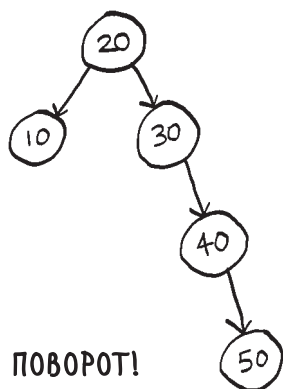
Мы выполнили поворот влево, теперь дерево снова сбалансировано.

Добавим еще один узел.

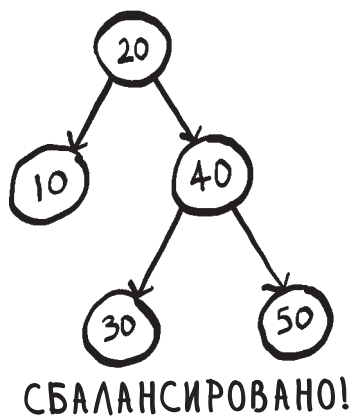


РАЗНОСТЬ 1
ДОПУСТИМА

И еще один.



Снова нужен поворот!

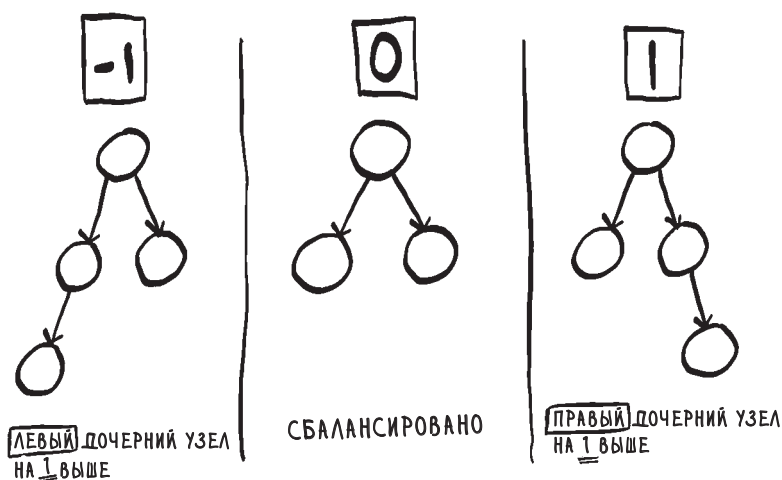


После поворота AVL-деревья сами перебалансируются. Напомним, что в последнем примере вместо узла 20 поворачивался узел 30. Следующий пример объясняет почему.

Как AVL-дерево узнает, что требуется поворот?

Мы видим, что дерево не сбалансировано — одна сторона длиннее другой. Но как дерево узнает об этом?

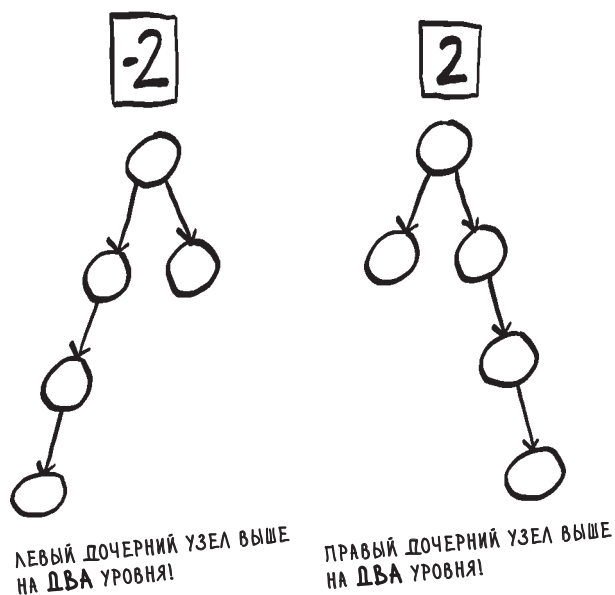
Чтобы дерево знало, когда требуется самобалансировка, оно должно хранить дополнительную информацию. В каждом узле хранится один или два вида информации: значение высоты или значение, которое иногда называют *коэффициентом балансировки*. Этот коэффициент должен быть равен -1 , 0 или 1 .



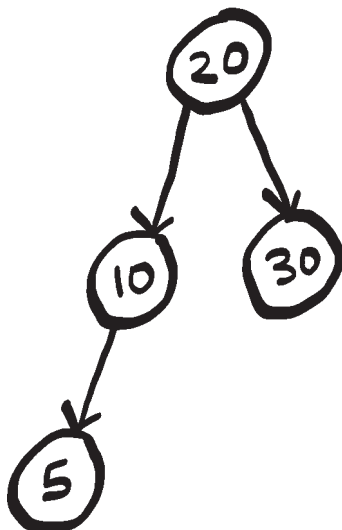
На рисунке приведены коэффициенты балансировки только для корневых узлов, но вам может понадобиться хранить коэффициент балансировки для каждого узла (пример такого рода будет приведен ниже).

Коэффициент балансировки сообщает, какой дочерний узел выше и насколько. По коэффициенту балансировки дерево может определить, когда следует проводить перебалансировку. Значение 0 означает, что дерево сбалансировано. Со значениями -1 или 1 тоже все нормально, потому что, напомним, AVL-деревья не обязаны быть идеально сбалансированы: разность 1 допустима.

Но если коэффициент балансировки падает ниже -1 или поднимается выше 1 , дерево нуждается в перебалансировке. Ниже изображены два дерева, нуждающихся в перебалансировке.



Как я уже говорил, в каждом узле должна храниться либо высота, либо коэффициент балансировки. В моем примере будет храниться и то и другое, чтобы вы видели, как они изменяются. Но если вам известны высоты каждого поддерева, вы легко вычислите коэффициент балансировки. Рассмотрим пример. Представьте следующее дерево.



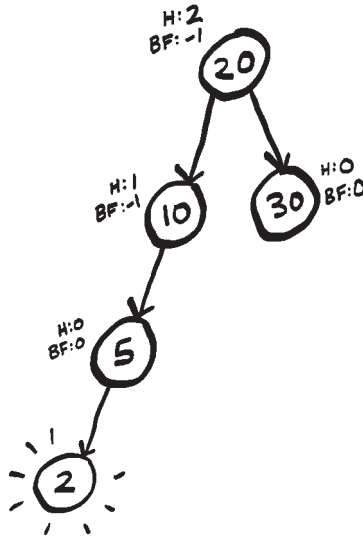
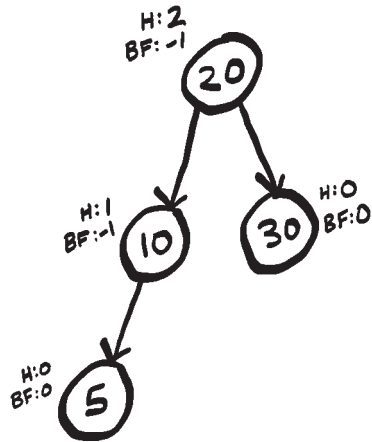
В него добавляется узел:

2

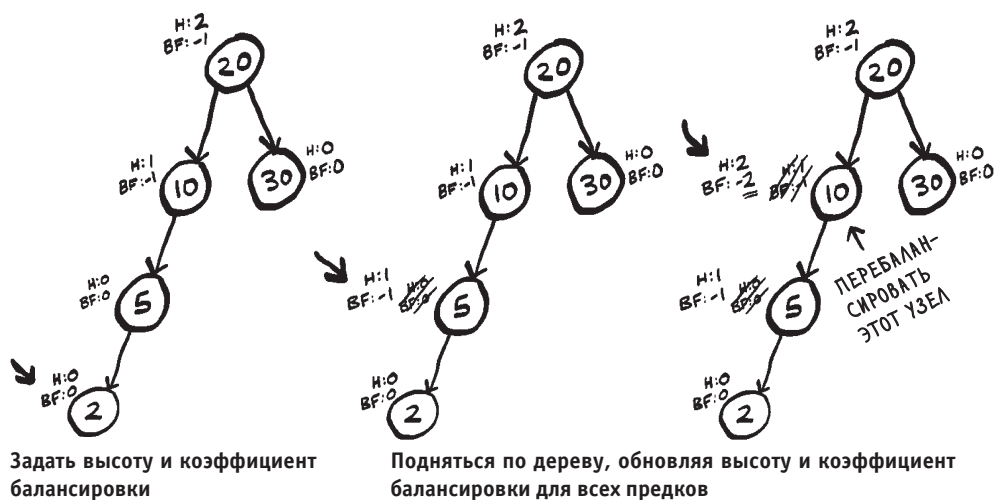
Для начала запишем высоту и коэффициент балансировки для каждого узла. На этой схеме H — высота, а BF — коэффициент балансировки.

Еще раз: я храню оба значения, чтобы показать, как они изменяются, но на самом деле достаточно хранить только одно. Убедитесь, что эти числа вам понятны. Заметим, что у всех корневых узлов коэффициент балансировки равен 0: у них нет дочерних узлов, поэтому и поддерживать баланс не нужно.

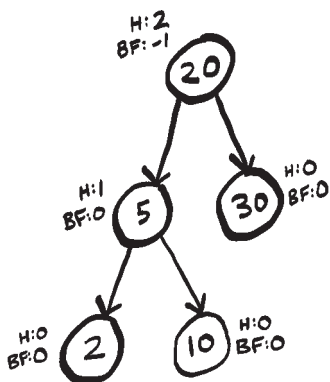
Добавим в дерево новый узел.



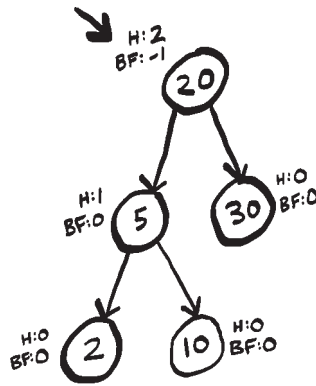
После того как узел будет добавлен, необходимо задать для него высоту и коэффициент балансировки. Затем следует подняться вверх по дереву, обновляя высоты и коэффициенты балансировки для всех его предков.



Ага! Мы только что присвоили коэффициент балансировки -2 , а это значит, что пришло время поворота! Далее мы вернемся к примеру и разберем его до конца, но главный вывод таков: после вставки вы обновляете коэффициенты балансировки для предков этого узла. AVL-дерево проверяет коэффициент балансировки, чтобы узнать, когда выполнять перебалансировку. В завершение повернем узел 10.



Теперь поддерево сбалансировано. Продолжим движение вверх по дереву.



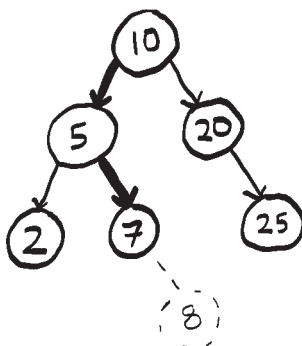
Ничего обновлять не нужно. На самом деле двигаться вверх по дереву необязательно, потому что AVL-деревья требуют не более одной перебалансировки.

AVL-деревья хороши, если требуется сбалансированное дерево BST. Вспомним главное, что мы узнали:

- Бинарные деревья — это одна из разновидностей деревьев.
- В бинарных деревьях каждый узел имеет не более двух дочерних узлов.
- BST — это разновидность бинарного дерева, в которой все значения в левом поддереве меньше значения узла, а все значения в правом поддереве больше значения узла.
- BST может дать превосходную производительность, если получится обеспечить их высоту равной $O(\log n)$.
- AVL-деревья являются деревьями BST; это гарантирует, что их высота будет равна $O(\log n)$.
- AVL-деревья балансируют себя с помощью поворотов.

Впрочем, это еще не все. Мы рассмотрели один сценарий поворотов, но существуют и другие. Мы не будем останавливаться на них подробно, так как вам вряд ли придется реализовывать AVL-дерево самостоятельно.

Теперь вы знаете, что AVL-деревья обеспечивают производительность поиска $O(\log n)$. А что же вставка? Вставка сводится к поиску места для вставки узла и добавления указателя, как в связанном списке. Например, если вы захотели вставить 8 в это дерево, необходимо определить, где добавить указатель.



Таким образом, вставки тоже выполняются за время $O(\log n)$.

В начале этой главы мы искали структуру данных, обеспечивающую как быстрый поиск, так и быструю вставку. И такая волшебная структура данных есть: это сбалансированное дерево BST!

	ПОИСК	ВСТАВКА
ОТСОРТИРОВАННЫЙ МАССИВ	$O(\log n)$	$O(n)$
СВЯЗАННЫЙ СПИСОК	$O(n)$	$O(1)$
BST	$O(\log n)$	$O(\log n)$

Косые деревья

АВЛ-деревья — хорошие базовые сбалансированные деревья BST, обеспечивающие время $O(\log n)$ для некоторых операций.

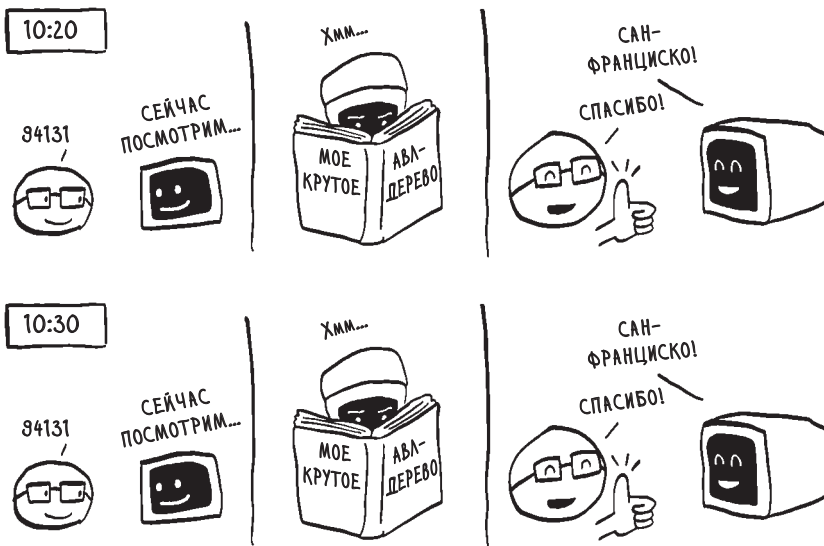


Косые деревья (или Splay-деревья) представляют другой подход к сбалансированным деревьям BST. Их самое замечательное свойство в том, что если вы недавно искали какой-то элемент, то следующий его поиск будет быстрее. Это свойство интуитивно воспринимается положительно. Например, представьте, что вы работаете с программой, которая получает почтовый индекс и ищет по нему город.

Взаимодействие можно представить следующим образом:



Теперь представьте, что вы многократно проводите поиск одного индекса.





Выглядит немного глупо.

Программа только что проводила поиск этого индекса; почему бы ей не запомнить результат? На самом деле все должно быть так:

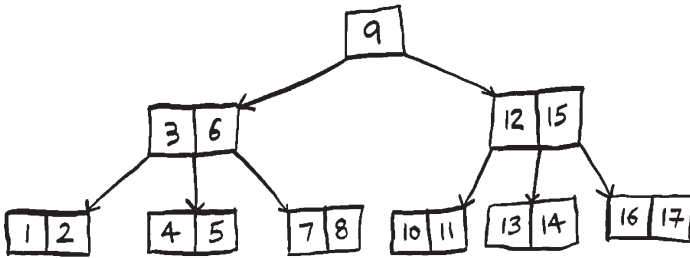


Косые деревья позволяют это сделать. Когда вы ищете узел в косом дереве, он становится новым корнем, так что при повторном поиске он будет найден сразу же. В общем случае узлы, которые вы недавно искали, группируются в начале и находятся быстрее. С другой стороны, дерево заведомо не будет сбалансировано, а значит, *некоторые* операции поиска будут занимать время, большее $O(\log n)$, и даже достигать линейного времени! Кроме того, при выполнении поиска вам, возможно, придется повернуть узел до корневой позиции, если он еще некорневой, а это тоже требует времени.

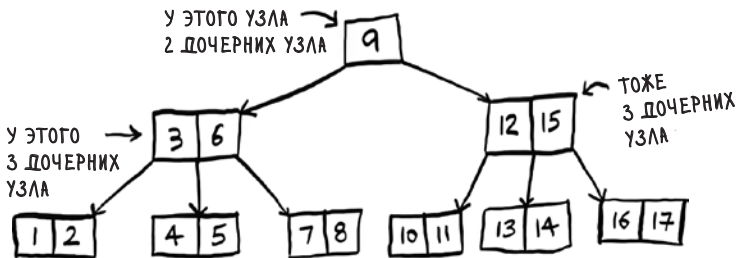
Но нас устраивает, что дерево не будет сбалансировано постоянно. Главное, что при проведении n операций поиска общее время $O(n \log n)$ *гарантировано* — то есть $O(\log n)$ на один поиск. Таким образом, хотя один поиск может занять время, превышающее $O(\log n)$, в среднем все операции поиска сходятся ко времени $O(\log n)$, а наша цель как раз и есть ускорение поиска.

В-деревья

В-деревья представляют собой обобщенную форму бинарных деревьев. Они часто используются для построения баз данных. Ниже представлен пример В-дерева.

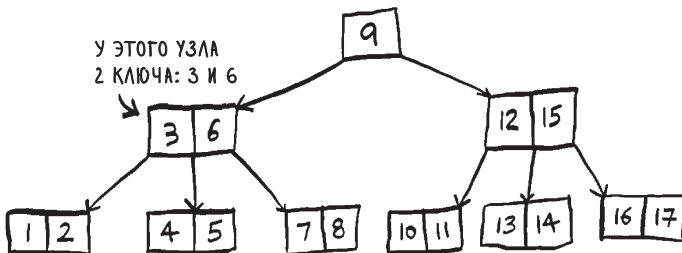


Выглядит необычно, правда? Можно заметить, что некоторые узлы имеют более двух дочерних узлов.



В отличие от бинарных деревьев, в В-деревьях узлы могут иметь много дочерних узлов.

Вероятно, вы уже заметили, что, в отличие от предыдущих деревьев, в В-дереве многие узлы имеют два ключа.



Итак, узлы в В-деревьях не только могут иметь более двух дочерних узлов — они также могут иметь несколько ключей! Именно поэтому я сказал, что В-деревья — это обобщенная форма BST.

Какие преимущества есть у В-деревьев?

Для В-деревьев существует оптимизация, интересная тем, что применяется на физическом уровне. При поиске по дереву получение данных требует перемещения механических компонентов оборудования (например, читающей головки на HDD). Время получения данных называется *временем поиска*. Оно может быть важным фактором, определяющим, насколько быстро или медленно работает алгоритм.

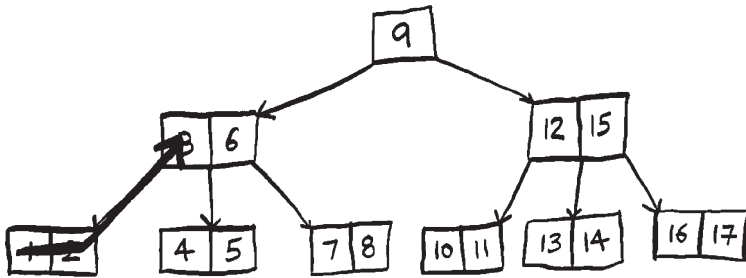
Ситуацию можно сравнить с посещением магазина. Вы можете покупать продукты по одному. Представьте, что вы решаете купить молоко. Вернувшись домой, вы понимаете, что было бы неплохо купить хлеба, и возвращаетесь в магазин. По возвращении вы видите, что у вас закончился кофе, и снова идете в магазин... Как же это неэффективно! Намного лучше зайти в магазин один раз и купить все необходимое, находясь в нем. В этом примере время похода и возвращения в магазин представляет собой время поиска.

Фундаментальная концепция В-деревьев заключается в том, что *после выполнения поиска можно прочитать дополнительные данные в память*. Другими словами, находясь в магазине, вы можете купить все необходимое, вместо того чтобы возвращаться в него снова и снова.

В В-деревьях используются большие узлы; каждый узел может иметь больше ключей и дочерних узлов, чем бинарное дерево. Таким образом, чтение каждого узла занимает больше времени. *С другой стороны, поиск ускоряется, потому что за один раз читается больший объем данных*. Именно это обстоятельство обеспечивает высокую скорость работы В-деревьев.

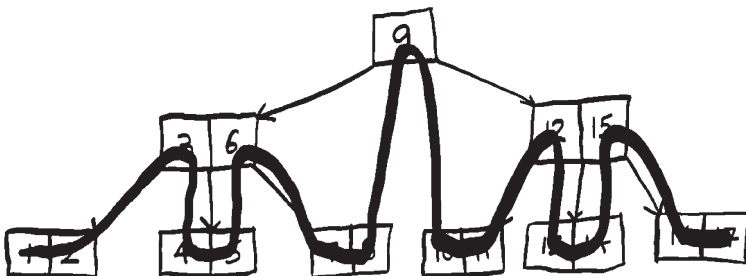
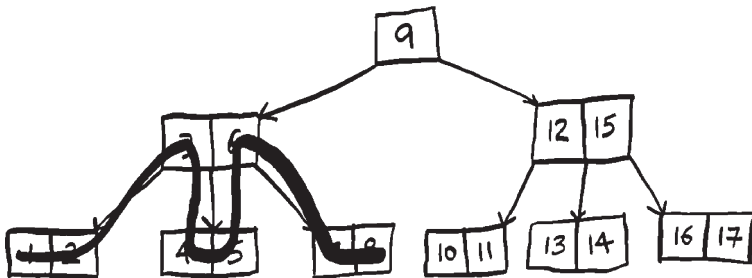
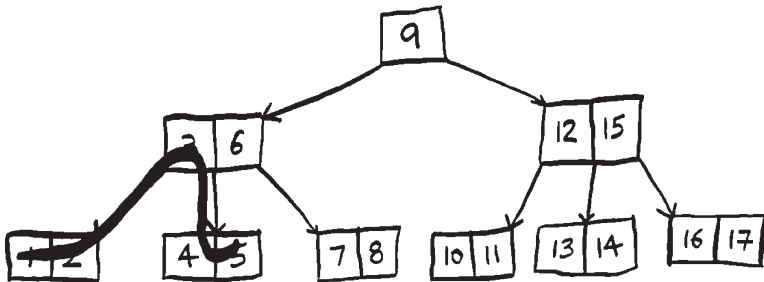
Структура В-деревьев популярна в реализациях баз данных. И это не удивительно, так как в базах данных много времени занимает чтение данных с диска.

Обратите внимание на то, как упорядочены узлы В-дерева; это тоже довольно интересно. Вы начинаете с левого нижнего узла.



Куда вы пойдете дальше?

А дальше вы змейкой обходите все дерево.



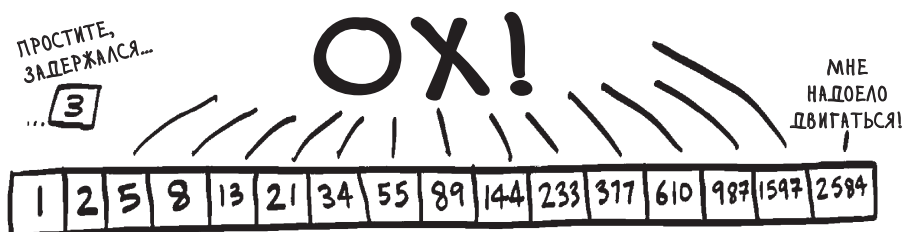
Обратите внимание: здесь все еще действует свойство деревьев BST, когда для каждого ключа значения ключей в левом поддереве меньше, а значения ключей в правом поддереве — больше него. Например, для ключа 3 левое поддерево состоит из ключей 1 и 2, а правое поддерево — из ключей 4 и 5.

Также обратите внимание, что количество дочерних узлов на 1 больше количества ключей. Таким образом, корневой узел имеет один ключ и два дочерних узла. Каждый из дочерних узлов имеет два ключа и три дочерних узла.

На этом завершаются две главы, посвященные деревьям. Вам вряд ли придется реализовывать дерево самостоятельно, но важно знать, что деревья — это разновидность графов и они обладают превосходной производительностью. В следующей главе мы вернемся к графам и рассмотрим их новую разновидность: *взвешенные графы*.

Шпаргалка

- Сбалансированные бинарные деревья поиска (BST) обеспечивают такую же производительность в нотации «О-большое», как массивы, и лучшую производительность вставки.
- Высота дерева влияет на его производительность.
- AVL-деревья — популярная разновидность сбалансированных деревьев BST. Как и большинство сбалансированных деревьев, AVL-деревья балансируются поворотом.
- В-деревья представляют собой обобщенные деревья BST, у которых каждый узел может иметь несколько ключей и несколько дочерних узлов.
- Время поиска можно сравнить со временем похода в магазин. В-деревья минимизируют время поиска за счет чтения большего объема данных за одну операцию.



9

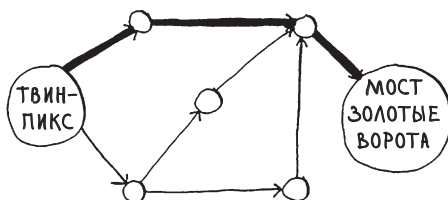
Алгоритм Дейкстры



В этой главе

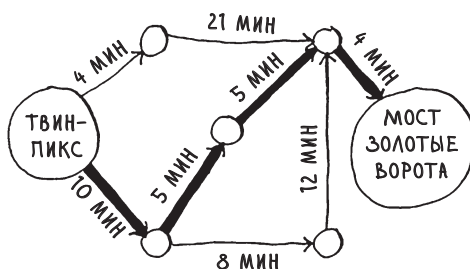
- ✓ Мы продолжим изучение графов и познакомимся со взвешенными графами, где некоторым ребрам назначаются большие или меньшие веса.
- ✓ Вы изучите алгоритм Дейкстры, который позволяет получить ответ на вопрос «Как выглядит кратчайший путь к X?» для взвешенных графов.
- ✓ Вы узнаете о циклах в графах, для которых алгоритм Дейкстры не работает.

В главе 6 вы узнали, как найти путь из точки А в точку В.



Найденный путь не обязательно окажется самым быстрым. Этот путь считается кратчайшим, потому что состоит из наименьшего количества

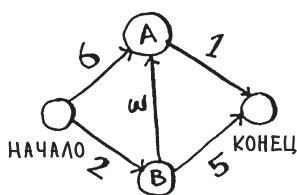
сегментов (трех). Но предположим, что с каждым сегментом связывается продолжительность перемещения. И тогда выясняется, что существует и более быстрый путь.



В предыдущей главе рассматривался поиск в ширину. Этот алгоритм находит путь с минимальным количеством сегментов (граф на предыдущей странице). А если вы захотите найти самый быстрый путь (второй граф)? Быстрее всего это делается при помощи другого алгоритма, который называется *алгоритмом Дейкстры*.

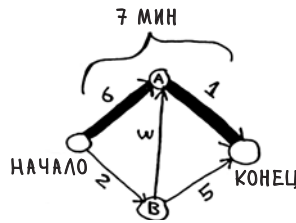
Работа с алгоритмом Дейкстры

Посмотрим, как этот алгоритм работает с графом.



Каждому ребру назначается время перемещения в минутах. Алгоритм Дейкстры используется для поиска пути от начальной точки к конечной за кратчайшее возможное время.

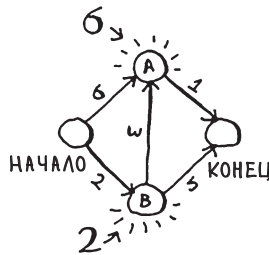
Применив к этому графу поиск в ширину, вы получите следующий кратчайший путь.



Этот путь занимает 7 минут. А может, существует путь, который займет меньше времени? Алгоритм Дейкстры состоит из четырех шагов.

1. Найти узел с наименьшей стоимостью (то есть узел, до которого можно добраться за минимальное время).
2. Обновить стоимости соседей этого узла (вскоре я объясню, что имеется в виду).
3. Повторять, пока это не будет сделано для всех узлов графа.
4. Вычислить итоговый путь.

Шаг 1: найти узел с наименьшей стоимостью. Вы стоите в самом начале и думаете, куда направиться: к узлу А или к узлу В. Сколько времени понадобится, чтобы добраться до каждого из этих узлов?

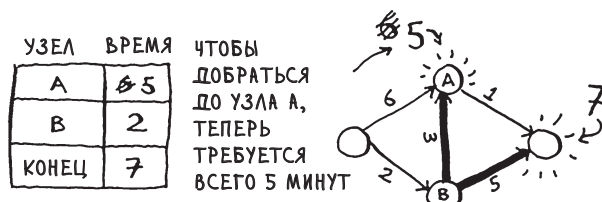


До узла А вы будете добираться 6 минут, а до узла В — 2 минуты. Что касается остальных узлов, мы о них пока ничего не знаем.

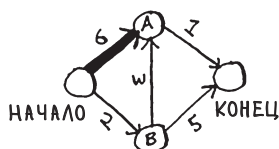
Так как время достижения конечного узла остается неизвестным, мы считаем, что оно бесконечно (вскоре вы увидите почему.) Узел В — ближайший... он находится всего в 2 минутах.

УЗЕЛ	ВРЕМЯ ПЕРЕХОДА К УЗЛУ
А	6
В	2
КОНЕЦ	∞

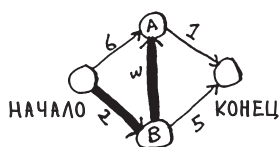
Шаг 2: вычислить, сколько времени потребуется для того, чтобы добраться до всех внешних соседей В при переходе по ребру из В.



Ого, да мы обнаружили более короткий путь к узлу A! Раньше для перехода к нему требовалось 6 минут.



А если идти через узел В, то существует путь, который занимает всего 5 минут!



Если вы нашли более короткий путь для соседа В, обновите его стоимость. В данном случае мы нашли:

- более короткий путь к А (сокращение с 6 минут до 5);
- более короткий путь к конечному узлу (сокращение от бесконечности до 7 минут).

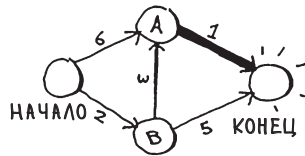
Шаг 3: повторяем!

Снова шаг 1: находим узел, для перехода к которому требуется наименьшее время. С узлом В работа закончена, поэтому наименьшую оценку времени имеет узел А.

УЗЕЛ ВРЕМЯ	
А	5
В	2
КОНЕЦ	7



Снова шаг 2: обновляем стоимости внешних соседей А.



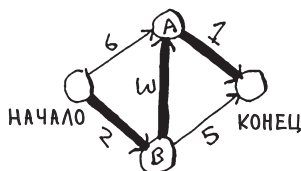
Путь до конечного узла теперь занимает всего 6 минут!

Алгоритм Дейкстры выполнен для каждого узла (выполнять его для конечного узла не нужно). К этому моменту вам известно следующее:

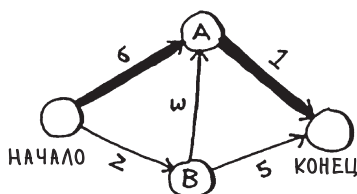
- Чтобы добраться до узла В, нужно 2 минуты.
- Чтобы добраться до узла А, нужно 5 минут.
- Чтобы добраться до конечного узла, нужно 6 минут.

УЗЕЛ ВРЕМЯ	
А	5
В	2
КОНЕЦ	6

Последний шаг — вычисление итогового пути — откладывается до следующего раздела. А пока я просто покажу, как выглядит итоговый путь.

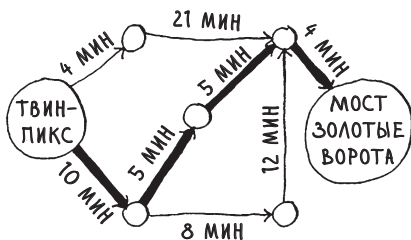


Алгоритм поиска в ширину не найдет этот путь как кратчайший, потому что он состоит из трех сегментов, а от начального узла до конечного можно добраться всего за два сегмента.

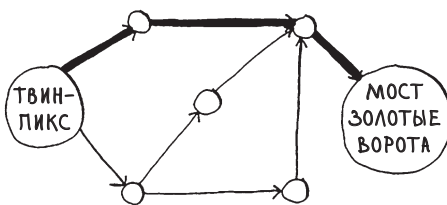


КРАТЧАЙШИЙ ПУТЬ
С ПОИСКОМ В ШИРИНУ

В предыдущей главе мы использовали поиск в ширину для нахождения кратчайшего пути между двумя точками. Тогда под «кратчайшим путем» понимался путь с минимальным количеством сегментов. С другой стороны, в алгоритме Дейкстры каждому сегменту присваивается число (вес), а алгоритм Дейкстры находит путь с наименьшим суммарным весом.



ВЗВЕШЕННЫЙ ГРАФ
(АЛГОРИТМ ДЕЙКСТРЫ)



НЕВЗВЕШЕННЫЙ ГРАФ
(ПОИСК В ШИРИНУ)

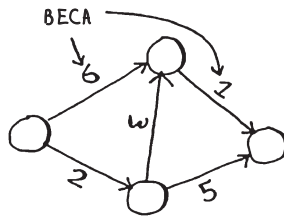
На всякий случай повторим: алгоритм Дейкстры состоит из четырех шагов.

1. Найти узел с наименьшей стоимостью (то есть узел, до которого можно добраться за минимальное время).
2. Проверить, существует ли более дешевый путь к соседям этого узла, и если существует, обновить их стоимости.
3. Повторять, пока это не будет сделано для всех узлов графа.
4. Вычислить итоговый путь (об этом в следующем разделе!).

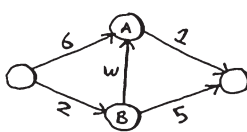
Терминология

Я хочу привести еще несколько примеров применения алгоритма Дейкстры. Но сначала стоит немного разобраться с терминологией.

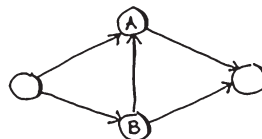
Когда вы работаете с алгоритмом Дейкстры, с каждым ребром графа связывается число, называемое *весом*.



Граф с весами называется *взвешенным графом*. Граф без весов называется *невзвешенным графом*.



ВЗВЕШЕННЫЙ ГРАФ

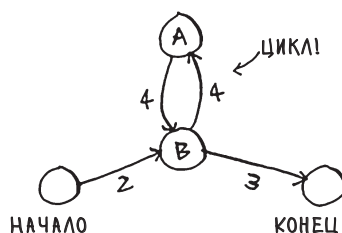


НЕВЗВЕШЕННЫЙ ГРАФ

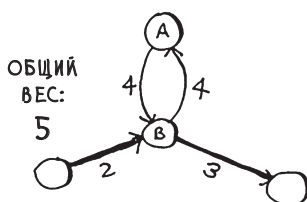
Для вычисления кратчайшего пути в невзвешенном графе используется *поиск в ширину*. Кратчайшие пути во взвешенном графе вычисляются по алгоритму Дейкстры. В графах также могут присутствовать *циклы*:



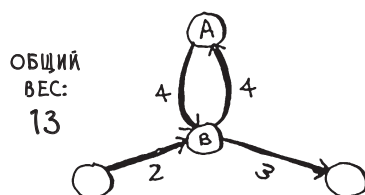
Это означает, что вы можете начать с некоторого узла, перемещаться по графу, а потом снова оказаться в том же узле. Предположим, вы ищете кратчайший путь в графе, содержащем цикл.



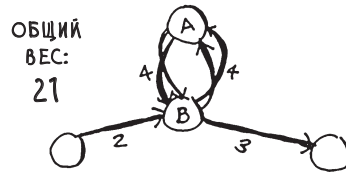
Есть ли смысл в перемещении по циклу? Что ж, вы можете использовать путь без прохождения цикла:



А можете пройти по циклу:



Вы в любом случае оказываетесь в узле А, но цикл добавляет лишний вес. Вы даже можете обойти цикл дважды, если вдруг захотите.

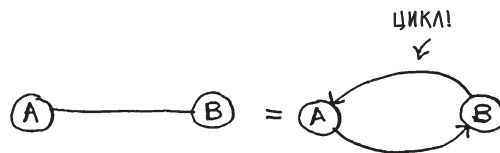


Но каждый раз, когда вы проходите по циклу, вы только увеличиваете суммарный вес на 8. Следовательно, путь с обходом цикла никогда не будет кратчайшим.

Наконец, вы еще не забыли наше обсуждение направленных и ненаправленных графов из главы 6?



Само понятие ненаправленного графа означает, что каждый из двух узлов фактически ведет к другому узлу. А это цикл!



В ненаправленном графе каждое новое ребро добавляет еще один цикл. Алгоритм Дейкстры работает только с графами, в которых нет циклов, где все ребра неотрицательны. Да, ребра графа могут иметь отрицательный вес! Но алгоритм Дейкстры не сработает — в этом случае понадобится алгоритм Белмана — Форда. О ребрах с отрицательным весом речь пойдет далее в этой главе.

История одного обмена

Но довольно терминологии, пора рассмотреть конкретный пример! Это Рама. Он хочет выменять свою книгу по музыке на пианино.

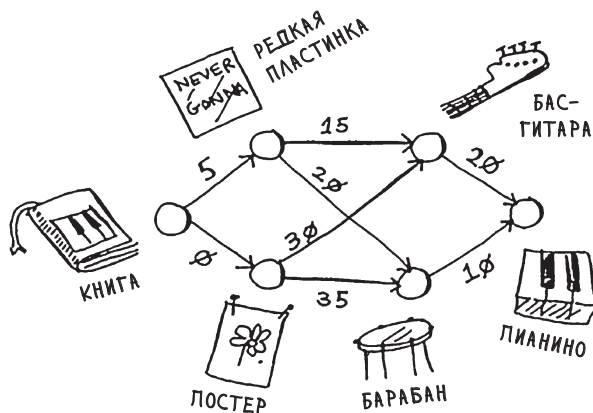
«Я тебе дам за книгу вот этот постер, — говорит Алекс. — Это моя любимая группа Destroyer. Или могу дать за книгу редкую пластинку Рика Этли и еще пять долларов». — «О, я слышала, что на этой пластинке есть отличные песни, — говорит Эми. — Готова отдать за постер или пластинку свою гитару или ударную установку».



«Всю жизнь мечтал играть на гитаре, — восклицает Бетховен. — Слушай, я отдам тебе свое пианино за любую из вещей Эми».

Прекрасно! Рама с небольшими дополнительными тратами может поменять свою книгу на настоящее пианино. Теперь остается понять, как ему потратить наименьшую сумму на цепочке обменов.

Изобразим полученные им предложения в виде графа:



Узлы графа — это предметы, на которые может поменяться Рама. Веса ребер представляют сумму доплаты за обмен. Таким образом, Рама может

поменять постер на гитару за \$30 или же пластинку на гитару за \$15. Как Раме вычислить путь от книги до пианино, при котором он потратит наименьшую сумму? На помощь приходит алгоритм Дейкстры! Вспомните, что алгоритм Дейкстры состоит из четырех шагов. В этом примере мы выполним все четыре шага, а в конце будет вычислен итоговый путь.

УЗЕЛ	СТОИМОСТЬ
ПЛАСТИНКА	5
ПОСТЕР	0
ГИТАРА	∞
БАРАБАН	∞
ПИАНИНО	∞

МЫ ЕЩЕ НЕ ДОХОДИЛИ ДО ЭТИХ УЗЛОВ ОТ НАЧАЛЬНОГО

Прежде чем начинать, необходимо немного подготовиться. Постройте таблицу со стоимостями всех узлов. (Стоимость узла определяет затраты на его достижение.)

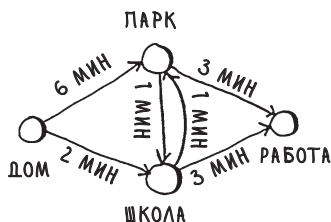
Таблица будет обновляться по мере работы алгоритма. Для вычисления итогового пути в таблицу также необходимо добавить столбец «родитель».

УЗЕЛ	РОДИТЕЛЬ
ПЛАСТИНКА	КНИГА
ПОСТЕР	КНИГА
ГИТАРА	—
БАРАБАН	—
ПИАНИНО	—

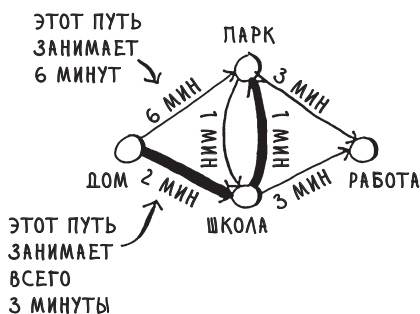
Вскоре я покажу, как работает этот столбец. А пока просто запустим алгоритм.

Шаг 1: найти узел с наименьшей стоимостью. В данном случае самый дешевый вариант обмена с доплатой \$0 — это постер. Возможно ли получить постер с меньшими затратами? Это очень важный момент, хорошенько по-

думайте над ним. Удастся ли вам найти серию обменов, при которой Рама получит постер менее чем за \$0? Продолжайте читать, когда будете готовы ответить на вопрос. Правильный ответ: нет, не удастся. *Так как постер является узлом с наименьшей стоимостью, до которого может добраться Рама, снизить его стоимость невозможно.* На происходящее можно взглянуть иначе: предположим, вы едете из дома на работу.



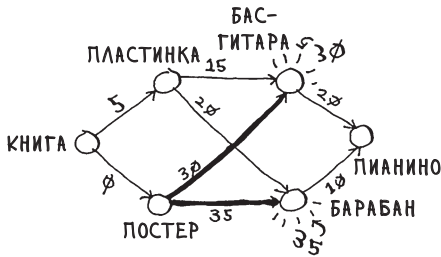
Если вы выберете путь к школе, это займет 2 минуты. Если вы выберете путь к парку, это займет 6 минут. Существует ли путь, при котором вы выбираете путь к парку и оказываетесь в школе менее чем за 2 минуты? Это невозможно, потому что только для того, чтобы попасть в парк, потребуется более 2 минут. Можно ли найти более быстрый путь в парк? Да, можно.



В этом заключается ключевая идея алгоритма Дейкстры: в графе ищется путь с наименьшей стоимостью. Пути к этому узлу с меньшими затратами не существует!

Возвращаемся к музыкальному примеру. Вариант с постером обладает наименьшей стоимостью.

Шаг 2: вычислить, сколько времени потребуется для того, чтобы добраться до всех его соседей (стоимость).



РОДИТЕЛЬ	УЗЕЛ	СТОИ-МОСТЬ
КНИГА		5
КНИГА	ПОСТЕР	0
ПОСТЕР		30
ПОСТЕР	БАРАБАН	35
—	ПИАНИНО	∞

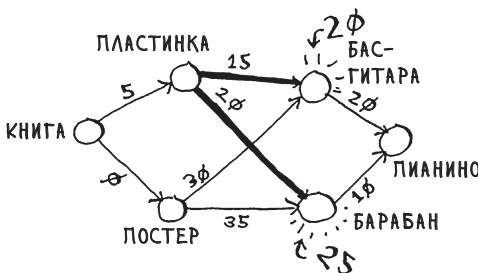
Стоимости бас-гитары и барабана заносятся в таблицу. Они были заданы при переходе через узел постера, поэтому постер указывается как их родитель. А это означает, что для того, чтобы добраться до бас-гитары, вы проходите по ребру от постера; то же самое происходит с барабаном.

К ЭТИМ УЗЛАМ ПЕРЕХОДИМ ОТ УЗЛА «ПОСТЕР»

РОДИТЕЛЬ	УЗЕЛ	СТОИ-МОСТЬ
КНИГА	ПЛАСТИНКА	5
КНИГА	ПОСТЕР	0
ПОСТЕР	ГИТАРА	30
ПОСТЕР	БАРАБАН	35
—	ПИАНИНО	∞

Снова шаг 1: пластинка — следующий по стоимости узел (\$5).

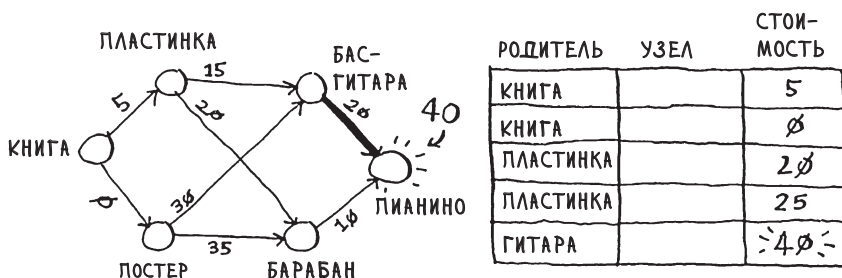
Снова шаг 2: обновляются значения всех его соседей.



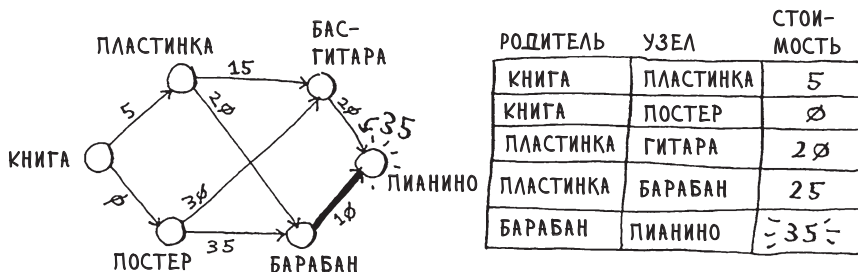
РОДИТЕЛЬ	УЗЕЛ	СТОИ-МОСТЬ
КНИГА	ПЛАСТИНКА	5
КНИГА	ПОСТЕР	0
ПЛАСТИНКА	ГИТАРА	30 20
ПЛАСТИНКА	БАРАБАН	35 25
—	ПИАНИНО	∞

Смотрите, стоимости барабана и гитары обновились! Это означает, что к барабану и гитаре дешевле перейти через ребро, идущее от пластинки. Соответственно, пластинка назначается новым родителем обоих инструментов.

Следующий по стоимости узел — бас-гитара. Обновите данные его соседей.



Хорошо, мы наконец-то вычислили стоимость для пианино при условии обмена гитары на пианино. Соответственно, гитара назначается родителем. Наконец, задается стоимость последнего узла — барабана.



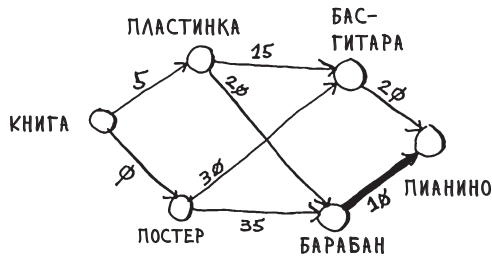
Оказывается, Рама может получить пианино еще дешевле, поменяв ударную установку на пианино. Таким образом, *самая дешевая цепочка обменов обойдется Рама в \$35.*

Теперь, как я и обещал, необходимо вычислить итоговый путь. К этому моменту вы уже знаете, что кратчайший путь обойдется в \$35, но как этот путь определить? Для начала возьмем родителя узла «пианино».

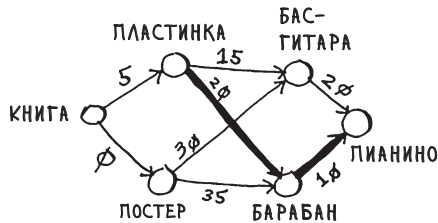
→

РОДИТЕЛЬ	УЗЕЛ
КНИГА	ПЛАСТИНКА
КНИГА	ПОСТЕР
ПЛАСТИНКА	ГИТАРА
ПЛАСТИНКА	БАРАБАН
БАРАБАН	ПИАНИНО

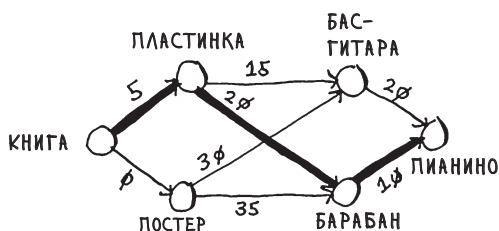
В качестве родителя узла «пианино» указан узел «барабан».



А в качестве родителя узла «барабан» указан узел «пластинка».



Следовательно, Рама обменивает пластинку на барабан. И конечно, в самом начале он меняет книгу на пластинку. Проходя по родительским узлам в обратном направлении, мы получаем полный путь.



Серия обменов, которую должен сделать Рама, выглядит так:

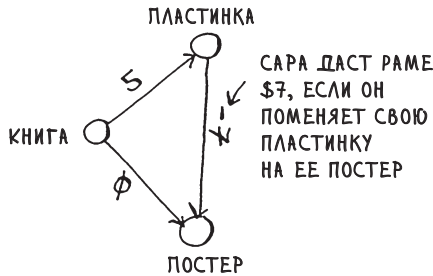
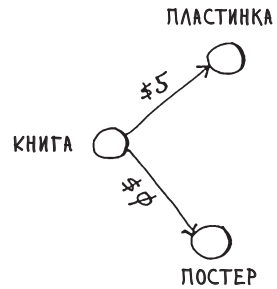


До сих пор я использовал термин «кратчайший путь» более или менее буквально, понимая под ним вычисление кратчайшего пути между двумя точками или двумя людьми. Надеюсь, этот пример показал, что кратчайший путь далеко не всегда связывается с физическим расстоянием: он может быть направлен на минимизацию какой-либо характеристики. В нашем примере Рама хотел свести к минимуму свои затраты при обмене. Спасибо Дейкстре!

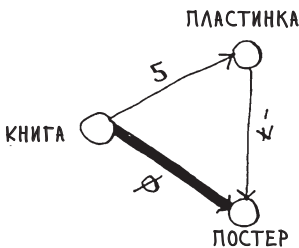
Ребра с отрицательным весом

В предыдущем примере Алекс предложил в обмен на книгу один из двух предметов.

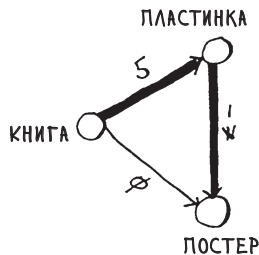
Предположим, Сара предложила обменять пластинку на постер и при этом она еще и *даст* Рама \$7. Рама ничего не тратит при этом обмене, вместо этого он получит \$7. Как изобразить это предложение на графе?



Ребро, ведущее от пластинки к постеру, имеет отрицательный вес! Если Рама пойдет на этот обмен, он получит \$7. Теперь к постеру можно добраться двумя способами.



ЕСЛИ РАМА ИДЕТ
ПО ЭТОМУ ПУТИ,
ОН ПОЛУЧАЕТ \$0



ЕСЛИ РАМА ИДЕТ
ПО ЭТОМУ ПУТИ,
ОН ПОЛУЧАЕТ \$2

А значит, во втором обмене появляется смысл — Рама получает \$2!

Теперь, если вы помните, Рама может обменять постер на барабан. И здесь возможны два пути.



Второй путь обойдется на \$2 дешевле, поэтому нужно выбрать этот путь, верно?

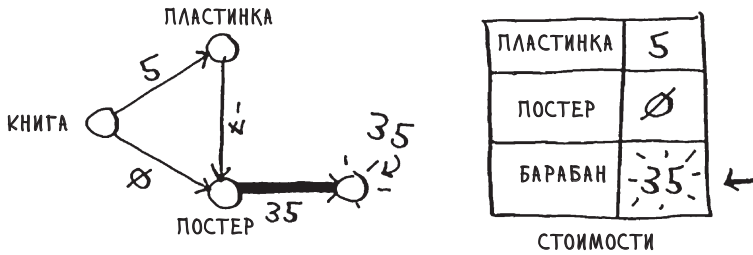
И знаете что? Если применить алгоритм Дейкстры к этому графу, Рама выберет неверный путь. Он пойдет по более длинному пути. Алгоритм Дейкстры не может использоваться при наличии ребер, имеющих отрицательный вес. Такие ребра нарушают работу алгоритма. Посмотрим, что произойдет, если попытаться применить алгоритм Дейкстры к этому графу. Все начинается с построения таблицы стоимостей.

ПЛАСТИНКА	5
ПОСТЕР	∅
БАРАБАН	∞

СТОИМОСТИ

Теперь найдем узел с наименьшей стоимостью и обновим стоимости его соседей. В этом случае постер оказывается узлом с наименьшей стоимостью. Итак, в соответствии с алгоритмом Дейкстры, к постеру *невозможно перейти*.

ти более дешевым способом, чем с оплатой \$0 (а вы знаете, что это неверно!)
Как бы то ни было, обновим стоимости его соседей.

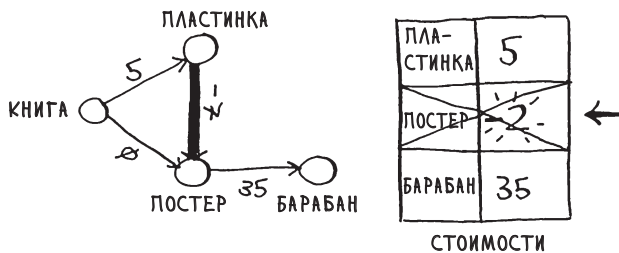


Получается, что теперь стоимость барабана составляет \$35.

Перейдем к следующему по стоимости узлу, который еще не был обработан.



Обновим стоимости его соседей.



Узел «постер» уже был обработан, однако вы обновляете его стоимость. Это очень тревожный признак — обработка узла означает, что к нему не-

возможно добраться с меньшими затратами. Но вы только что нашли более дешевый путь к постеру! У барабана соседей нет, поэтому работа алгоритма завершена. Ниже приведены итоговые стоимости.

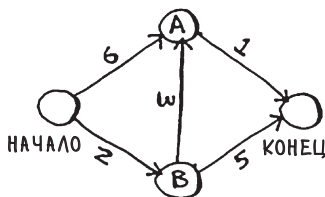
ПЛА- СТИНКА	5
ПОСТЕР	-2
БАРАБАН	35

ИТОГОВЫЕ
СТОИМОСТИ

Чтобы добраться до барабанов, Раме потребовалось \$35. Вы знаете, что существует путь, который стоит всего \$33, но алгоритм Дейкстры его не находит. Алгоритм Дейкстры предположил, что, поскольку вы обрабатываете узел «постер», к этому узлу невозможно добраться быстрее. Это предположение работает только в том случае, если ребер с отрицательным весом не существует. Следовательно, *использование алгоритма Дейкстры с графом, содержащим ребра с отрицательным весом, невозможно*. Если вы хотите найти кратчайший путь в графе, содержащем ребра с отрицательным весом, для этого существует специальный алгоритм, называемый *алгоритмом Беллмана — Форда*. Рассмотрение этого алгоритма выходит за рамки книги, но вы сможете найти хорошие описания в интернете.

Реализация

Посмотрим, как алгоритм Дейкстры реализуется в программном коде. Ниже изображен граф, который будет использоваться в этом примере.



Для реализации этого примера понадобятся три хеш-таблицы.

НАЧАЛО	А	6
	В	2
А	КОНЕЦ	1
В	А	3
		5
КОНЕЦ		—

ГРАФ
(GRAPH)

А	6
В	2
КОНЕЦ	∞

СТОИМОСТИ
(COSTS)

А	НАЧАЛО
В	НАЧАЛО
КОНЕЦ	—

РОДИТЕЛИ
(PARENTS)

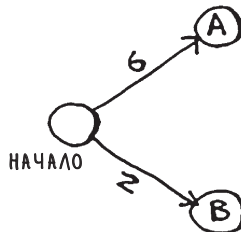
Хеш-таблицы стоимостей и родителей будут обновляться по ходу работы алгоритма. Сначала необходимо реализовать граф. Как и в главе 6, для этого будет использована хеш-таблица:

```
graph = {}
```

Ранее все соседи узла были сохранены в хеш-таблице:

```
graph["вы"] = ["Алиса", "Боб", "Клэр"]
```

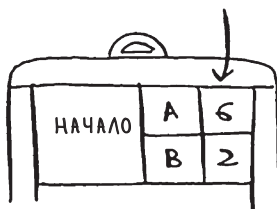
Но на этот раз необходимо сохранить как соседей, так и стоимость перехода к соседу. Предположим, у начального узла есть два соседа, А и В.



Как представить веса этих ребер? Почему бы не воспользоваться другой хеш-таблицей?

```
graph["начало"] = {}
graph["начало"]["a"] = 6
graph["начало"]["b"] = 2
```

ЭТА ХЕШ-ТАБЛИЦА
СОДЕРЖИТ ДРУГИЕ
ХЕШ-ТАБЛИЦЫ



Итак, `graph["начало"]` является хеш-таблицей. Для получения всех соседей начального узла можно воспользоваться следующим выражением:

```
>>> print graph["начало"].keys()
["a", "b"]
```

Одно ребро ведет из начального узла в А, а другое — из начального узла в В. А если вы захотите узнать веса этих ребер?

```
>>> print graph["начало"]["a"]
6
>>> print graph["начало"]["b"]
2
```

Включим в граф остальные узлы и их соседей:

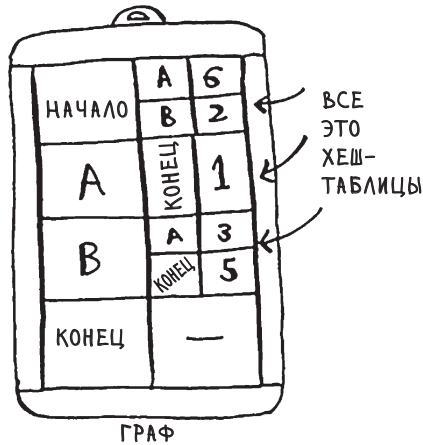
```
graph["a"] = {}
graph["a"]["fin"] = 1
```

```
graph["b"] = {}
graph["b"]["a"] = 3
```

```
graph["b"]["конец"] = 5
graph["конец"] = {}
```

У конечного узла нет соседей

Полная хеш-таблица графа выглядит так:



Также понадобится хеш-таблица для хранения стоимостей всех узлов.

Стоимость узла определяет, сколько времени потребуется для перехода к этому узлу от начального узла. Вы знаете, что переход от начального узла к узлу В занимает 2 минуты. Вы знаете, что для перехода к узлу А требуется 6 минут (хотя, возможно, вы найдете более быстрый

путь). Вы не знаете, сколько времени потребуется для достижения конечного узла. Если стоимость еще неизвестна, она считается бесконечной. Можно ли представить *бесконечность* (infinity) в Python? Оказывается, можно:

```
infinity = float("inf")
```

Код создания таблицы стоимостей costs:

```
infinity = float("inf")
costs = {}
costs["a"] = 6
costs["b"] = 2
costs["fin"] = infinity
```

Для родителей также создается отдельная таблица:

А	НАЧАЛО
В	НАЧАЛО
КОНЕЦ	—

РОДИТЕЛИ
(PARENTS)

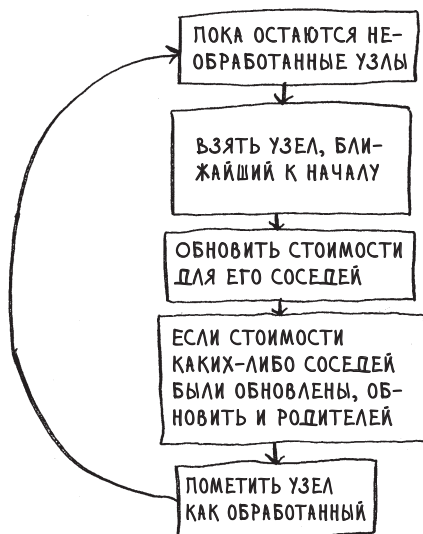
Код создания хеш-таблицы родителей:

```
parents = {}
parents["a"] = "начало"
parents["b"] = "начало"
parents["in"] = None
```

Наконец, вам нужен массив для отслеживания всех уже обработанных узлов, так как один узел не должен обрабатываться многократно:

```
processed = []
```

На этом подготовка завершается. Теперь обратимся к алгоритму.



Сначала я приведу код, а потом мы разберем его более подробно.

```

node = find_lowest_cost_node(costs)  ← ..... Найти узел с наименьшей стои-
while node is not None:             ← ..... мостью среди необработанных
    cost = costs[node]               ← ..... Если обработаны все узлы, цикл while завершен
    neighbors = graph[node]
    for n in neighbors.keys():       ← ..... Перебрать всех соседей (neighbors) текущего узла
        new_cost = cost + neighbors[n] ← ..... Если к соседу можно быстрее
        if costs[n] > new_cost:      ← ..... добраться через текущий узел...
            costs[n] = new_cost      ← ..... ...обновить стоимость для этого узла
            parents[n] = node         ← ..... Этот узел становится новым родителем для соседа
    processed.append(node)           ← ..... Узел помечается как обработанный
    node = find_lowest_cost_node(costs) ← ..... Найти следующий узел для
                                         обработки и повторить цикл

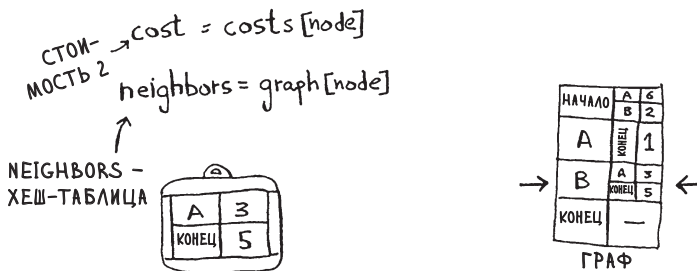
```

Так выглядит алгоритм Дейкстры на языке Python! Код функции будет приведен далее, а пока рассмотрим пример использования алгоритма в действии.

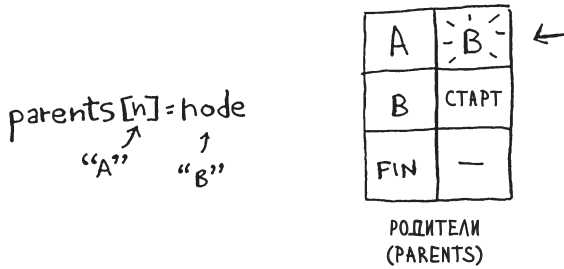
Найти узел с наименьшей стоимостью.



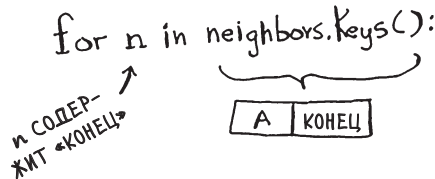
Получить стоимость и соседей этого узла.



Новый путь проходит через узел В, поэтому В назначается новым родителем.



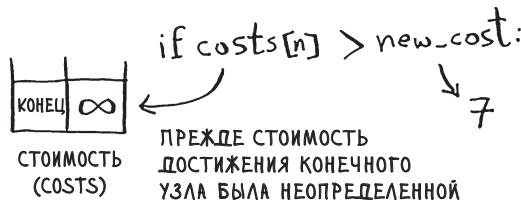
Мы снова вернулись к началу цикла. Следующим соседом в цикле for является конечный узел.



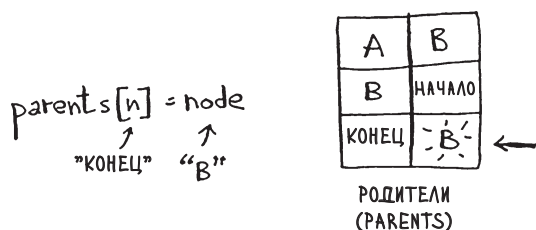
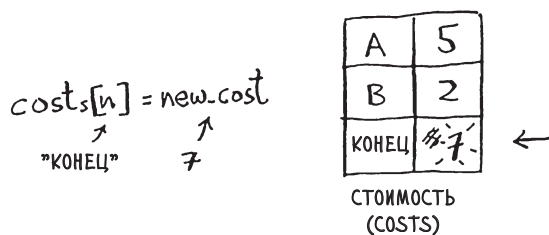
Сколько времени потребуется для достижения конечного узла, если идти через узел В?

$$\begin{array}{ccc}
 \text{new_cost} = \text{cost} + \text{neighbors}[n] & & \\
 \downarrow & \downarrow & \\
 2 & \text{РАССТОЯНИЕ} & \\
 & \text{ОТ В ДО КОНЦА:} & \\
 & 5 & \\
 & & \left. \vphantom{\begin{array}{c} 2 \\ 5 \end{array}} \right\} \begin{array}{l} 2+5 \\ = 7 \end{array}
 \end{array}$$

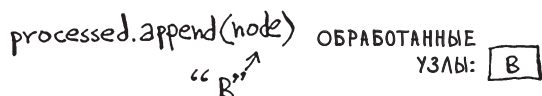
Потребуется 7 минут. Предыдущая стоимость была бесконечной, а 7 минут определенно меньше бесконечности.



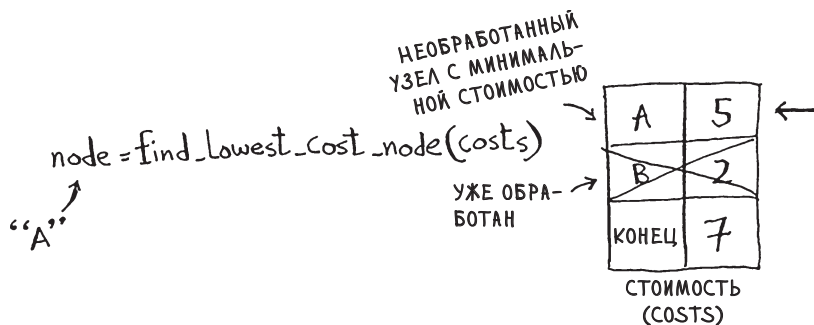
Конечному узлу назначается новая стоимость и новый родитель.



Порядок, мы обновили стоимости всех внешних соседей узла B. Узел помечается как обработанный.



Найти следующий узел для обработки.



Получить стоимость и внешних соседей узла A.

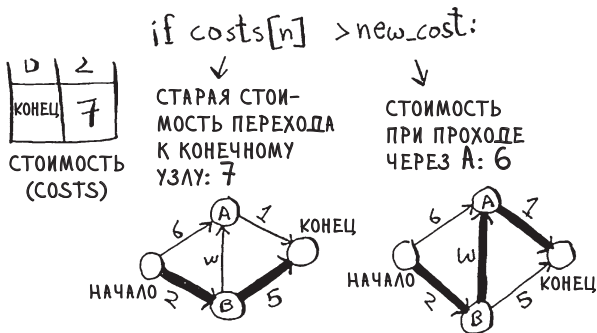
```
cost = costs[node]
5 ↑
neighbors = graph[node]
↑
[КОНЕЦ | 1]
```

У узла A всего один внешний сосед: конечный узел.

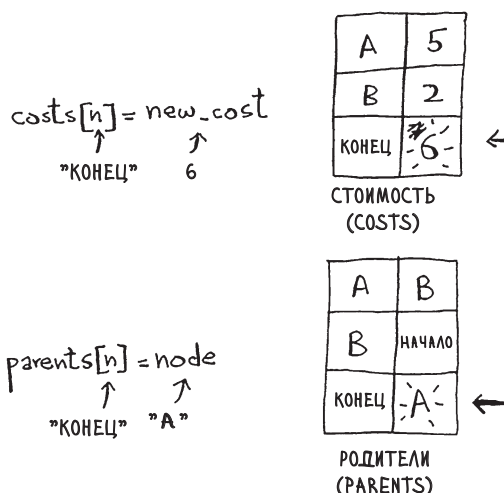
```
for n in neighbors.keys():
    ↑
    "КОНЕЦ"
    { "КОНЕЦ" }
```

Время достижения конечного узла составляет 7 минут. Сколько времени потребуется для достижения конечного узла, если идти через узел A?

```
new_cost = cost + neighbors[n]
      ↓      ↓
      ↓      ↓
СТОИМОСТЬ ПЕРЕХОДА К А    СТОИМОСТЬ ОТ А
ОТ НАЧАЛА: 5              ДО КОНЕЧНОГО
                          УЗЛА: 1
      } 5 + 1
      } = 6
```



Через узел A можно добраться быстрее! Обновим стоимость и родителя.



После того как все узлы будут обработаны, алгоритм завершается. Надеюсь, этот пошаговый разбор помог вам чуть лучше понять алгоритм. С функцией `find_lowest_cost_node` узел с наименьшей стоимостью находится проще простого. Код выглядит так:

```
def find_lowest_cost_node(costs):
    lowest_cost = float("inf")
    lowest_cost_node = None
    for node in costs:
        cost = costs[node]
        if cost < lowest_cost and node not in processed:
            lowest_cost = cost
            lowest_cost_node = node
    return lowest_cost_node
```

Перебрать все узлы

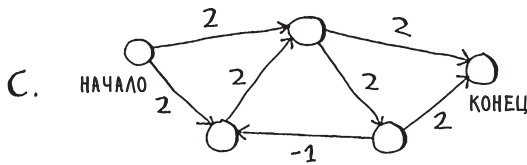
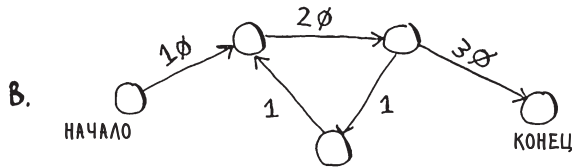
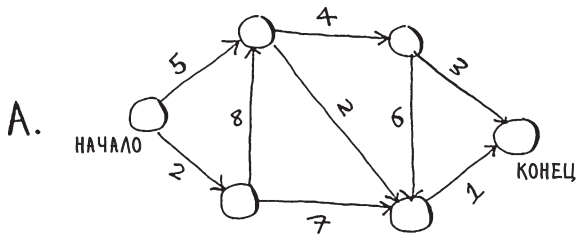
...он назначается новым узлом с наименьшей стоимостью

Если это узел с наименьшей стоимостью из уже виденных и он еще не был обработан...

Чтобы найти узел с наименьшей стоимостью, мы каждый раз перебираем все узлы. Существует более эффективный вариант этого алгоритма. Он использует структуру данных, называемую очередью с приоритетом. Эта очередь строится на основе другой структуры данных — кучи. Если вам интересны очереди с приоритетом и кучи, ознакомьтесь с разделом о кучах в последней главе книги.

УПРАЖНЕНИЕ

9.1 Каков вес кратчайшего пути от начала до конца в каждом из следующих графов?



Шпаргалка

- Поиск в ширину вычисляет кратчайший путь в невзвешенном графе.
- Алгоритм Дейкстры вычисляет кратчайший путь во взвешенном графе.
- Алгоритм Дейкстры работает только в том случае, если все веса положительны.
- При наличии отрицательных весов используйте алгоритм Беллмана — Форда.

10

Жадные алгоритмы

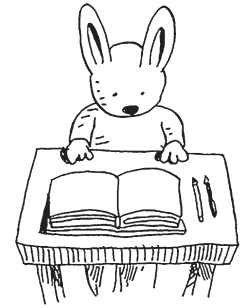


В этой главе

- ✓ Вы узнаете о жадной стратегии — очень простой стратегии решения задач.
- ✓ Вы узнаете, как браться за невозможные задачи, не имеющие быстрого алгоритмического решения (NP-трудные задачи).
- ✓ Вы познакомитесь с приближенными алгоритмами, которые могут использоваться для быстрого нахождения приближенного решения NP-полных задач.

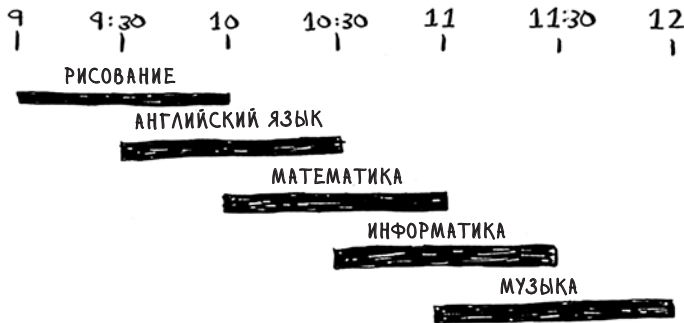
Задача составления расписания

Допустим, имеется учебный класс, в котором нужно провести как можно больше уроков. Вы получаете список уроков.



ПРЕДМЕТ	С	ДО
РИС.	8:00	10:00
АНГЛ.	8:30	10:30
МАТ-КА	10:00	11:00
ИНФ-КА	10:30	11:30
МУЗЫКА	11:00	12:00

Провести в классе *все* уроки не получится, потому что некоторые из них перекрываются по времени.



Требуется провести в классе как можно больше уроков. Как отобрать уроки, чтобы полученный набор оказался самым большим из возможных?

Вроде бы сложная задача, верно? На самом деле алгоритм оказывается на удивление простым. Вот как он работает:

1. Выбрать урок, завершающийся раньше всех. Это первый урок, который будет проведен в классе.
2. Затем выбирается урок, начинающийся после завершения первого урока. И снова следует выбрать урок, который завершается раньше всех остальных. Он становится вторым уроком в расписании.

Продолжайте действовать по тому же принципу — и вы получите ответ! Давайте попробуем. Рисование заканчивается раньше всех уроков (в 10:00), поэтому мы выбираем именно его.

РИС.	8:00	10:00	✓
АНГЛ.	8:30	10:30	
МАТ-КА	10:00	11:00	
ИНФ-КА	10:30	11:30	
МУЗЫКА	11:00	12:00	

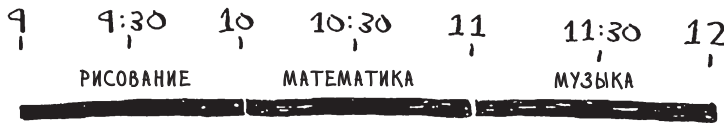
Теперь нужно найти следующий урок, который начинается после 10:00 и завершается раньше остальных.

РИС.	8:00	10:00	✓
АНГЛ.	8:30	10:30	✗
МАТ-КА	10:00	11:00	✓
ИНФ-КА	10:30	11:30	
МУЗЫКА	11:00	12:00	

Английский язык отпадает — он пересекается по времени с рисованием, но математика подходит. Наконец, информатика пересекается по времени с математикой, но музыка подходит.

РИС.	8:00	10:00	✓
АНГЛ.	8:30	10:30	✗
МАТ-КА	10:00	11:00	✓
ИНФ-КА	10:30	11:30	✗
МУЗЫКА	11:00	12:00	✓

Итак, эти три урока должны проводиться в классе.

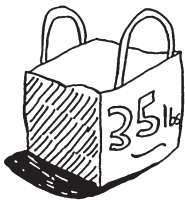


Я очень часто слышу, что этот алгоритм подозрительно прост. Он слишком очевиден, а значит, должен быть неправильным. Но в этом и заключается красота жадных алгоритмов: они просты! Жадный алгоритм прост: на каждом шаге он выбирает оптимальный вариант. В нашем примере выбирается тот урок, который завершается раньше других. В технической терминологии: на каждом шаге выбирается *локально-оптимальное решение*, а в итоге вы получаете глобально-оптимальное решение. Хотите верьте, хотите нет, но этот простой алгоритм успешно находит оптимальное решение задачи составления расписания!

Конечно, жадные алгоритмы работают не всегда. Но они так просто реализуются! Рассмотрим другой пример.

Задача о рюкзаке

Представьте, что вы жадный воришка. Вы забрались в магазин с рюкзаком, и перед вами множество товаров, которые вы можете украсть. Однако емкость рюкзака небесконечна: он выдержит не более 35 фунтов.



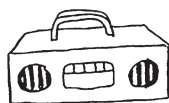
Требуется подобрать набор товаров максимальной стоимости, которые можно сложить в рюкзак. Какой алгоритм вы будете использовать?



И снова жадная стратегия выглядит очень просто:

1. Выбрать самый дорогой предмет, который поместится в рюкзаке.
2. Выбрать следующий по стоимости предмет, который поместится в рюкзаке... И так далее.

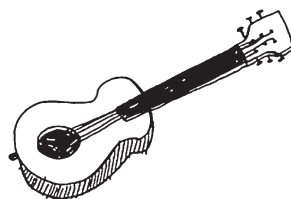
Вот только на этот раз она не работает! Предположим, есть три предмета.



МАГНИТОФОН
\$ 3000
30 ФУНТОВ

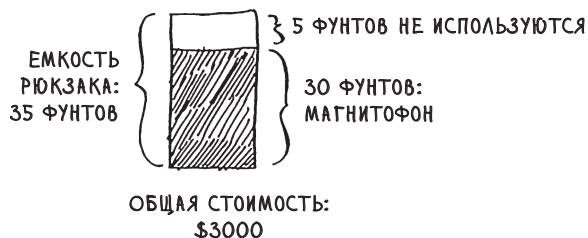


НОУТБУК
\$ 2000
20 ФУНТОВ

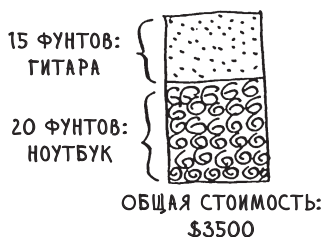


ГИТАРА
\$ 1500
15 ФУНТОВ

В рюкзаке поместятся товары общим весом не более 35 фунтов. Самый дорогой товар — магнитофон, вы выбираете его. Теперь ни для чего другого места уже не осталось.



Вы набрали товаров на \$3000. Погодите-ка! Если бы вместо магнитофона вы выбрали ноутбук и гитару, то стоимость добычи составила бы \$3500!



Очевидно, жадная стратегия не дает оптимального решения. Впрочем, результат не так уж далек от оптимума. В следующей главе я расскажу, как вычислить правильное решение. Но вор, забравшийся в магазин, вряд ли станет стремиться к идеалу. «Достаточно хорошего» решения должно хватить.

Второй пример приводит нас к следующему выводу: *иногда идеальное — враг хорошего*. В некоторых случаях достаточно алгоритма, способного решить задачу достаточно хорошо. И в таких областях жадные алгоритмы работают отлично, потому что они просто реализуются, а полученное решение обычно близко к оптимуму.

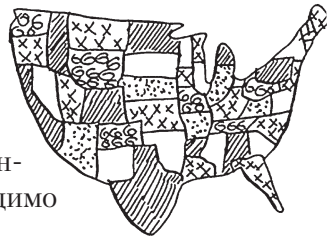
УПРАЖНЕНИЯ

- 10.1** Вы работаете в фирме по производству мебели и поставляете мебель по всей стране. Коробки с мебелью размещаются в грузовике. Все коробки имеют разный размер, и вы стараетесь наиболее эффективно использовать доступное пространство. Как выбрать коробки для того, чтобы загрузка имела максимальную эффективность? Предложите жадную стратегию. Будет ли полученное решение оптимальным?
- 10.2** Вы едете в Европу, и у вас есть семь дней на знакомство с достопримечательностями. Вы присваиваете каждой достопримечательности стоимость в баллах (насколько вы хотите ее увидеть) и оцениваете продолжительность поездки. Как обеспечить максимальную стоимость (увидеть все самое важное) во время поездки? Предложите жадную стратегию. Будет ли полученное решение оптимальным?

Рассмотрим еще один пример, в котором без жадных алгоритмов практически не обойтись.

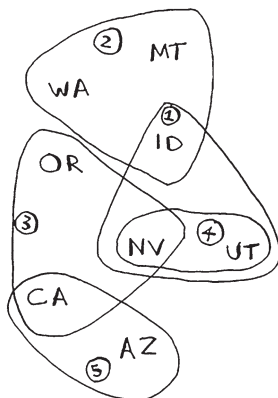
Задача о покрытии множества

Вы открываете собственную авторскую программу на радио и хотите, чтобы вас слушали во всех 50 штатах США. Нужно решить, на каких радиостанциях должна транслироваться ваша передача. Каждая станция стоит денег, поэтому количество станций необходимо свести к минимуму. Имеется список станций.



РАДИО- СТАНЦИЯ	ДОСТУПНА В ШТАТАХ
KONE	ID, NV, UT
KTWO	WA, ID, MT
KTHREE	OR, NV, CA
KFOUR	NV, UT
KFIVE	CA, AZ
... и т. д. ...	

Каждая станция покрывает определенный набор штатов, эти наборы перекрываются.



Как найти минимальный набор станций, который покрывал бы все 50 штатов? Вроде бы простая задача, верно? Оказывается, она чрезвычайно сложна. Вот как это делается:

1. Составить список всех возможных подмножеств станций — так называемое *степенное множество*. В нем содержатся 2^n возможных подмножеств.
2. Из этого списка выбирается множество с наименьшим набором станций, покрывающих все 50 штатов.



Проблема в том, что вычисление всех возможных подмножеств станций займет слишком много времени. Для n станций оно потребует времени $O(2^n)$. Если станций немного, скажем от 5 до 10, — это допустимо. Но подумайте, что произойдет во всех рассмотренных примерах при большом количестве элементов. Предположим, вы можете вычислять по 10 подмножеств в секунду.

Не существует алгоритма, который будет вычислять подмножества с приемлемой скоростью! Что же делать?

КОЛИЧЕСТВО СТАНЦИЙ	НЕОБХОДИМОЕ ВРЕМЯ
5	3.2 с
10	102.4 с
32	13.6 ГОДА
100	4×10^{21} ГОДА

Приближенные алгоритмы

На помощь приходят жадные алгоритмы! Вот как выглядит жадный алгоритм, который выдает результат, достаточно близкий к оптимуму:

1. Выбрать станцию, покрывающую наибольшее количество штатов, еще не входящих в покрытие. Если станция будет покрывать некоторые штаты, уже входящие в покрытие, это нормально.
2. Повторять, пока остаются штаты, не входящие в покрытие.

Этот алгоритм является *приближенным*. Когда вычисление точного решения занимает слишком много времени, применяется приближенный алгоритм. Эффективность приближенного алгоритма оценивается по:

- скорости;
- близости полученного решения к оптимальному.

Жадные алгоритмы хороши не только тем, что они обычно легко формулируются, но и тем, что простота обычно оборачивается скоростью выполнения. В данном случае жадный алгоритм выполняется за время $O(n^2)$, где n — количество радиостанций.

А теперь посмотрим, как эта задача выглядит в программном коде.

Подготовительный код

В этом примере для простоты будет использоваться небольшое подмножество штатов и станций.

Сначала составьте список штатов:

```
states_needed = set(["mt", "wa", "or", "id", "nv", "ut",
                    "ca", "az"])  ←..... Переданный массив преобразуется в множество
```

В этой реализации я использовал множество. Эта структура данных похожа на список, но каждый элемент может встречаться во множестве не более одного раза. *Множества не содержат дубликатов*. Предположим, имеется следующий список:

```
>>> arr = [1, 2, 2, 3, 3, 3]
```

Этот список преобразуется во множество:

```
>>> set(arr)
set([1, 2, 3])
```

Значения 1, 2 и 3 встречаются в списке по одному разу.

$[1, 2, 2, 3, 3, 3] \rightarrow \begin{matrix} \text{ПРЕОБРА-} \\ \text{ЗУЕТСЯ ВО} \\ \text{МНОЖЕСТВО} \end{matrix} \rightarrow \begin{matrix} (1, 2, 3) \\ \text{МНОЖЕСТВО} \end{matrix}$

Также понадобится список станций, из которого будет выбираться покрытие. Я решил воспользоваться хешем:

```
stations = {}
stations["kone"] = set(["id", "nv", "ut"])
stations["ktwo"] = set(["wa", "id", "mt"])
stations["kthree"] = set(["or", "nv", "ca"])
stations["kffour"] = set(["nv", "ut"])
stations["kfive"] = set(["ca", "az"])
```

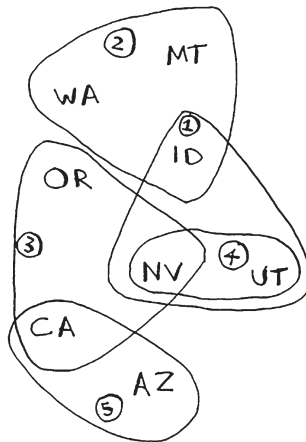
Ключи — названия станций, а значения — сокращенные обозначения штатов, входящих в зону охвата. Таким образом, в данном примере станция *kone* вещает в штатах Айдахо (*id*), Невада (*nv*) и Юта (*ut*). Все значения являются множествами. Как вы вскоре увидите, хранение данных во множествах упрощает работу.

Наконец, нам понадобится структура данных для хранения итогового набора станций:

```
final_stations = set()
```

Вычисление ответа

Теперь необходимо вычислить набор используемых станций. Взгляните на диаграмму и попробуйте предсказать, какие станции следует использовать.



Учтите, что правильных решений может быть несколько. Вы перебираете все станции и выбираете ту, которая обслуживает больше всего штатов, не входящих в текущее покрытие. Будем называть ее `best_station`:

```
best_station = None
states_covered = set()
for station, states_for_station in stations.items():
```

Множество `states_covered` содержит все штаты, обслуживаемые этой станцией, которые еще не входят в текущее покрытие. Помните, мы ищем станцию, которая будет обслуживать большинство еще не охваченных штатов. Цикл `for` перебирает все станции и находит среди них наилучшую. Рассмотрим тело цикла `for`:

```
covered = states_needed & states_for_station
if len(covered) > len(states_covered):
    best_station = station
    states_covered = covered
```

Новый синтаксис! Эта операция называется пересечением множеств

В коде встречается необычная строка:

```
covered = states_needed & states_for_station
```

Что здесь происходит?

Множества

Допустим, имеется множество с названиями фруктов.



Также имеется множество с названиями овощей.



С двумя множествами можно выполнить ряд интересных операций.



- Объединение множеств означает слияние элементов обоих множеств.
- Под операцией пересечения множеств понимается поиск элементов, входящих в оба множества (в данном случае — только помидор).
- Под разностью множеств понимается исключение из одного множества элементов, присутствующих в другом множестве.

Пример:

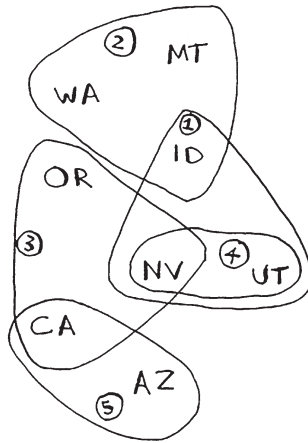
```
>>> fruits = set(["авокадо", "помидор", "банан"])
>>> vegetables = set(["свекла", "морковь", "помидор"])
>>> fruits | vegetables  <----- Объединение множеств
set(["авокадо", "свекла", "морковь", "помидор", "банан"])
>>> fruits & vegetables  <----- Пересечение множеств
set(["помидор"])
>>> fruits - vegetables  <----- Разность множеств
set(["авокадо", "банан"])
>>> vegetables - fruits  <----- Как вы думаете, как будет выглядеть результат?
```

Еще раз напомним основные моменты:

- множества похожи на списки, но не содержат дубликатов;
- с множествами можно выполнять различные интересные операции — вычислять их объединение, пересечение и разность.

Вернемся к коду

Продолжим рассматривать исходный пример.



Пересечение множеств:

```
covered = states_needed & states_for_station
```

Множество `covered` содержит штаты, присутствующие как в `states_needed`, так и в `states_for_station`. Таким образом, `covered` — множество штатов, не входящих в покрытие, которые покрываются текущей станцией! Затем мы проверяем, покрывает ли эта станция больше штатов, чем текущая станция `best_station`:

```
if len(covered) > len(states_covered):
    best_station = station
    states_covered = covered
```

Если условие выполняется, то станция сохраняется в `best_station`. Наконец, после завершения цикла `best_station` добавляется в итоговый список станций:

```
final_stations.add(best_station)
```

Также необходимо обновить содержимое `states_needed`. Те штаты, которые входят в зону покрытия станции, больше не нужны:

```
states_needed -= states_covered
```

Цикл продолжается, пока множество `states_needed` не станет пустым. Полный код цикла `for` выглядит так:

```
while states_needed:
    best_station = None
    states_covered = set()
    for station, states in stations.items():
        covered = states_needed & states
        if len(covered) > len(states_covered):
            best_station = station
            states_covered = covered
    states_needed -= states_covered
    final_stations.add(best_station)
```

Остается вывести содержимое `final_stations`:

```
>>> print final_stations
set(['ktwo', 'kthree', 'kone', 'kfive'])
```

Этот результат совпадает с вашими ожиданиями? Вместо станций 1, 2, 3 и 5 можно было выбрать станции 2, 3, 4 и 5. Сравним время выполнения жадного алгоритма со временем точного алгоритма.

	$O(2^n)$	$O(n^2)$
КОЛИЧЕСТВО СТАНЦИЙ	ТОЧНЫЙ АЛГОРИТМ	ЖАДНЫЙ АЛГОРИТМ
5	3.2 с	2.5 с
10	102.4 с	10 с
32	13.6 ГОДА	102.4 с
100	4×10^{24} ГОДА	16.67 МИН

Жадный алгоритм не всегда даст точный ответ, но он очень быстр. Задача о покрытии множеств относится к NP-трудным задачам. Если вы хотите узнать немного больше об NP-трудных задачах, обратитесь к приложению Б.

Шпаргалка

- Жадные алгоритмы стремятся к локальной оптимизации в расчете на то, что в итоге будет достигнут глобальный оптимум.
- Если у вас имеется NP-трудная задача, лучше всего воспользоваться приближенным алгоритмом.
- Жадные алгоритмы легко реализовать и быстро выполнить, поэтому из них получаются хорошие приближенные алгоритмы.

11

Динамическое программирование



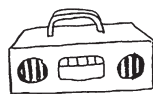
В этой главе

- ✓ Вы освоите динамическое программирование — метод решения сложных задач, разбиваемых на подзадачи, которые решаются в первую очередь.
- ✓ Рассматриваются примеры, которые научат вас искать решения новых задач, основанные на методе динамического программирования.

И снова задача о рюкзаке

Вернемся к задаче о рюкзаке из главы 10. У вас есть рюкзак, в котором можно унести товары общим весом до 4 фунтов.

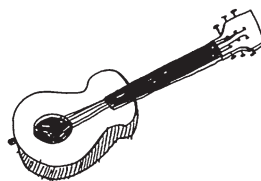




МАГНИТОФОН
\$ 3000
4 ФУНТА



НОУТБУК
\$ 2000
3 ФУНТА



ГИТАРА
\$ 1500
1 ФУНТ

Есть три предмета, которые можно уложить в рюкзак.

Какие предметы следует положить в рюкзак, чтобы стоимость добычи была максимальной?

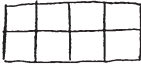
Простое решение

Простой алгоритм выглядит так: вы перебираете все возможные множества товаров и находите множество с максимальной стоимостью.



Такое решение работает, но очень медленно. Для 3 предметов приходится обработать 8 возможных множеств, для 4 — 16 и т. д. С каждым добавляемым предметом количество множеств удваивается! Этот алгоритм выполняется за время $O(2^n)$, что очень, очень медленно.

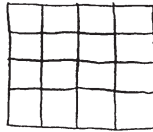
3 ПРЕДМЕТА:



8

ВОЗМОЖНЫХ
МНОЖЕСТВ

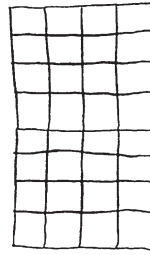
4 ПРЕДМЕТА:



16

ВОЗМОЖНЫХ
МНОЖЕСТВ

5 ПРЕДМЕТОВ:



32

ВОЗМОЖНЫХ
МНОЖЕСТВ

32 ПРЕДМЕТА =
~4 МИЛЛИАРДА
ВОЗМОЖНЫХ
МНОЖЕСТВ

Для любого сколько-нибудь значительного количества предметов это неприемлемо. В главе 8 вы видели, как вычисляются *приближенные* решения. Такие решения близки к оптимальным, но могут не совпадать с ними.

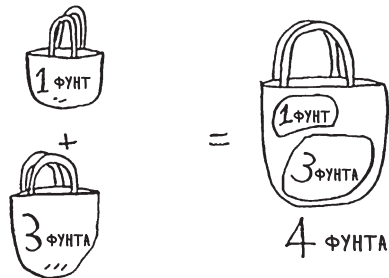
Как же вычислить оптимальное решение?

Динамическое программирование

Ответ: с помощью динамического программирования! Давайте посмотрим, как работает этот метод. Процедура начинается с решения подзадач с постепенным переходом к решению полной задачи.

В задаче о рюкзаке начать следует с решения задачи для меньшего рюкзака (или «подрюкзака»), а потом на этой основе попытаться решить исходную задачу.

Динамическое программирование — достаточно сложная концепция; не огорчайтесь, если после первого прочтения что-то останется непонятным. Примеры помогут вам разобраться в теме.



Для начала я покажу вам алгоритм в действии. После этого у вас наверняка появится много вопросов! Я постараюсь ответить на них.

Каждый алгоритм динамического программирования начинается с таблицы. Вот как выглядит таблица для задачи о рюкзаке.

СТОЛБЦЫ ПРЕДСТАВЛЯЮТ
РАЗМЕРЫ РЮКЗАКА
ОТ 1 ДО 4 ФУНТОВ

		1	2	3	4
ПО ОДНОЙ СТРОКЕ ДЛЯ КАЖДОГО ПРЕДМЕТА	ГИТАРА				
	МАГНИТОФОН				
	НОУТБУК				

Строки таблицы представляют предметы, а столбцы — емкость рюкзака от 1 до 4 фунтов. Все эти столбцы нужны, потому что они упрощают вычисление стоимостей «подрюкзаков».

В исходном состоянии таблица пуста. Нам предстоит заполнить каждую ячейку таблицы. После того как таблица будет заполнена, вы получите ответ на свою задачу. Пожалуйста, внимательно разберитесь в происходящем. Нарисуйте собственную таблицу, а мы вместе ее заполним.

Гитара

Точная формула для вычисления значений в таблице будет приведена позднее, а пока ограничимся общим описанием. Начнем с первой строки.

		1	2	3	4
ГИТАРА					
МАГНИТОФОН					
НОУТБУК					

Строка снабжена пометкой «гитара»; это означает, что вы пытаетесь уложить гитару в рюкзак. В каждой ячейке принимается простое решение: класть гитару в рюкзак или нет? Помните: мы пытаемся найти множество элементов с максимальной стоимостью.

В первой ячейке емкость рюкзака равна 1 фунту. Гитара также весит 1 фунт — значит, она поместится в рюкзак! Итак, стоимость этой ячейки составляет \$1500, а в рюкзаке лежит гитара.

Начнем заполнять ячейку.

	1	2	3	4
ГИТАРА	\$1500 Г			
МАГНИТОФОН				
НОУТБУК				

По тому же принципу каждая ячейка в таблице содержит список всех элементов, которые помещаются в рюкзак на данный момент.

Посмотрим на следующую ячейку. На этот раз емкость рюкзака составляет 2 фунта. Понятно, что гитара здесь поместится!

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г		
МАГНИТОФОН				
НОУТБУК				

Процедура повторяется для остальных ячеек строки. Вспомните, что текущей является первая строка, поэтому выбирать приходится *только* из одного предмета — гитары. Считайте, что два других предмета пока недоступны.

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН				
НОУТБУК				

Возможно, к этому моменту вы слегка сбиты с толку. Почему все это делается для рюкзаков с емкостью 1, 2 и т. д., если в задаче речь идет о рюкзаке с емкостью 4 фунта? Помните, что я говорил ранее? Метод динамического программирования начинается с малых задач, а затем переходит к большой задаче. Вы решаете подзадачи, которые помогут в решении большой задачи. Читайте дальше, и ситуация постепенно прояснится.

После того как первая строка будет заполнена, таблица будет выглядеть так:

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН				
НОУТБУК				

Помните, что мы стремимся обеспечить максимальную стоимость предметов в рюкзаке. *Эта строка представляет текущую лучшую оценку максимума.* Итак, на данный момент из этой строки следует, что для рюкзака с емкостью 4 фунта максимальная стоимость предметов составит \$1500.

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН				
НОУТБУК				

← НАША ТЕКУЩАЯ
ОЦЕНКА ТОГО,
ЧТО СЛЕДУЕТ
КЛАСТЬ: ГИТАРУ
ЗА \$1500

Вы знаете, что это решение неокончательно. В процессе работы алгоритма оценка будет уточняться.

Магнитофон

Займемся следующей строкой, которая относится к магнитофону. Теперь, когда вы перешли ко второй строке, появляется выбор между магнитофоном и гитарой. В каждой строке можно взять предмет этой строки или предметы, находящиеся в верхних строках. Таким образом, сейчас нельзя выбрать ноутбук, но можно выбрать магнитофон и/или гитару. Начнем с первой ячейки (рюкзак с емкостью 1 фунт). Текущая максимальная стоимость предметов, которые можно положить в рюкзак с емкостью 1 фунт, составляет \$1500.

ТЕКУЩАЯ МАКСИМАЛЬНАЯ ОЦЕНКА ДЛЯ РЮКЗАКА С ЕМКОСТЬЮ 1 ФУНТ

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН				
НОУТБУК				

НОВЫЙ МАКСИМУМ ДЛЯ РЮКЗАКА С ЕМКОСТЬЮ 1 ФУНТ

Брать магнитофон или нет?

Емкость рюкзака составляет 1 фунт. Поместится туда магнитофон? Нет, он слишком тяжел! Так как магнитофон не помещается в рюкзак, максимальная оценка для 1-фунтового рюкзака *остается* равной \$1500.

	1	2	3	4
ГИТАРА	\$1500 ↓ Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН	\$1500 Г			
НОУТБУК				

То же самое происходит со следующими двумя клетками. Емкость этих рюкзаков составляет 2 и 3 фунта соответственно. Старая максимальная стоимость для обеих ячеек была равна \$1500.

	1	2	3	4
ГИТАРА	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	
НОУТБУК				

Магнитофон все равно не помещается, так что оценка остается неизменной.

А если емкость рюкзака увеличивается до 4 фунтов? Ага, магнитофон наконец-то войдет в рюкзак! Старая максимальная стоимость была равна \$1500, но если вместо гитары положить магнитофон, она увеличится до \$3000! Берем магнитофон.

	1	2	3	4
ГИТАРА	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 Г М
НОУТБУК				

Оценка только что обновилась! Имея рюкзак емкостью 4 фунта, вы можете положить в него товары стоимостью по крайней мере \$3000. Из таблицы видно, что оценка постепенно возрастает.

	1	2	3	4	
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г	← СТАРАЯ ОЦЕНКА
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М	← НОВАЯ ОЦЕНКА
НОУТБУК					← ИТОГОВОЕ РЕШЕНИЕ

Ноутбук

А теперь сделаем то же для ноутбука! Ноутбук весит 3 фунта, поэтому он не поместится в рюкзак с емкостью 1 или 2 фунта. Оценка для первых двух ячеек остается на уровне \$1500.

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
НОУТБУК	\$1500 Г	\$1500 Г		

Для 3 фунтов старая оценка составляла \$1500. Но теперь вы можете выбрать ноутбук, который стоит \$2000. Следовательно, новая максимальная оценка равна \$2000!

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
НОУТБУК	\$1500 Г	\$1500 Г	\$2000 Н	

При 4 фунтах ситуация становится по-настоящему интересной. Это очень важная часть. В настоящее время оценка составляет \$3000. В рюкзак можно положить ноутбук, но он стоит всего \$2000.

$$\begin{array}{cc} \$3000 & \text{или} & \$2000 \\ \text{МАГНИТОФОН} & & \text{НОУТБУК} \end{array}$$

Так-так, старая оценка была лучше. Но постойте! Ноутбук весит всего 3 фунта, так что 1 фунт еще свободен! На это место можно еще что-нибудь положить.

$$\begin{array}{cc} \$3000 & \text{или} & \left(\$2000 + \frac{???}{\text{СВОБОДНОЕ МЕСТО НА 1 ФУНТ}} \right) \\ \text{МАГНИТОФОН} & & \text{НОУТБУК} \end{array}$$

Какую максимальную стоимость можно разместить в 1 фунте? Да вы же уже вычислили ее!

МАКСИМАЛЬНАЯ СТОИМОСТЬ ДЛЯ 1 ФУНТА →

	1	2	3	4
	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г
	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г	\$3000 М
	\$1500 ↓ Г	\$1500 ↓ Г	\$2000 ↓ Н	

В соответствии с последней оценкой в свободном месте емкостью в 1 фунт можно разместить гитару стоимостью \$1500. Следовательно, настоящее сравнение выглядит так:

$$\begin{array}{cc} \$3000 & \text{или} & \left(\$2000 + \$1500 \right) \\ \text{МАГНИТОФОН} & & \text{НОУТБУК} \quad \text{ГИТАРА} \end{array}$$

Вы удивлялись, зачем мы вычисляем максимальную стоимость для рюкзаков меньшей емкости? Надеюсь, теперь все стало на свои места! Если в рюкзаке остается свободное место, вы можете использовать ответы на эти подзадачи для определения того, чем заполнить это пространство. Вместо магнитофона лучше взять ноутбук + гитару за \$3500.

В завершающем состоянии таблица выглядит так:

	1	2	3	4
ГИТАРА	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г
МАГНИТОФОН	\$1500 ↓ Г	\$1500 ↓ Г	\$1500 ↓ Г	\$3000 ↓ М
НОУТБУК	\$1500 ↓ Г	\$1500 ↓ Г	\$2000 ↓ Н	\$3500 ↓ Н Г

↑
ОТВЕТ!

Итак, мы получили ответ: максимальная стоимость товаров, которые поместятся в рюкзак, равна \$3500 — для гитары и ноутбука.

Возможно, вы подумали, что я воспользовался другой формулой для вычисления стоимости последней ячейки. Это связано с тем, что я опустил некоторые лишние сложности при заполнении предыдущих ячеек. Стоимость каждой ячейки вычисляется по постоянной формуле, которая выглядит так:

$$\begin{array}{c} \text{СТОЛ-} \\ \text{СТРОКА} \\ \downarrow \\ \text{CELL}[i][j] \end{array} = \text{МАКСИМУМ} \left\{ \begin{array}{l} 1. \text{ ПРЕДЫДУЩИЙ МАКСИМУМ (ЗНАЧЕНИЕ В CELL}[i-1][j]) \\ \text{ИЛИ} \\ 2. \text{ СТОИМОСТЬ ТЕКУЩЕГО ЭЛЕМЕНТА +} \\ \text{СТОИМОСТЬ ОСТАВШЕГОСЯ ПРОСТРАНСТВА} \\ \uparrow \\ \text{CELL}[i-1][j] - \text{ВЕС ПРЕДМЕТА} \end{array} \right.$$

Применяя эту формулу к каждой ячейке таблицы, вы получите такую же таблицу, как у меня. Помните, что я говорил о решении подзадач?

Вы объединили решения двух подзадач для решения еще одной, большей задачи.



Задача о рюкзаке: вопросы

Вам все еще кажется, что это какой-то фокус? В этом разделе я отвечу на некоторые часто задаваемые вопросы.



IPHONE
\$2000
1 ФУНТ

Что произойдет при добавлении элемента?

Представьте, что вы увидели четвертый предмет, который тоже можно засунуть в рюкзак! Вместе со всем предыдущим добром можно также украсть iPhone.

Придется ли пересчитывать все заново с новым предметом? Нет. Напомню, что динамическое программирование последовательно строит решение на основании вашей оценки. К настоящему моменту максимальные стоимости выглядят так:

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
НОУТБУК	\$1500 Г	\$1500 Г	\$2000 Н	\$3500 Н Г

Это означает, что в рюкзак с емкостью 4 фунта можно упаковать товары стоимостью до \$3500. И вы полагали, что это итоговый максимум. Но давайте добавим новую строку для iPhone.

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
НОУТБУК	\$1500 Г	\$1500 Г	\$2000 Н	\$3500 Н Г
IPHONE				

↑
НОВЫЙ ОТВЕТ

Оказывается, в таблице появляется новый максимум! Попробуйте заполнить последнюю строку, прежде чем читать дальше.

Начнем с первой ячейки. iPhone сам по себе помещается в рюкзак с емкостью 1 фунт. Старый максимум был равен \$1500, но iPhone стоит \$2000. Значит, берем iPhone.

	1	2	3	4
ГИТАРА	\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
МАГНИТОФОН	\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
НОУТБУК	\$1500 Г	\$1500 Г	\$2000 Н	\$3500 Н Г
IPHONE	\$2000 И			

В следующей ячейке можно разместить iPhone и гитару.

\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
\$1500 Г	\$1500 Г	\$2000 Н	\$3500 Н Г
\$2000 I	\$3500 I Г		

Для ячейки 3 ничего лучшего, чем снова взять iPhone вместе с гитарой, все равно не найдется, поэтому оставим этот вариант.

А вот в последней ячейке ситуация становится более интересной. Текущий максимум равен \$3500. Вы снова можете взять iPhone, и у вас еще останется свободное место на 3 фунта.

$$\begin{array}{ccc}
 \$3500 & \text{или} & \left(\$2000 + \frac{???}{\text{свободное место на 3 фунта}} \right) \\
 \text{НОУТБУК+ГИТАРА} & & \begin{array}{c} \text{IPHONE} \\ \text{свободное место на 3 фунта} \end{array}
 \end{array}$$

Но эти 3 фунта можно заполнить на \$2000! \$2000 от iPhone + \$2000 из старой подзадачи: получается \$4000. Новый максимум!

Вот как выглядит новая завершающая таблица.

\$1500 Г	\$1500 Г	\$1500 Г	\$1500 Г
\$1500 Г	\$1500 Г	\$1500 Г	\$3000 М
\$1500 Г	\$1500 Г	\$2000 Н	\$3500 Н Г
\$3500 I	\$3500 I Г	\$3500 I Г	\$4000 I Н

↑
НОВЫЙ ОТВЕТ

Вопрос: может ли значение в столбце *уменьшиться*? Такое возможно?

МАКСИМАЛЬНАЯ
СТОИМОСТЬ
УМЕНЬШАЕТСЯ
В ПРОЦЕССЕ
РАБОТЫ

	1	2	3	4
	\$1500	\$1500	\$1500	\$1500
	∅	∅	∅	\$3000

Подумайте над ответом, прежде чем продолжить чтение.

Ответ: нет. При каждой итерации сохраняется текущая оценка максимума. Эта оценка ни при каких условиях не может быть меньше предыдущей!

УПРАЖНЕНИЕ

11.1 Предположим, к предметам добавился еще один: механическая клавиатура. Она весит 1 фунт и стоит \$1000. Стоит ли ее брать?

Что произойдет при изменении порядка строк?

Изменится ли ответ? Допустим, строки заполняются в другом порядке: магнитофон, ноутбук, гитара. Как будет выглядеть таблица? Заполните таблицу самостоятельно, прежде чем двигаться дальше.

Таблица должна выглядеть так:

	1	2	3	4
МАГНИТОФОН	∅	∅	∅	\$3000 М
НОУТБУК	∅	∅	\$2000 Н	\$3000 М
ГИТАРА	\$1500 Г	\$1500 Г	\$2000 Н	\$3500 Н Г

Ответ не изменился. Он не зависит от порядка строк.

Можно ли заполнять таблицу по столбцам, а не по строкам?

Попробуйте сами! В данной задаче это ни на что не влияет, но в других задачах возможны изменения.

Что произойдет при добавлении меньшего элемента?

Допустим, вы можете выбрать ожерелье, которое весит 0,5 фунта и стоит \$1000. Пока структура таблицы предполагает, что все веса являются целыми числами. Теперь вы решаете взять ожерелье. Остается еще 3,5 фунта. Какую максимальную стоимость можно разместить в объеме 3,5 фунта? Неизвестно! Вы вычисляли стоимость только для рюкзаков с емкостью 1, 2, 3 и 4 фунта. Теперь придется определять стоимость для рюкзака на 3,5 фунта.

Из-за ожерелья приходится повысить точность представления весов, поэтому таблица должна измениться.

	0.5	1	1.5	2	2.5	3	3.5	4
ГИТАРА								
МАГНИТОФОН								
НОУТБУК								
ОЖЕРЕЛЬЕ								

Можно ли взять часть предмета?

Допустим, вы наполняете рюкзак в продуктовом магазине. Вы можете украсть мешки с чечевицей и рисом. Если весь мешок не помещается, его можно открыть и отсыпать столько, сколько унесете. В этом случае вы уже

не действуете по принципу «все или ничего» — можно взять только часть предмета. Как решить такую задачу методом динамического программирования?

Ответ: никак. В решении, полученном методом динамического программирования, вы либо берете предмет, либо не берете. Алгоритм не предусматривает возможности взять половину предмета.

Однако проблема легко решается с помощью жадного алгоритма! Сначала вы берете самый ценный предмет — настолько большую его часть, насколько возможно. Когда самый ценный предмет будет исчерпан, вы берете максимально возможную часть следующего по ценности предмета и т. д.

Допустим, вы можете выбрать из следующих товаров.



Фунт киноа стоит дороже, чем фунт любого другого товара. А раз так, набирайте столько киноа, сколько сможете унести! И если вам удастся набить им свой рюкзак, то это и будет лучшее из возможных решений.



Если киноа кончится, а в рюкзаке еще остается свободное место, возьмите следующий по ценности товар и т. д.

Оптимизация туристического маршрута

Представьте, что вы приехали в Лондон на выходные. У вас два дня, а мест, которые хочется посетить, слишком много. Побывать везде не получится, поэтому вы составляете список.

ДОСТОПРИМЕЧАТЕЛЬНОСТЬ	ВРЕМЯ	ОЦЕНКА
ВЕСТМИНСТЕРСКОЕ АББАТСТВО	$\frac{1}{2}$ ДНЯ	7
ТЕАТР «ГЛОБУС»	$\frac{1}{2}$ ДНЯ	6
НАЦИОНАЛЬНАЯ ГАЛЕРЕЯ	1 ДЕНЬ	9
БРИТАНСКИЙ МУЗЕЙ	2 ДНЯ	9
СОБОР СВ. ПАВЛА	$\frac{1}{2}$ ДНЯ	8

Для каждой достопримечательности, которую вы захотите увидеть, вы указываете, сколько времени займет осмотр и насколько сильно вы хотите ее увидеть. Сможете ли вы построить оптимальный туристический маршрут на основании этого списка?

Да это все та же задача о рюкзаке! Вместо ограниченной емкости рюкзака — ограниченное время. Вместо магнитофонов и ноутбуков — список мест, которые вы хотите посетить. Нарисуйте таблицу динамического программирования для списка, прежде чем двигаться дальше.

Вот как должна выглядеть эта таблица:

	$\frac{1}{2}$	1	$1\frac{1}{2}$	2
ВЕСТМИНСТЕР				
ТЕАТР «ГЛОБУС»				
НАЦИОНАЛЬНАЯ ГАЛЕРЕЯ				
БРИТАНСКИЙ МУЗЕЙ				
СОБОР СВ. ПАВЛА				

Вы изобразили ее правильно? Теперь заполните. Какие достопримечательности вы выберете? Ответ:

	$\frac{1}{2}$	1	$1\frac{1}{2}$	2
ВЕСТМИНСТЕР	7 _w	7 _w	7 _w	7 _w
ТЕАТР «ГЛОБУС»	7 _w	13 _{wg}	13 _{wg}	13 _{wg}
НАЦИОНАЛЬНАЯ ГАЛЕРЕЯ	7 _w	13 _{wg}	16 _{wn}	22 _{wgn}
БРИТАНСКИЙ МУЗЕЙ	7 _w	13 _{wg}	16 _{wn}	22 _{wgn}
СОБОР СВ. ПАВЛА	8 _s	15 _{ws}	21 _{wgs}	24 _{wns}

↑
ОТВЕТ:
ВЕСТМИНСТЕРСКОЕ АББАТСТВО,
НАЦИОНАЛЬНАЯ ГАЛЕРЕЯ,
СОБОР СВ. ПАВЛА

Взаимозависимые элементы

Предположим, вы хотите посетить Париж и добавили в свой список пару элементов.

ЭЙФЕЛЕВА БАШНЯ	$1\frac{1}{2}$ DAY	8
ЛУВР	$1\frac{1}{2}$ DAY	9
НОТР-ДАМ	$1\frac{1}{2}$ DAY	7

На их посещение потребуется много времени, потому что сначала придется приехать из Лондона в Париж. Переезд отнимает полдня. Если вы захотите посмотреть все 3 достопримечательности, осмотр займет 4,5 дня.

Стоп, небольшая поправка. Вам не обязательно приезжать в Париж ради каждой достопримечательности. После того как вы там окажетесь, каждый

последующий элемент займет всего один день. Следовательно, потребуется 1 день на каждую достопримечательность + 1 день на переезды = 3,5 дня, а не 4,5.

Если вы положите Эйфелеву башню в свой «рюкзак», то Лувр станет «дешевле» — он займет всего 1 день вместо 1,5. Как смоделировать это обстоятельство в динамическом программировании?

Никак. Динамическое программирование — мощный метод, способный решать подзадачи и использовать полученные ответы для решения большой задачи. *Динамическое программирование работает только в том случае, если каждая подзадача автономна, то есть не зависит от других подзадач.* Из этого следует, что учесть поездки в Париж в алгоритме динамического программирования не удастся.

Может ли оказаться, что решение требует более двух «подрюкзаков»?

Может оказаться, что в лучшем решении должны отбираться больше двух элементов. В текущем варианте алгоритма объединяются не более двух «подрюкзаков» — больше двух их не бывает. Однако вполне возможно, что у этих «подрюкзаков» будут собственные «подрюкзаки».



Возможно ли, что при лучшем решении в рюкзаке остается пустое место?



БРИЛЛИАНТ.
СТОИМОСТЬ \$1 000 000,
ВЕС 3,5 ФУНТА

Да. Представьте, что вы можете также положить в рюкзак бриллиант.

Бриллиант очень крупный: он весит 3,5 фунта и стоит 1 миллион долларов — намного больше, чем любые другие предметы. Безусловно, нужно брать именно его! Но в рюкзаке остается еще пустое место на 0,5 фунта, и в нем ничего не поместится.

УПРАЖНЕНИЕ

11.2 Предположим, что вы собираетесь в турпоход. Емкость вашего рюкзака составляет 6 фунтов, и вы можете взять предметы из следующего списка. У каждого предмета имеется стоимость; чем она выше, тем важнее предмет:

- вода, 3 фунта, 10;
- книга, 1 фунт, 3;
- еда, 2 фунта, 9;
- куртка, 2 фунта, 5;
- камера, 1 фунт, 6

Как выглядит оптимальный набор предметов для похода?

Самая длинная общая подстрока

Мы рассмотрели одну задачу динамического программирования. Какие выводы из нее можно сделать?

- Динамическое программирование применяется для оптимизации какой-либо характеристики при заданных ограничениях. В задаче о рюкзаке



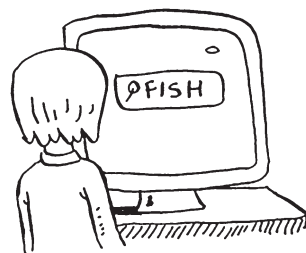
требуется максимизировать стоимость отобранных предметов с ограничениями по емкости рюкзака.

- Динамическое программирование работает только в ситуациях, в которых задача может быть разбита на автономные подзадачи, не зависящие друг от друга.

Построить решение на базе динамического программирования бывает не просто. В этом разделе мы сосредоточимся на этой теме. Несколько общих рекомендаций:

- в каждом решении из области динамического программирования строится таблица;
- значения ячеек таблицы обычно соответствуют оптимизируемой характеристике. Для задачи о рюкзаке значения представляли общую стоимость товаров;
- каждая ячейка представляет подзадачу, поэтому вы должны подумать о том, как разбить задачу на подзадачи. Это поможет вам определиться с осями.

Рассмотрим еще один пример. Допустим, вы открыли сайт *dictionary.com*. Пользователь вводит слово, а сайт возвращает определение. Но если пользователь ввел несуществующее слово, нужно предположить, какое слово имелось в виду. Алекс ищет определение «fish», но он случайно ввел «hish». Такого слова в словаре нет, но зато у вас есть список похожих слов.



СЛОВА, ПОХОЖИЕ НА «HISH»:

- FISH
- VISTA

(Это несерьезный пример, поэтому список ограничен всего двумя словами. Вероятно, на практике такой список будет состоять из тысяч слов.)

Итак, Алекс ввел строку *hish*. Какое слово он хотел ввести на самом деле: *fish* или *vista*?

Построение таблицы

Как должна выглядеть таблица для этой задачи? Вы должны ответить на следующие вопросы.

- Какие значения должны содержаться в ячейках?
- Как разбить эту задачу на подзадачи?
- Каков смысл осей таблицы?

В динамическом программировании вы пытаетесь максимизировать некоторую характеристику. В данном случае ищется самая длинная подстрока, общая в двух словах. Какую общую подстроку содержат *hish* и *fish*? А как насчет *hish* и *vista*? Именно это требуется вычислить.

Как говорилось ранее, значения в ячейках обычно представляют ту характеристику, которую вы пытаетесь оптимизировать. Вероятно, в данном случае этой характеристикой будет число: длина самой длинной подстроки, общей для двух строк.

Как разделить эту задачу на подзадачи? Например, можно заняться сравнением подстрок. Вместо того чтобы сравнивать *hish* и *fish*, можно сначала сравнить *his* и *fis*. Каждая ячейка будет содержать длину самой длинной подстроки, общей для двух подстрок. Такое решение также подсказывает, что строками и столбцами таблицы, вероятно, будут два слова. А значит, таблица будет выглядеть примерно так:

	h	i	s	h
f				
i				
s				
h				

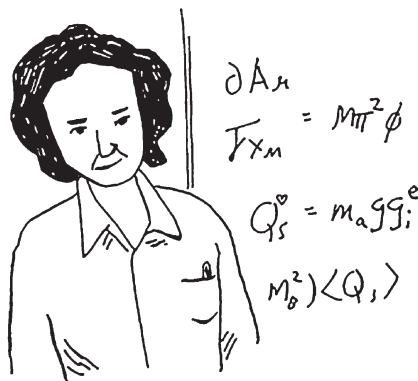
Если у вас голова идет кругом, не огорчайтесь. Это сложный материал — собственно, именно поэтому я объясняю его в конце книги! Ниже будет приведено упражнение, чтобы вы могли самостоятельно потренироваться в динамическом программировании.

Заполнение таблицы

Сейчас вы уже достаточно хорошо представляете, как должна выглядеть таблица. По какой формуле заполняются ячейки таблицы? Мы можем немного упростить свою задачу, потому что уже знаем решение — у *hish* и *fish* имеется общая подстрока длины 3: *ish*.

Однако этот факт ничего не говорит о том, какая формула должна использоваться. Программисты иногда шутят об использовании алгоритма Фейнмана. *Алгоритм Фейнмана*, названный по имени известного физика Ричарда Фейнмана, работает так:

1. Записать формулировку задачи.
2. Хорошенько подумать.
3. Записать решение.



Да, программисты — большие шутники!

По правде говоря, простого способа вычислить формулу для данного случая не существует. Вам придется экспериментировать и искать работоспособ-

ное решение. Иногда алгоритм предоставляет не точный рецепт, а основу, на которую вы наращиваете свою идею.

Попробуйте предложить решение этой задачи самостоятельно. Даю подсказку — часть таблицы выглядит так:

	Н	І	Ѕ	Н
F	0	0		
І				
Ѕ			2	0
Н				3

Чему равны другие значения? Вспомните, что каждая ячейка содержит значение *подзадачи*. Почему ячейка (3, 3) содержит значение 2? Почему ячейка (3, 4) содержит значение 0?

Попробуйте вывести формулу самостоятельно, прежде чем продолжить читать. Даже если вам не удастся получить правильный ответ, мои объяснения покажутся вам намного более понятными.

Решение

Итоговая версия таблицы выглядит так:

	Н	І	Ѕ	Н
F	0	0	0	0
І	0	1	0	0
Ѕ	0	0	2	0
Н	0	0	0	3

А это моя формула для заполнения ячеек:

1. ЕСЛИ БУКВЫ
НЕ СОВПАДАЮТ,
ЗНАЧЕНИЕ
РАВНО 0

	Н	І	Ѕ	Н
Н	0	0	0	0
І	0	1	0	0
Ѕ	0	0	2	0
Н	0	0	0	3

2. ЕСЛИ ОНИ СОВПАДАЮТ,
ТО ЗНАЧЕНИЕ РАВНО
ЗНАЧЕНИЮ ЯЧЕЙКИ
НАВЕРХУ СЛЕВА + 1

На псевдокоде эта формула реализуется так:

```

if word_a[i] == word_b[j]:  <----- Буквы совпадают
    cell[i][j] = cell[i-1][j-1] + 1
else:                       <----- Буквы не совпадают
    cell[i][j] = 0
  
```

Аналогичная таблица для строк *hish* и *vista*:

	В	І	Ѕ	Т	А
Н	0	0	0	0	0
І	0	1	0	0	0
Ѕ	0	0	2	0	0
Н	0	0	0	0	0

ОКОНЧАТЕЛЬНЫЙ ОТВЕТ ↑ НЕОКОНЧАТЕЛЬНЫЙ ОТВЕТ

Важный момент: в этой задаче окончательное решение далеко не всегда находится в последней ячейке! В задаче о рюкзаке последняя ячейка всегда содержит окончательное решение. Но в задаче поиска самой длинной общей подстроки решение определяется самым большим числом в таблице — и это может быть не последняя, а какая-то другая ячейка.

Вернемся к исходному вопросу: какая строка ближе к *hish*? У строк *hish* и *fish* есть общая подстрока длиной в три буквы. У *hish* и *vista* есть общая подстрока из двух букв. Скорее всего, Алекс хотел ввести строку *fish*.

Самая длинная общая подпоследовательность

Предположим, Алекс ввел строку *fosh*. Какое слово он имел в виду: *fish* или *fort*?

Сравним строки по формуле самой длинной общей подстроки.

	F	O	S	H
F	1	0	0	0
O	0	2	0	0
R	0	0	0	0
T	0	0	0	0

или

	F	O	S	H
F	1	0	0	0
I	0	0	0	0
S	0	0	1	0
H	0	0	0	2

Длина подстрок одинакова: две буквы! Но *fosh* при этом ближе к *fish*:

$$\begin{array}{cccc}
 F & O & S & H \\
 \downarrow & \downarrow & \downarrow & \\
 F & I & S & H
 \end{array} = 3$$

$$\begin{array}{cccc}
 F & O & S & H \\
 \downarrow & \downarrow & & \\
 F & O & R & T
 \end{array} = 2$$

Мы сравниваем самую длинную общую *подстроку*, а на самом деле нужно сравнивать самую длинную общую *подпоследовательность*: количество букв в последовательности, общих для двух слов. Как вычислить самую длинную общую подпоследовательность?

Ниже приведена частично заполненная таблица для *fish* и *fosh*.

	F	O	S	H
F	1	1		
I	1			
S			1	2
H				

Сможете ли вы определить формулу для этой таблицы? Самая длинная общая подпоследовательность имеет много общего с самой длинной общей подстрокой, и их формулы тоже очень похожи. Попробуйте решить задачу самостоятельно, а я приведу ответ ниже.

Самая длинная общая подпоследовательность — решение

Окончательная версия таблицы:

	F	O	S	H
F	1 → 1 → 1 → 1			
O	↓ 1	2 → 2 → 2		
R	↓ 1	↓ 2	2 → 2 → 2	
T	↓ 1	↓ 2	↓ 2	2

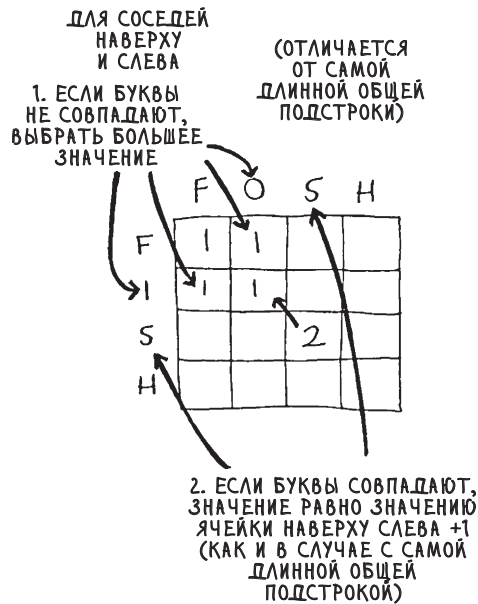
vs

	F	O	S	H
F	1 → 1 → 1 → 1			
I	↓ 1	↓ 1	↓ 1	1
S	↓ 1	↓ 1	2 → 2	
H	↓ 1	↓ 1	↓ 2	3

САМАЯ ДЛИННАЯ ОБЩАЯ
ПОДПОСЛЕДОВА-
ТЕЛЬНОСТЬ = 2

САМАЯ ДЛИННАЯ ОБЩАЯ
ПОДПОСЛЕДОВА-
ТЕЛЬНОСТЬ = 3

А теперь моя формула для заполнения каждой ячейки:



На псевдокоде эта формула реализуется так:

```

if word_a[i] == word_b[j]:  <..... Буквы совпадают
    cell[i][j] = cell[i-1][j-1] + 1
else:  <..... Буквы не совпадают
    cell[i][j] = max(cell[i-1][j], cell[i][j-1])
    
```

Поздравляю — вы справились! Безусловно, это была одна из самых сложных глав в книге. Находит ли динамическое программирование практическое применение? Да, находит.

- Биологи используют самую длинную общую подпоследовательность для выявления сходства в цепях ДНК. По этой метрике можно судить о сходстве двух видов животных, двух заболеваний и т. д. Самая длинная общая подпоследовательность используется для поиска лекарства от рассеянного склероза.

- Вы когда-нибудь пользовались ключом `diff` (например, в команде `git diff`)? Этот ключ выводит информацию о различиях между двумя файлами, а для этого он использует динамическое программирование.
- Мы также упоминали о сходстве строк. *Расстояние Левенштейна* оценивает, насколько похожи две строки, а для его вычисления применяется динамическое программирование. Расстояние Левенштейна используется в самых разных областях, от проверки орфографии до выявления отправки пользователем данных, защищенных авторским правом.

УПРАЖНЕНИЕ

11.3 Нарисуйте и заполните таблицу для вычисления самой длинной общей подстроки между строками *blue* и *clues*.

Шпаргалка

- Динамическое программирование применяется при оптимизации некоторой характеристики.
- Динамическое программирование работает только в ситуациях, в которых задача может быть разбита на автономные подзадачи.
- В каждом решении из области динамического программирования строится таблица.
- Значения ячеек таблицы обычно соответствуют оптимизируемой характеристике.
- Каждая ячейка представляет подзадачу, поэтому вы должны подумать о том, как разбить задачу на подзадачи.
- Не существует единой формулы для вычисления решений методом динамического программирования.

12 Алгоритм k ближайших соседей



В этой главе

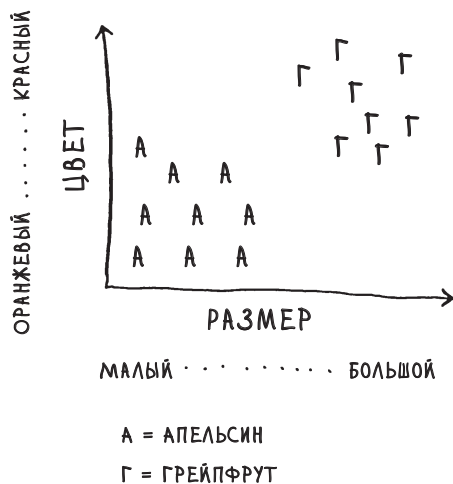
- ✓ Вы научитесь строить системы классификации на базе алгоритма k ближайших соседей.
- ✓ Вы узнаете об извлечении признаков.
- ✓ Вы узнаете о регрессии: прогнозировании чисел (например, завтрашних биржевых котировок или успеха фильма у зрителей).
- ✓ Вы познакомитесь с типичными сценариями использования и ограничениями алгоритма k ближайших соседей.

Апельсины и грейпфруты

Взгляните на этот фрукт. Что это, апельсин или грейпфрут? Я слышал, что грейпфруты обычно крупнее, а их кожура имеет красноватый оттенок.



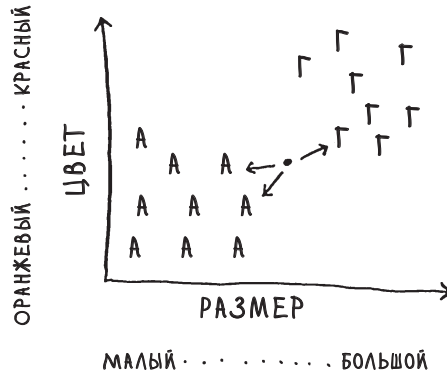
Мой мыслительный процесс выглядит примерно так: у меня в мозге существует некое подобие графика.



Как правило, крупные и красные фрукты оказываются грейпфрутами. Этот фрукт большой и красный, поэтому, скорее всего, это грейпфрут. Но что, если вам попадется фрукт вроде такого?



Как классифицировать этот фрукт? Один из способов — рассмотреть соседей этой точки. Возьмем ее трех ближайших соседей.



Среди соседей больше апельсинов, чем грейпфрутов. Следовательно, этот фрукт, скорее всего, является апельсином. Поздравляем: вы только что применили алгоритм *k* ближайших соседей для классификации! В целом алгоритм работает по довольно простому принципу.

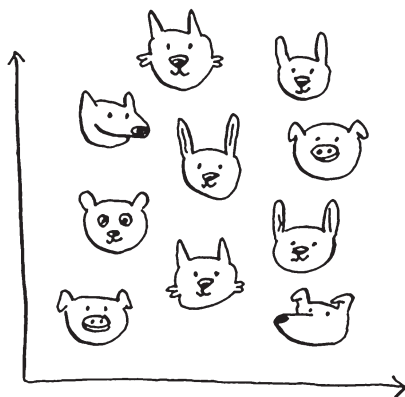


Алгоритм *k* ближайших соседей прост, но полезен! Если вы пытаетесь выполнить классификацию чего-либо, сначала попробуйте применить алгоритм *k* ближайших соседей. Рассмотрим более реалистичный пример.

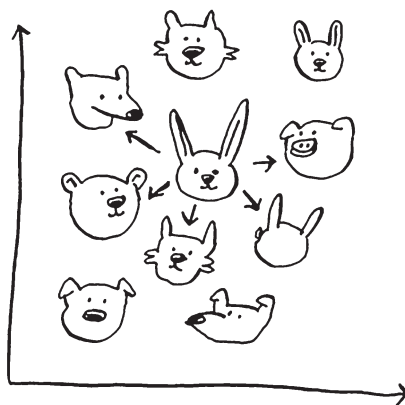
Построение рекомендательной системы

Представьте, что вы работаете на сайте Netflix и хотите построить систему, которая будет рекомендовать фильмы для ваших пользователей. На высоком уровне эта задача похожа на задачу с грейпфрутами!

Информация о каждом пользователе наносится на график.

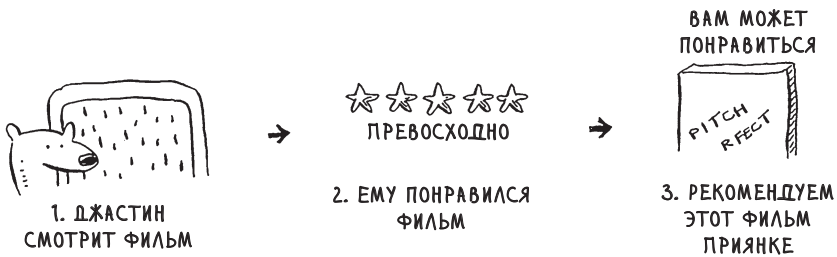


Положение пользователя определяется его вкусами, поэтому пользователи с похожими вкусами располагаются недалеко друг от друга. Предположим, вы хотите порекомендовать фильмы Приянке. Найдите 5 пользователей, ближайших к ней.



У Джастина, Джей-Си, Джозефа, Ланса и Криса похожие вкусы. Значит, те фильмы, которые нравятся *им*, с большой вероятностью понравятся и Приянке!

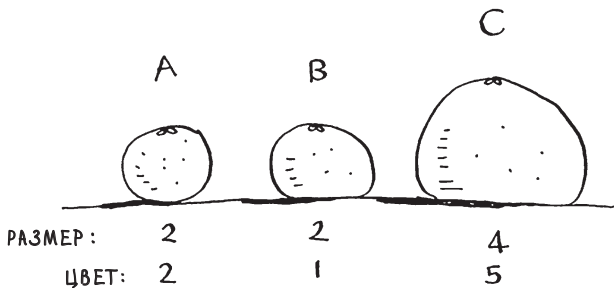
После того как у вас появится такая диаграмма, построить рекомендательную систему будет несложно. Если Джастину нравится какой-нибудь фильм, порекомендуйте этот фильм Приянке.



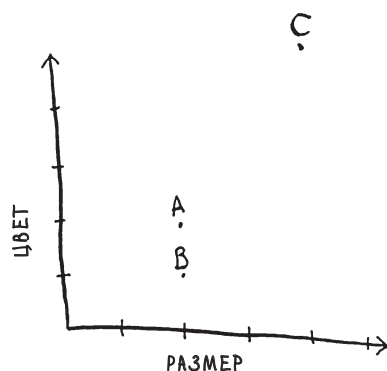
Однако в картине не хватает одного важного фрагмента. Вы сравнивали вкусы двух пользователей. Но как определить, насколько они близки?

Извлечение признаков

В примере с грейпфрутами мы сравнивали фрукты на основании их размера и цвета кожуры. Размер и цвет — *признаки*, по которым ведется сравнение. Теперь предположим, что у вас есть три фрукта. Вы можете извлечь из них информацию, то есть провести извлечение признаков.



Данные трех фруктов наносятся на график.



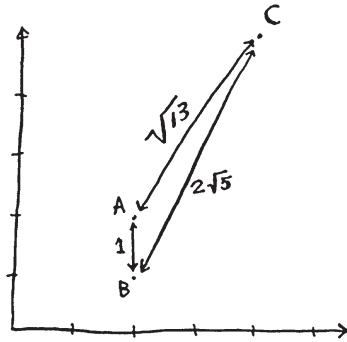
Из диаграммы хорошо видно, что фрукты А и В похожи. Давайте измерим степень их сходства. Для вычисления расстояния между двумя точками применяется формула Пифагора.

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Например, расстояние между А и В вычисляется так:

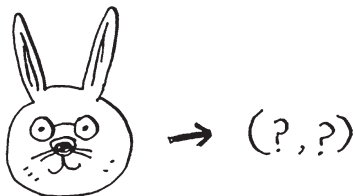
$$\begin{aligned} & \sqrt{(2-2)^2 + (2-1)^2} \\ &= \sqrt{0 + 1} \\ &= \sqrt{1} \\ &= 1 \end{aligned}$$

Расстояние между A и B равно 1. Другие расстояния вычисляются аналогично.



Формула расстояния подтверждает то, что мы видим: между фруктами A и B есть сходство.

Допустим, вместо фруктов вы сравниваете пользователей Netflix. Пользователей нужно будет как-то нанести на график. Следовательно, каждого пользователя нужно будет преобразовать в координаты — так же, как это было сделано для фруктов.



Когда вы сможете нанести пользователей на график, вы также сможете измерить расстояние между ними.

Начнем с преобразования пользователей в набор чисел. Когда пользователь регистрируется на Netflix, предложите ему оценить несколько категорий

фильмов: нравятся они лично ему или нет. Таким образом у вас появляется набор оценок для каждого пользователя!

			
	ПРИЯНКА	ДЖАСТИН	МОРФЕУС
КОМЕДИЯ	3	4	2
БОЕВИК	4	3	5
ДРАМА	4	5	1
УЖАСЫ	1	1	3
МЕЛОДРАМА	4	5	1

Приянка и Джастин обожают мелодрамы и терпеть не могут ужасы. Морфеусу нравятся боевики, но он не любит мелодрамы (хороший боевик не должен прерываться слащавой романтической сценой). Помните, как в задаче об апельсинах и грейпфрутах каждый фрукт представлялся двумя числами? Здесь каждый пользователь представляется набором из пяти чисел.

$$\text{Апельсин} \rightarrow (2, 2)$$

$$\text{Приянка} \rightarrow (3, 4, 4, 1, 4)$$

Математик скажет, что вместо вычисления расстояния в двух измерениях вы теперь вычисляете расстояние в пяти измерениях. Тем не менее формула расстояния остается неизменной.

$$\sqrt{(a_1 - a_2)^2 + (b_1 - b_2)^2 + (c_1 - c_2)^2 + (d_1 - d_2)^2 + (e_1 - e_2)^2}$$

Просто на этот раз используется набор из пяти чисел вместо двух.

Формула расстояния универсальна: даже если вы используете набор из миллиона чисел, расстояние вычисляется по той же формуле. Естественно спросить: какой смысл передает метрика *расстояния* с пятью числами? Она сообщает, насколько близки между собой эти наборы из пяти чисел.

$$\begin{aligned}
 & \sqrt{(3-4)^2 + (4-3)^2 + (4-5)^2 + (1-1)^2 + (4-5)^2} \\
 &= \sqrt{1 + 1 + 1 + 0 + 1} \\
 &= \sqrt{4} \\
 &= 2
 \end{aligned}$$

Это расстояние между Приянкой и Джастином.

ПРИМЕЧАНИЕ

Немного терминологии, которая вам пригодится. Массивы чисел, такие как (2, 2) для грейпфрута или (3, 4, 4, 1, 4) для вкусов Приянки в отношении кино, называются векторами. Если вам встретится статья о машинном обучении и вы увидите, что авторы говорят о векторах, знайте, что они имеют в виду подобный массив чисел.

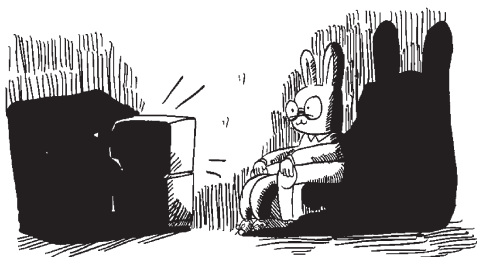
Вкусы Приянки и Джастина похожи. А насколько различаются вкусы Приянки и Морфеуса? Вычислите расстояние между ними, прежде чем продолжить чтение.

Сколько у вас получилось? Приянка и Морфеус находятся на расстоянии 24. По этому расстоянию можно понять, что у Приянки больше общего с Джастином, чем с Морфеусом.

Прекрасно! Теперь порекомендовать фильм Приянке будет несложно: если Джастину понравился какой-то фильм, мы рекомендуем его Приянке, и наоборот. Вы только что построили систему, рекомендующую фильмы.

Если вы являетесь пользователем Netflix, то он постоянно напоминает вам: «Пожалуйста, оценивайте больше фильмов. Чем больше фильмов вы оце-

ните, тем точнее будут наши рекомендации». Теперь вы знаете почему: чем больше фильмов вы оцениваете, тем точнее Netflix определяет, с какими пользователями у вас общие вкусы.



УПРАЖНЕНИЯ

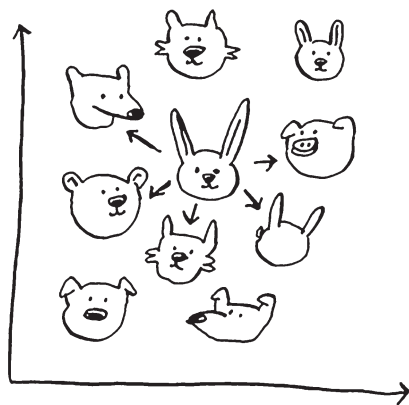
12.1 В примере с Netflix сходство между двумя пользователями оценивалось по формуле расстояния. Но не все пользователи оценивают фильмы одинаково. Допустим, есть два пользователя, Йоги и Пинки, вкусы которых совпадают. Но Йоги ставит 5 баллов любому фильму, который ему понравился, а Пинки более разборчива и ставит пятерки только самым лучшим фильмам. Вроде бы вкусы одинаковые, но по метрике расстояния они не являются соседями. Как учесть различия в стратегиях выставления оценок?



12.2 Предположим, Netflix определяет группу «авторитетов». Скажем, Квентин Тарантино и Уэс Андерсон относятся к числу авторитетов Netflix, поэтому их оценки оказывают более сильное влияние, чем оценки рядовых пользователей. Как изменить систему рекомендаций, чтобы она учитывала повышенную ценность оценок авторитетов?

Регрессия

А теперь предположим, что просто порекомендовать фильм недостаточно: вы хотите спрогнозировать, какую оценку Приянка поставит фильму. Возьмите 5 пользователей, находящихся вблизи от нее.



Кстати, я уже не в первый раз говорю о «ближайших пяти». В числе «5» нет ничего особенного: с таким же успехом можно взять 2 ближайших пользователей, 10 или 10 000. Поэтому-то алгоритм и называется «алгоритмом k ближайших пользователей», а не «алгоритмом 5 ближайших пользователей»!

Допустим, вы пытаетесь угадать оценку Приянки для фильма «Идеальный голос». Как этот фильм оценили Джастин, Джей-Си, Джозеф, Ланс и Крис?

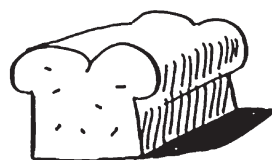
ДЖАСТИН : 5
 ДЖЕЙ-СИ : 4
 ДЖОЗЕФ : 4
 ЛАНС : 5
 КРИС : 3

Если вычислить среднее арифметическое их оценок, вы получите 4,2. Такой метод прогнозирования называется *регрессией*. У алгоритма k ближайших соседей есть два основных применения — классификация и регрессия:

- классификация = распределение по категориям;
- регрессия = прогнозирование ответа (в числовом выражении).

Регрессия чрезвычайно полезна. Представьте, что вы открыли маленькую булочную в Беркли и каждый день выпекаете свежий хлеб. Вы пытаетесь предсказать, сколько буханок следует испечь на сегодня. Есть несколько признаков:

- погода по шкале от 1 до 5 (1 = плохая, 5 = отличная);
- праздник или выходной (1, если сегодня праздник или выходной, 0 в противном случае);
- проходят ли сегодня спортивные игры (1 = да, 0 = нет).



И вы знаете, сколько буханок хлеба было продано в прошлом при разных сочетаниях признаков.

$$\boxed{\text{A.}} (5, 1, \emptyset) = \underset{\text{БУХАНОК}}{300} \quad \boxed{\text{B.}} (3, 1, 1) = \underset{\text{БУХАНОК}}{225}$$

$$\boxed{\text{C.}} (1, 1, \emptyset) = \underset{\text{БУХАНОК}}{75} \quad \boxed{\text{D.}} (4, \emptyset, 1) = \underset{\text{БУХАНОК}}{200}$$

$$\boxed{\text{E.}} (4, \emptyset, \emptyset) = \underset{\text{БУХАНОК}}{150} \quad \boxed{\text{F.}} (2, \emptyset, \emptyset) = \underset{\text{БУХАНОК}}{50}$$

Сегодня выходной и хорошая погода. Сколько буханок вы продадите на основании только что приведенных данных? Используем алгоритм k ближайших соседей для $k = 4$. Сначала определим четырех ближайших соседей для этой точки.

$$(4, 1, \emptyset) = ?$$

Ниже перечислены расстояния. Точки A, B, D и E являются ближайшими.

- A. 1 ←
- B. 2 ←
- C. 9
- D. 2 ←
- E. 1 ←
- F. 5

Вычисляя среднее арифметическое продаж в эти дни, вы получаете 218,75. Значит, именно столько буханок нужно выпекать на сегодня!

БЛИЗОСТЬ КОСИНУСОВ

До сих пор мы использовали формулу расстояния для вычисления степени сходства двух пользователей. Но является ли эта формула лучшей? На практике также часто применяется метрика близости косинусов. Допустим, два пользователя похожи, но один из них более консервативен в своих оценках. Обоим пользователям понравился фильм Манмохана Десаи «Амар Акбар Антони». Пол поставил фильму оценку 5 звезд, но Роуэн оценил его только в 4 звезды. Если использовать формулу расстояния, эти два пользователя могут не оказаться соседями, несмотря на сходство вкусов.

Метрика близости косинусов не измеряет расстояния между двумя векторами. Вместо этого она сравнивает углы двух векторов и в целом лучше подходит для подобных случаев. Тема метрики близости косинусов выходит за рамки этой книги, вам стоит самостоятельно поискать информацию, если вы будете применять алгоритм k ближайших соседей!

Выбор признаков

Чтобы подобрать рекомендации, вы предлагаете пользователям ставить оценки категориям фильмов. А если бы вы вместо этого предлагали им ставить оценки картинкам с котиками? Наверное, вам удалось бы найти пользователей, которые ставили похожие оценки этим картинкам. Однако у вас получилась бы самая плохая рекомендательная система в мире, потому что эти «признаки» не имеют никакого отношения к их вкусам в области кино!



Или представьте, что вы предлагаете пользователям оценить фильмы для формирования рекомендаций — но только «Историю игрушек», «Историю игрушек-2» и «Историю игрушек-3». Эти оценки ничего не скажут вам о вкусах пользователей.

Когда вы работаете с алгоритмом k ближайших соседей, очень важно правильно выбрать признаки для сравнения. Под правильным выбором признаков следует понимать:

- признаки, напрямую связанные с фильмами, которые вы пытаетесь рекомендовать;
- признаки, не содержащие смещения (например, если предлагать пользователям оценивать только комедии, вы не получите никакой информации об их отношении к боевикам).

Как вы думаете, оценки хорошо подходят для рекомендации фильмов? Возможно, я поставил «Прослушке» более высокую оценку, чем «Охотникам за недвижимостью», но на самом деле я провел больше времени за просмотром «Охотников». Как улучшить рекомендательную систему Netflix?

Возвращаясь к примеру с пекарней: сможете ли вы придумать два хороших и два плохих признака, которые можно было бы выбрать для прогнозирования объема выпечки? Возможно, нужно выпечь побольше хлеба после рекламы в газете. Или увеличить объем производства по понедельникам.

В том, что касается выбора хороших признаков, не существует единственно правильного ответа. Тщательно продумайте все факторы, которые необходимо учесть при прогнозировании.

УПРАЖНЕНИЕ

12.3 У сервиса Netflix миллионы пользователей. В приведенном ранее примере рекомендательная система строилась для пяти ближайших соседей. Пять — это слишком мало? Слишком много?

Знакомство с машинным обучением

Мало того, что алгоритм k ближайших соседей полезен — он открывает путь в волшебный мир машинного обучения! Суть машинного обучения — сделать ваш компьютер более разумным. Вы уже видели один пример машинного обучения: построение рекомендательной системы. В этом разделе будут рассмотрены другие примеры.



OCR

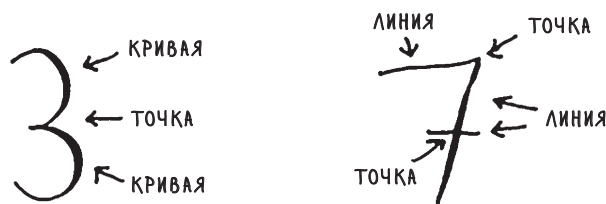
Сокращение OCR означает «Optical Character Recognition», то есть «оптическое распознавание текста». Иначе говоря, вы берете фотографию страницы текста, а компьютер автоматически преобразует изображение в текст. Google использует OCR для оцифровки книг. Как работает OCR? Для примера возьмем следующую цифру:

7

Как автоматически определить, что это за цифра? Можно воспользоваться алгоритмом k ближайших соседей:

1. Переберите изображения цифр и извлеките признаки.
2. Получив новое изображение, извлеките признаки и проверьте ближайших соседей.

По сути, это та же задача, что и задача классификации апельсинов и грейпфрутов. В общем случае алгоритмы OCR основаны на выделении линий, точек и кривых.



Затем при получении нового символа из него можно извлечь те же признаки.

Извлечение признаков в OCR происходит намного сложнее, чем в примере с фруктами. Однако важно понимать, что даже сложные технологии строятся на основе простых идей (таких, как алгоритм k ближайших соседей). Те же принципы могут использоваться для распознавания речи или лиц. Когда вы отправляете фотографию на Facebook, иногда сайту хватает сообразительности для автоматической пометки людей на фото. Да это машинное обучение в действии!

Первый шаг OCR, в ходе которого перебираются изображения цифр и происходит извлечение признаков, называется *тренировкой*. В большинстве алгоритмов машинного обучения присутствует фаза тренировки: прежде чем компьютер сможет решить свою задачу, его необходимо натренировать. В следующем примере рассматривается создание спам-фильтров, и в нем тоже есть шаг тренировки.

Построение спам-фильтра

Спам-фильтры используют другой простой алгоритм, называемый *наивным классификатором Байеса*. Сначала наивный классификатор Байеса тренируется на данных.

ТЕМА	СПАМ?
«ИЗМЕНИТЕ ПАРОЛЬ»	НЕ СПАМ
«ВЫ ВЫИГРАЛИ МИЛЛИОН»	СПАМ
«СООБЩИТЕ СВОЙ ПАРОЛЬ»	СПАМ
«НИГЕРИЙСКИЙ ПРИНЦ ГОТОВ ПЕРЕВЕСТИ ВАМ МИЛЛИОН»	СПАМ
«С ДНЕМ РОЖДЕНИЯ!»	НЕ СПАМ

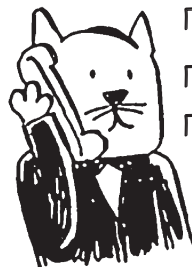
Предположим, вы получили сообщение с темой «Получите свой миллион прямо сейчас!» Это спам? Предложение можно разбить на слова, а затем для каждого слова проверить вероятность присутствия этого слова в спамовом сообщении. Например, в нашей очень простой модели слово «миллион» встречается только в спаме. Наивный классификатор Байеса вычисляет вероятность того, что сообщение с большой вероятностью является спамом. На практике он применяется примерно для тех же целей, что и алгоритм *k* ближайших соседей.

Например, наивный классификатор Байеса может использоваться для классификации фруктов: есть большой и красный фрукт. Какова вероятность того, что он окажется грейпфрутом? Это простой, но весьма эффективный алгоритм — из тех, что нам нравятся больше всего!

Прогнозы на биржевых торгах

Есть одна задача, в которой трудно добиться успеха машинным обучением: точно спрогнозировать курсы акций на бирже. Как выбрать хорошие признаки? Предположим, вы говорите, что если курс акций рос вчера, то он будет расти и сегодня. Хороший это признак или нет? Или, предположим, вы утверждаете, что курс всегда снижается в мае. Сработает или нет? Не

существует гарантированного способа прогнозировать будущее на основании прошлых данных. Прогнозирование будущего — сложное дело, а при таком количестве переменных оно становится почти невозможным.



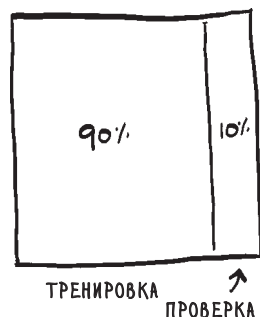
ПРОДАВАЙТЕ!
ПРОДАВАЙТЕ!
ПРОДАВАЙТЕ!

Тренировка модели МО: общий обзор

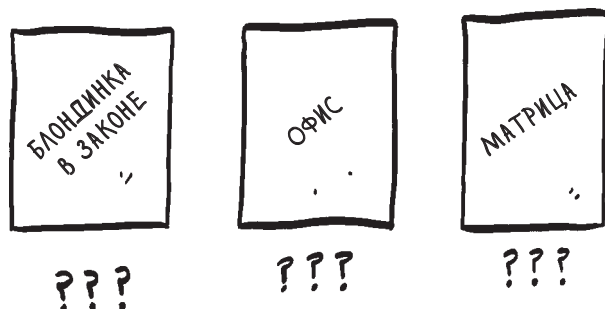


Мы разобрали несколько примеров; теперь рассмотрим последовательность действий по тренировке модели МО. Все начинается со сбора данных. В примере с Netflix данными были оценки пользователей. Затем необходимо провести очистку данных, то есть удаление неподходящих данных. Например, какой-нибудь пользователь не пожелал оценивать фильмы, поэтому расставил случайные оценки и перешел к следующему экрану. Такие данные необходимо удалить из набора. После этого требуется извлечь из данных признаки.

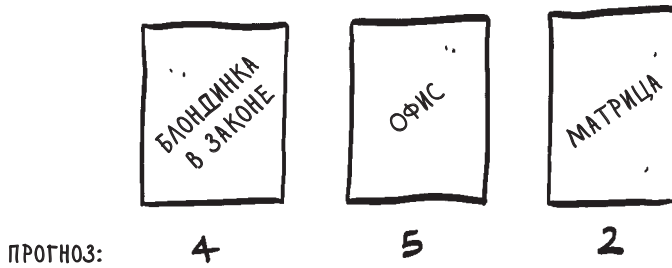
Получив признаки, можно переходить к тренировке модели. Выберите модель — k ближайших соседей, SVM, нейросеть — и проведите ее тренировку на 90 % данных. Оставьте 10 % для проверки модели. После того как модель пройдет обучение, протестируйте ее, предложив ей составить прогноз. Чтобы оценить качество прогноза, используйте зарезервированные 10 % данных.



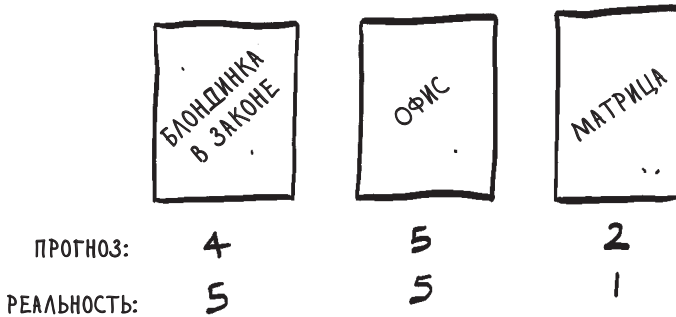
Например, вы хотите проверить рекомендательную модель Netflix. Можно спросить, понравились ли Приянке некоторые фильмы и сериалы:



Наша модель выдает прогнозы.



Мы знаем, какие фильмы нравятся Приянке, — эта информация содержится в 10 % данных, которые мы зарезервировали. Можно сравнить ее с прогнозами модели.



Неплохо! В таком случае можно сказать, что модель выдала хороший прогноз, потому что числа достаточно близки к реальным оценкам Приянки. Этот шаг тестирования модели называется проверкой (или валидацией) модели.

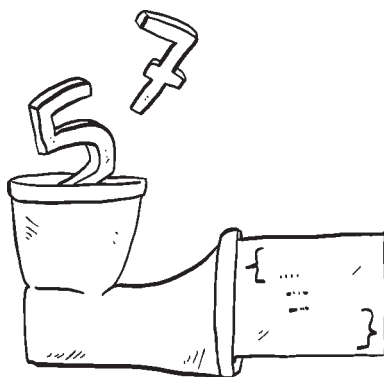
После проверки можно вернуться к модели и скорректировать ее. Представьте, что вы построили модель k ближайших соседей, в которой $k = 5$. Можно опробовать ее с $k = 7$, чтобы посмотреть, не даст ли она лучший результат. Это называется *подгонкой параметров*.

После того как вы завершите тренировку и оценку модели, ваша модель будет готова к работе. Так выглядит процедура построения модели МО в общем случае.

Шпаргалка

Надеюсь, вы хотя бы в общих чертах поняли, что можно сделать с помощью алгоритма k ближайших соседей и машинного обучения! Машинное обучение — интересная область, и при желании в нее можно зайти достаточно глубоко.

- Алгоритм k ближайших соседей применяется для классификации и регрессии. В нем используется проверка k ближайших соседей.
- Классификация = распределение по категориям.
- Регрессия = прогнозирование результата (например, в виде числа).
- «Извлечением признаков» называется преобразование элемента (например, фрукта или пользователя) в список чисел, которые могут использоваться для сравнения.
- Качественный выбор признаков — важная часть успешного алгоритма k ближайших соседей.



ИЗВЛЕЧЕНИЕ ПРИЗНАКОВ

13

Что дальше?



В этой главе

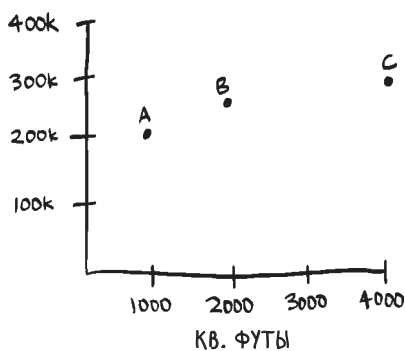
- ✓ Краткий обзор 10 алгоритмов, которые не рассматривались в книге. Вы узнаете, для чего нужны эти алгоритмы.
- ✓ Рекомендация книг, которые стоит читать дальше в зависимости от того, какие темы вам интересны.

Линейная регрессия

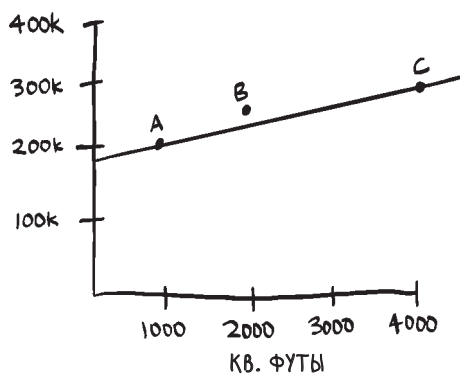
Представьте, что вы хотите продать свой дом. Его площадь составляет 3000 квадратных футов. Вы изучаете описания домов, недавно проданных по соседству.



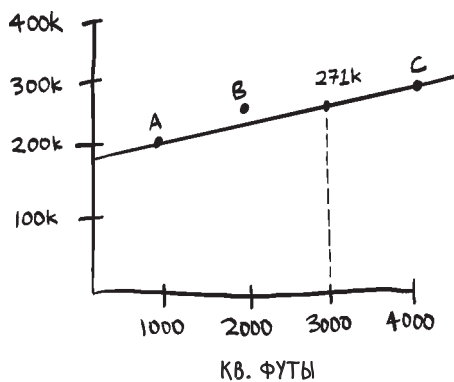
Какую цену вы бы назначили за свой дом, исходя из этой информации? Ниже представлено одно из возможных решений. Сначала вы наносите все точки на график.



Затем соединяете их линией.



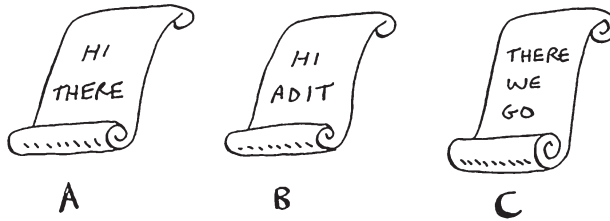
Теперь вы видите, где на этой линии находится точка 3000 квадратных футов. Соответствующую ей цену можно сделать начальной ценой.



Так работает линейная регрессия. Она выстраивает линию из заданной группы точек, а затем использует эту линию для вычисления прогнозов. Линейная регрессия уже давно применяется в статистике, а теперь и в машинном обучении, поскольку это метод, с которого легко начинать. Линейная регрессия полезна, если значения непрерывны. Если вам нужно сделать прогноз, проще всего использовать именно линейную регрессию.

Инвертированные индексы

Перед вами сильно упрощенное объяснение того, как работает поисковая система. Допустим, имеются три веб-страницы с простым содержимым.



Построим хеш-таблицу для этого содержимого.

Ключами хеш-таблицы являются слова, а значения указывают, на каких страницах встречается каждое слово. Теперь предположим, что пользователь ищет слово *hi*. Посмотрим, на каких страницах это слово встречается.

HI	A, B
THERE	A, C
ADIT	B
WE	C
GO	C

HI	A, B
----	------

Ага, слово встречается на страницах A и B. Выведем эти страницы в результатах поиска. Или предположим, что пользователь ищет слово *there*. Вы знаете, что это слово встречается на страницах A и C. Несложно, верно?

Это очень полезная структура данных: хеш-таблица, связывающая слова с местами, в которых эти слова встречаются. Такая структура данных, называемая *инвертированным индексом*, часто используется для построения поисковых систем. Если вас интересует область поиска, эта тема станет хорошей отправной точкой для дальнейшего изучения.

Преобразование Фурье

Преобразование Фурье — действительно выдающийся алгоритм: великолепный, элегантный и имеющий миллион практических применений. Лучшая аналогия для преобразования Фурье приводится на сайте Better Explained (отличный веб-сайт, на котором просто объясняется математическая теория): если у вас есть коктейль, преобразование Фурье сообщает, из каких ингредиентов он состоит¹. Или для заданной песни преобразование разделяет ее на отдельные частоты.

Оказывается, эта простая идея находит множество практических применений. Например, если песню можно разложить на частоты, вы можете усилить тот диапазон, который вас интересует, — скажем, усилить низкие частоты и приглушить высокие. Преобразование Фурье прекрасно подходит для обработки сигналов. Также оно может применяться для сжатия музыки: сначала звуковой файл разбивается на составляющие. Преобразование Фурье сообщает, какой вклад вносит каждая составляющая в музыку, что позволяет исключить несущественные составляющие. Собственно, именно так работает музыкальный формат MP3!

Музыка — не единственный вид цифровых сигналов. Графический формат JPG также использует сжатие и работает по тому же принципу. Кроме того, преобразование Фурье применяется для прогнозирования землетрясений и анализа ДНК.

С его помощью можно построить аналог Shazam — приложение, которое находит песни по отрывкам. Преобразование Фурье очень часто применяется на практике. Почти наверняка вы с ним еще столкнетесь!

¹ Kalid, «An Interactive Guide to the Fourier Transform,» Better Explained, <http://mng.bz/dd9N>.

Параллельные алгоритмы

Следующие три темы связаны с масштабируемостью и обработкой больших объемов данных. Когда-то компьютеры становились все быстрее и быстрее. Если вы хотели, чтобы ваш алгоритм работал быстрее, можно было подождать несколько месяцев и запустить программу на более мощном компьютере. Но сейчас этот период подошел к концу. Современные компьютеры и ноутбуки оснащаются многоядерными процессорами. Чтобы алгоритм заработал быстрее, необходимо преобразовать его в форму, подходящую для параллельного выполнения сразу на всех ядрах!

Рассмотрим простой пример. Лучшее время выполнения для алгоритма сортировки равно приблизительно $O(n \log n)$. Известно, что массив невозможно отсортировать за время $O(n)$, *если только не воспользоваться параллельным алгоритмом!* Существует параллельная версия быстрой сортировки, которая сортирует массив за время $O(n)$.

Параллельный алгоритм трудно разработать. И так же трудно убедиться в том, что он работает правильно, и понять, какой прирост скорости он обеспечивает. Одно можно заявить твердо: выигрыш по времени нелинеен. Следовательно, если процессор вашего компьютера имеет два ядра вместо одного, из этого не следует, что ваш алгоритм по волшебству заработает вдвое быстрее. Это объясняется несколькими причинами.

- *Затраты ресурсов на управление параллелизмом* — допустим, нужно отсортировать массив из 1000 элементов. Как разбить эту задачу для выполнения на двух ядрах? Выделить каждому ядру 500 элементов, а затем объединить два отсортированных массива в один большой отсортированный массив? Слияние двух массивов требует времени.
- *Закон Амдала* — представьте, что вы пишете картину. Это занимает очень много времени, допустим около 20 часов. В идеале вам хотелось бы уложиться в 10 часов. Вы решаете оптимизировать процесс. Для этого разбиваете его на два шага: (1) эскиз и (2) рисование красками. Исходный эскиз не обязательно рисовать вручную — конечно, с помощью трейсинга (обводки) это делается быстрее. Но при следующей попытке на картину уходит 19 часов и 5 минут! Что произошло? Раньше эскиз занимал 1 час. Вы сократили его до 5 минут, что стало значительным улучшением. Тем не менее большая часть времени приходится на рисование красками, так что оптимизация оказалась незначительной.

Это проявление так называемого закона Амдала. Он гласит, что при оптимизации одной части системы выигрыш в производительности ограничивается тем, какую долю времени занимает эта часть. В данном случае время создания эскиза сокращается до $1/12$ исходного значения. При этом экономятся 55 минут. Если бы нам удалось сократить время рисования красками в той же степени, это сэкономило бы 1045 минут! Ускоряя алгоритм за счет параллелизации, подумайте, какую его часть следует распараллелить — создание эскиза или рисование?

- *Распределение нагрузки* — допустим, необходимо выполнить 10 задач, и вы назначаете каждому ядру 5 задач. Однако ядру А достаются все простые задачи, поэтому оно выполняет свою работу за 10 секунд, тогда как ядро В справится со сложными задачами только за минуту. Это означает, что ядро А целых 50 секунд простаивает, пока ядро В выполняет всю работу! Как организовать равномерное распределение работы, чтобы оба ядра трудились с одинаковой интенсивностью?

Если вас интересует теоретическая сторона производительности и масштабируемости, возможно, параллельные алгоритмы — именно то, что вам нужно!

map/reduce

Существует разновидность параллельных алгоритмов, которая в последнее время активно набирает популярность: *распределенные алгоритмы*. Конечно, параллельный алгоритм удобно запускать на компьютере, если для его выполнения требуется от двух до четырех ядер. Но если нужны сотни? Тогда алгоритм записывают так, чтобы он мог выполняться на множестве машин. Компания Google популяризировала распределенный алгоритм, который называется MapReduce, но его функции известны уже давно.

Для чего нужны распределенные алгоритмы? Предположим, имеется таблица с миллиардами или триллионами записей и к ней нужно сделать сложный запрос SQL. В MySQL это сделать не удастся, потому что MySQL начнет «тормозить» уже после нескольких миллиардов записей. Используйте MapReduce!

Или, предположим, вам нужно обработать длинный список задач. Обработка каждой задачи занимает 10 секунд, а всего в списке 1 миллион задач.

Если выполнять эту работу на одном компьютере, она займет несколько месяцев! Если бы ее можно было выполнить на 100 машинах, работа завершилась бы за несколько дней.

Распределенные алгоритмы хорошо работают в ситуациях, когда нужно выполнить большой объем работы за максимально короткое время.

Фильтры Блума и HyperLogLog

Представьте себя на месте сайта Reddit. Когда пользователь публикует ссылку, нужно проверить, публиковалась ли эта ссылка ранее. Истории, которые еще не публиковались, считаются более ценными.

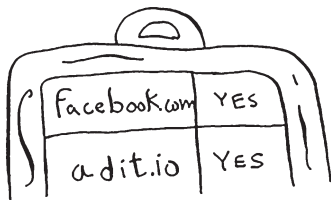
Или представьте себя на месте поискового бота Google. Обрабатывать веб-страницу нужно только в том случае, если она еще не обрабатывалась ранее. Итак, нужно проверить, обрабатывалась ли страница ранее.

Или представьте себя на месте *bit.ly* — сервиса сокращения URL. Пользователи не должны перенаправляться на вредоносные сайты. У вас имеется набор URL-адресов, которые считаются вредоносными. Теперь нужно выяснить, не направляется ли пользователь на URL-адрес из этого набора.

Во всех этих примерах возникает одна проблема. Имеется очень большой набор данных.



Появляется новый объект, и вы хотите узнать, содержится ли он в существующем наборе. Эта задача быстро решается при помощи хеша. Например, представьте, что Google создает большой хеш, ключами которого являются все обработанные страницы.



Как узнать, обрабатывался ли сайт *adit.io*? Нужно заглянуть в хеш.

adit.io → YES

У *adit.io* имеется свой ключ в хеше, а значит, адрес уже обрабатывался. Среднее время обращения к элементам в хеш-таблице составляет $O(1)$. Таким образом, вы узнали о том, что страница *adit.io* уже проиндексирована за постоянное время. Неплохо!

Вот только этот хеш получится просто *огромным*. Google индексирует триллионы веб-страниц. Если хеш содержит все URL-адреса, индексируемые Google, он займет слишком много места. У Reddit и *bit.ly* возникает аналогичная проблема. Сталкиваясь с такими объемами данных, приходится действовать более изобретательно!

Фильтры Блума

Для решения проблемы можно воспользоваться *вероятностными структурами данных*, которые называются *фильтрами Блума*. Они дают ответ, который может оказаться ложным, но с большой вероятностью является правильным. Вместо того чтобы обращаться к хешу, вы спрашиваете у фильтра Блума, обрабатывался ли этот URL-адрес ранее. Хеш-таблица

даст точный ответ. Фильтр Блума дает ответ, правильный с высокой вероятностью:

- возможны ложноположительные срабатывания. Фильтр скажет: «Этот сайт уже обрабатывался», хотя этого не было;
- ложноотрицательные срабатывания исключены. Если фильтр утверждает, что сайт не обрабатывался, вы можете быть в этом уверены.

Фильтры Блума хороши тем, что занимают очень мало места. Хеш-таблице пришлось бы хранить все URL-адреса, обрабатываемые Google, а фильтру Блума это не нужно. Фильтр Блума очень удобен тогда, когда нет необходимости в точном ответе (как во всех приведенных примерах). Например, про *bit.ly* он может сказать: «Мы полагаем, что сайт может оказаться вредоносным, будьте особенно внимательны».

HyperLogLog

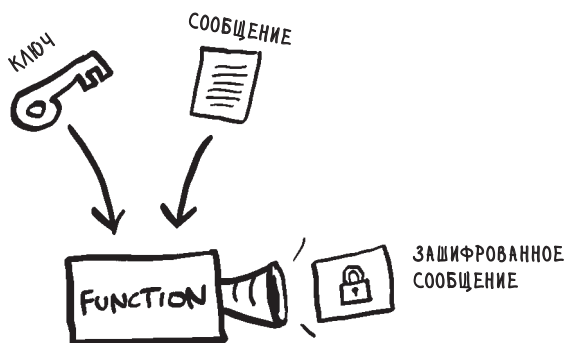
Примерно так же действует другой алгоритм, который называется HyperLogLog. Предположим, Google хочет подсчитать количество *уникальных* поисков, выполненных пользователями. Или Amazon хочет подсчитать количество уникальных предметов, просмотренных пользователями за сегодняшний день. Для получения ответов на эти вопросы потребуется очень много места! Так, в примере с Google придется вести журнал всех уникальных вариантов поиска. Когда пользователь что-то ищет, вы сначала проверяете, присутствует ли условие в журнале, и если нет, добавляете его. Даже для одного дня этот журнал получится гигантским.

HyperLogLog аппроксимирует количество уникальных элементов в множестве. Как и фильтры Блума, он не дает точного ответа, но выдает достаточно близкий результат с использованием малой части памяти, которую обычно занимает такая задача.

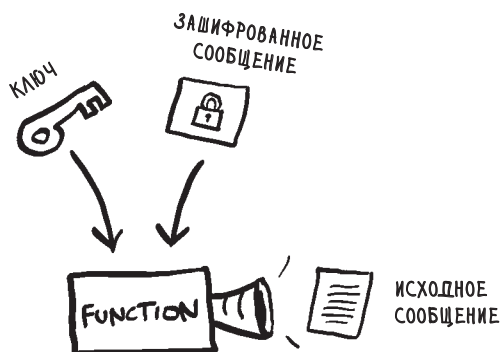
Если вы используете большие объемы данных и вас устраивают приближенные ответы — воспользуйтесь вероятностными алгоритмами!

HTTPS и обмен ключами Диффи — Хеллмана

HTTPS — опора современного интернета. Этот протокол обеспечивает безопасные онлайн-транзакции, от ввода пароля до оплаты покупок. Работа HTTPS основана на шифровании сообщений между клиентом и сервером. Как же работает шифрование? Функции передаются сообщение и секретный ключ. Затем функция генерирует зашифрованное сообщение.



Чтобы расшифровать его, передайте зашифрованное сообщение и *тот же* ключ функции, и вы получите исходное сообщение.



Когда вы отправляете данные серверу, ваш браузер зашифровывает сообщение, после чего сервер расшифровывает его. Просто, не так ли? Да, кроме

одного: как проверить, что браузер и сервер используют один и тот же ключ?



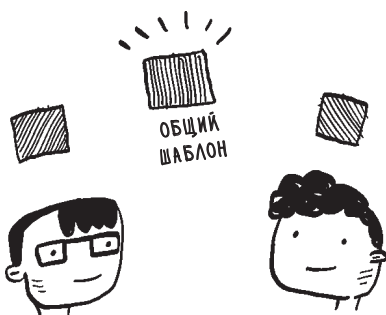
Помните, что для работы HTTPS обе стороны должны иметь одинаковые ключи. Но как задать ключ, чтобы никто его не видел? Если отправить ключ серверу, злоумышленник сможет перехватить этот ключ. Как задать ключ, чтобы только браузер и сервер знали его? Задача кажется неразрешимой, но это не так! Для ее решения существует очень умный алгоритм, называемый обменом ключей Диффи — Хеллмана. Ниже описано, как он работает.

На шаге 1 мы генерируем собственные ключи. Я — клиент, и я генерирую ключ для себя. Сервис также генерирует ключ. Эти ключи разные. Мы не знаем ключи друг друга, они закрытые.



На иллюстрации ключи обозначены схематично, чтобы наглядно показать, что происходит. В реальности ключи представляют собой байты.

На шаге 2 генерируется общий шаблон.



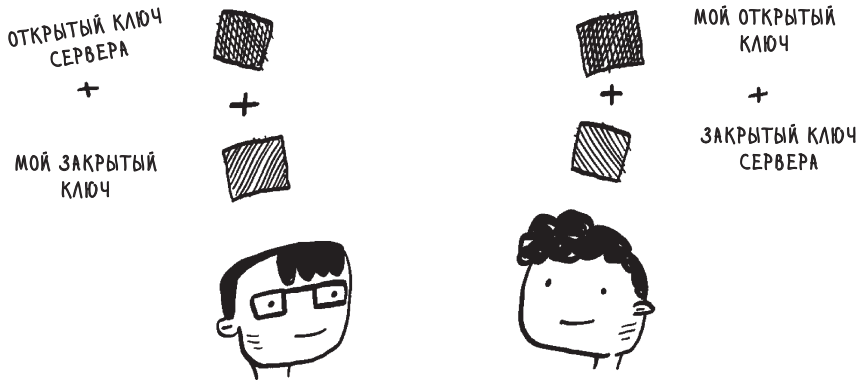
Это открытый шаблон. Он виден обеим сторонам и всем остальным. Нам неважно, кто его видит.

На шаге 3 этот шаблон совмещается с закрытыми ключами.



В результате мы получаем *открытые ключи*. Раз они открытые, нас не интересует, кто их видит. Сервер видит мой открытый ключ, а я вижу его открытый ключ.

Наконец, на шаге 4 я получаю открытый ключ сервера и совмещаю его со своим закрытым ключом. Сервис делает то же самое с моим открытым ключом.



Ура! Теперь у нас одинаковые ключи! У каждого из нас есть ключ, объединяющий три шаблона.



Нам удалось задать ключ, не отправляя ключи друг другу. Заданный ключ называется *общим секретом*; так работает обмен ключами Диффи — Хеллмана.



HTTPS — интересная и важная часть нашей повседневной жизни. Вот немного терминологии, связанной с HTTPS:

- *TLS* — протокол безопасности транспортного уровня (Transport Layer Security). TLS используется для установления безопасного соединения.
- *SSL* — прежнее название TLS, но люди часто не различают их. Если вы слышите, как кто-то упоминает SSL, то, скорее всего, речь идет о TLS. В этом протоколе часто находят дефекты безопасности, поэтому его постоянно приходится обновлять. Протокол TLS впервые появился в 1999 году. Все версии протокола SSL, выпущенные до TLS, небезопасны.
- *Шифрование с симметричным ключом* — в нашем примере обе стороны используют одинаковые ключи. Но существует и *асимметричное шифрование с открытым ключом*, когда обе стороны используют разные ключи. Я говорю о шифровании с симметричным ключом, потому что оно используется HTTPS.

HTTPS использует измененную версию обмена ключей Диффи — Хеллмана, называемую *эфемерным обменом ключей Диффи — Хеллмана*. Она работает так же, как было показано, кроме того, что закрытые ключи генерируются заново для каждого соединения. Это означает, что даже если злоумышленнику станет известен один из закрытых ключей, он сможет расшифровать сообщения только одного соединения.

Криптография — интересная и глубокая тема. Если вам захочется узнать больше, я рекомендую другую книгу издательства Manning: «Real-World Cryptography» Дэвида Вона (David Wong) (<https://www.manning.com/books/real-world-cryptography>).



Локально-чувствительное хеширование

Многие хеш-функции имеют одну важную особенность: они являются *локально-нечувствительными*. Предположим, имеется строка, для которой генерируется хеш-код с помощью алгоритма SHA-256:

dog → cd6357

Если изменить в строке всего один символ, а потом сгенерировать хеш заново, строка полностью изменяется!

dot → e392da

И это хорошо, потому что сравнение хешей не позволит атакующему определить, насколько он близок к взлому пароля.

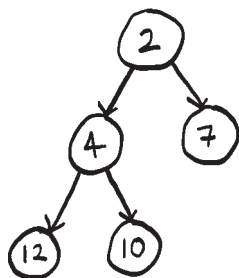
Иногда требуется обратный результат: локально-чувствительная функция хеширования. Здесь на помощь приходит алгоритм *Simhash*. При незначительном изменении строки Simhash генерирует хеш-код, который почти не отличается от исходного. Это позволяет сравнивать хеш-коды и определять, насколько похожи две строки, — весьма полезная возможность!

- Google использует Simhash для выявления дубликатов в процессе индексирования.
- Преподаватель может использовать Simhash для обнаружения плагиата (копирования рефератов из интернета).
- Scribd позволяет пользователям загружать документы или книги, чтобы они стали доступны для других пользователей. Но Scribd не хочет, чтобы пользователи размещали информацию, защищенную авторским правом! С помощью Simhash сайт может обнаружить, что отправленная информация похожа на книгу о Гарри Поттере, и при обнаружении сходства автоматически запретить ее размещение.

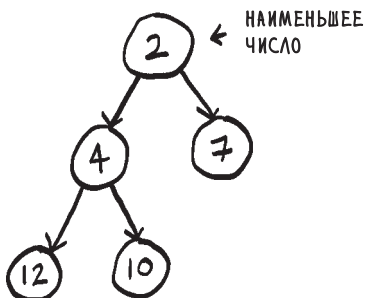
Simhash используется для выявления сходства между фрагментами текста.

Минимальные кучи и приоритетные очереди

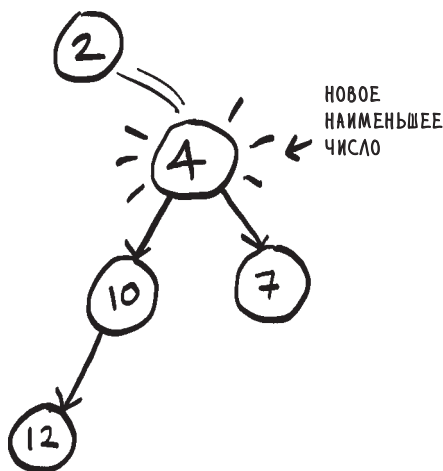
Минимальная куча — структура данных, построенная на базе деревьев. В таких кучах хранятся числа. Справа приведен пример минимальной кучи.



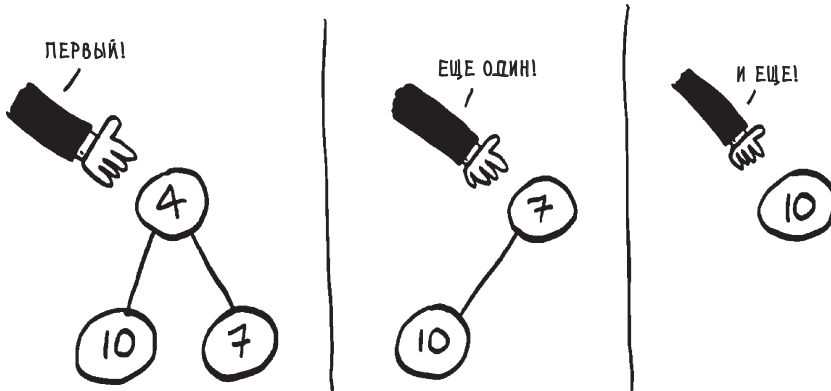
Минимальные кучи позволяют быстро найти наименьший элемент кучи, потому что наименьшее значение всегда находится в корне. Это самое важное свойство минимальной кучи, и наименьший элемент находится за время $O(1)$.



Или же за время $O(\log n)$ его можно удалить из кучи, заменив новым минимумом:

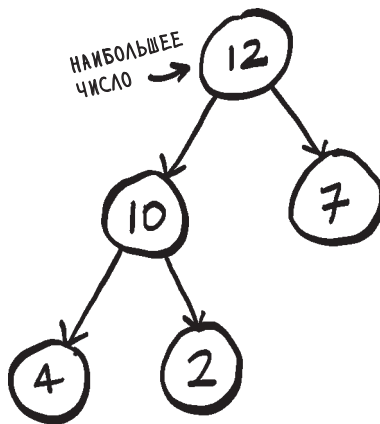


Благодаря кучам проводить сортировку очень легко — просто продолжайте запрашивать минимальное значение.



И сохраняйте значения по порядку. В конце дерево будет пустым, а у вас появится отсортированный список чисел! Такой алгоритм называется *пирамидальной сортировкой*, или *сортировкой кучей*.

Максимальные кучи (невозрастающие пирамиды) очень похожи на минимальные, но у них в корне хранится наибольшее значение.



Кучи отлично подходят для реализации приоритетных очередей. Мы обсуждали структуру данных очереди в главе 6. Помните, что очереди относятся к категории FIFO (First In, First Out — «первым пришел, первым

вышел»). С другой стороны, стек относится к структуре данных LIFO (Last In, First Out — «последним пришел, первым вышел»). Приоритетная очередь похожа на обычную очередь, не считая того, что при запросе элемента приоритетная очередь выдает элемент с наибольшим приоритетом. Приоритетная очередь отлично подходит для реализации списка задач. Сначала вы сохраняете в списке задачи, а затем запрашиваете задачу для работы; приоритетная очередь выдает задачу с наибольшим приоритетом. Приоритетные очереди также используются для реализации эффективной версии алгоритма Дейкстры.

Линейное программирование

Лучшее я приберег напоследок. Линейное программирование — одна из самых интересных областей, которые мне известны.

Линейное программирование используется для максимизации некоторой характеристики при заданных ограничениях. Предположим, ваша компания выпускает два продукта: рубашки и сумки. На рубашку требуется 1 метр ткани и 5 пуговиц. На изготовление сумки необходимо 2 метра ткани и 2 пуговицы. У вас есть 11 метров ткани и 20 пуговиц. Рубашка приносит прибыль \$2, а сумка — \$3. Сколько рубашек и сумок следует изготовить для получения максимальной прибыли?

Здесь мы пытаемся максимизировать прибыль, а ограничения определяют количество имеющихся материалов.

Другой пример: вы политик, пытающийся получить максимальное количество голосов. Исследования показали, что на каждый голос жителя Сан-Франциско требуется примерно 1 час работы (маркетинг, исследования и т. д.), а на каждый голос жителя Чикаго — 1,5 часа. Вам нужны голоса как минимум 500 жителей Сан-Франциско и как минимум 300 жителей Чикаго. В вашем распоряжении 50 дней. Кроме того, затраты на жителя Сан-Франциско составляют \$2, а на жителя Чикаго — \$1. Ваш бюджет составляет \$1500. Какое максимальное количество голосов вы сможете получить (Сан-Франциско+Чикаго)?

На этот раз вы стремитесь к максимуму голосов при ограничениях по времени и деньгам.

Возможно, вы думаете: «В этой книге много говорилось о вопросах оптимизации. Как они связаны с линейным программированием?» Все алгоритмы, работающие с графами, могут быть реализованы средствами линейного программирования. Линейное программирование — намного более общая область, а задачи с графами составляют ее подмножество.

В линейном программировании используется *симплекс-метод*. Этот алгоритм достаточно сложен, поэтому я не привожу его в книге. Если вы интересуетесь задачами оптимизации, поищите информацию о линейном программировании!

Эпилог

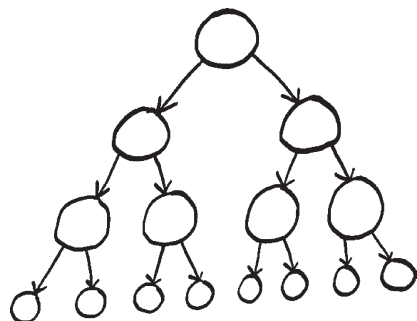
Надеюсь, этот краткий обзор показал, как много вам еще предстоит узнать. Я считаю, что лучший способ узнать что-то — найти тему, которая вас интересует, и изучить ее. Надеюсь, книга закладывает для этого достаточно надежную основу.

А Производительность АВЛ-деревьев

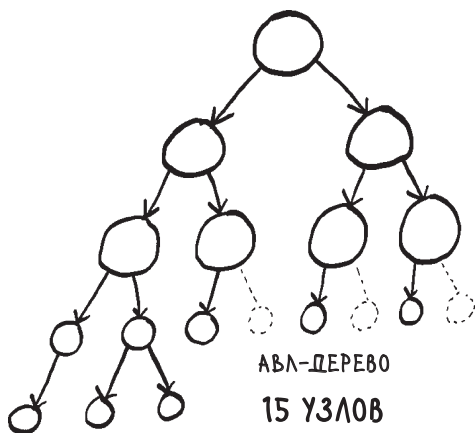


В этом приложении рассматривается производительность АВЛ-деревьев, речь о которых шла в главе 8. Прежде чем читать приложение, необходимо изучить эту главу.

Вспомните, что АВЛ-деревья обеспечивают производительность поиска $O(\log n)$. Однако не все так просто. На иллюстрации изображены два дерева. Оба обеспечивают производительность поиска $O(\log n)$, но имеют разную высоту!



ИДЕАЛЬНО СБАЛАНСИРОВАННОЕ ДЕРЕВО
15 УЗЛОВ
ВЫСОТА: 3



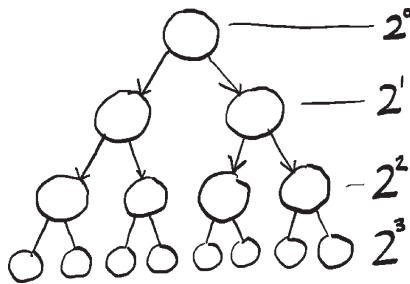
АВЛ-ДЕРЕВО
15 УЗЛОВ
ВЫСОТА: 4

(Пунктиром обозначены недостающие узлы в дереве.)

AVL-деревья допускают разность в единицу по высоте. Вот почему, хотя оба дерева состоят из 15 узлов, у идеально сбалансированного дерева высота равна 3, а у AVL-дерева — 4. Идеально сбалансированное дерево — это наше интуитивное представление «сбалансированного дерева», в котором каждый уровень заполнен узлами перед добавлением нового узла. Но AVL-дерево также считается сбалансированным, хотя в нем присутствуют пропуски — позиции, в которых мог бы находиться узел.

Вспомните, что в дереве производительность тесно связана с высотой. Как эти деревья при разной высоте обеспечивают одинаковую производительность? Что ж, мы нигде не обсуждали, по какому основанию вычисляется логарифм!

Идеально сбалансированное дерево имеет производительность $O(\log n)$, где \log обозначает логарифм по основанию 2, как при бинарном поиске. Это видно на иллюстрации. Каждый новый уровень удваивает количество узлов плюс 1. Таким образом, идеально сбалансированное дерево с высотой 1 состоит из 3 узлов, с высотой 2 — из 7 узлов ($3 \times 2 + 1$), с высотой 3 — из 15 узлов ($7 \times 2 + 1$) и т. д. Также можно считать, что каждый новый уровень добавляет количество узлов, равное степени 2.



Таким образом, идеально сбалансированное дерево имеет производительность $O(\log n)$, где \log обозначает логарифм по основанию 2.

В AVL-дереве имеются пропуски. В нем каждый новый уровень добавляет *меньше* узлов, чем удвоенное количество узлов предыдущего уровня. Как выясняется, AVL-дерево обеспечивает производительность $O(\log n)$, но в этом случае \log обозначает логарифм по основанию ϕ («золотое сечение», или $\sim 1,618$). Это небольшое, но интересное отличие — AVL-деревья обеспечивают производительность, несколько уступающую идеально сбалансированным деревьям, что объясняется другим основанием логарифма. Но в обоих случаях производительность остается очень близкой, так как она подчиняется оценке $O(\log n)$. Просто помните, что речь идет не о точном равенстве.

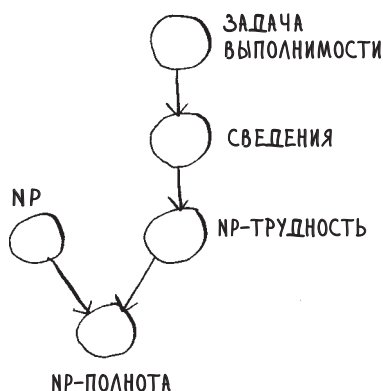
Б

NP-трудные задачи



У задачи о покрытии множества и задачи о коммивояжере есть нечто общее: их трудно решать. Чтобы найти наименьшее покрытие множества или кратчайший маршрут, приходится проверять все возможные варианты.

Обе задачи относятся к категории NP-трудных. Термины *NP*, *NP-трудность* и *NP-полнота* могут породить путаницу; конечно, когда-то они запутали и меня. В этом приложении я постараюсь объяснить, что означают все эти термины, но сначала необходимо представить ряд других концепций. Ниже перечислены понятия, которые вы узнаете, и их взаимосвязь:



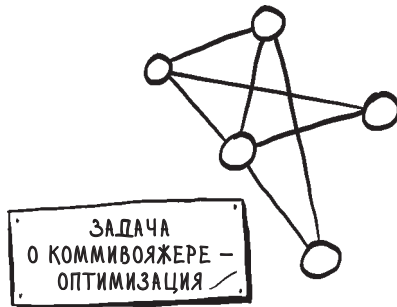
Но сначала необходимо объяснить, что такое *задача разрешимости*, потому что все задачи, которые мы рассмотрим в этом приложении, относятся к такой категории задач.

Задачи разрешимости

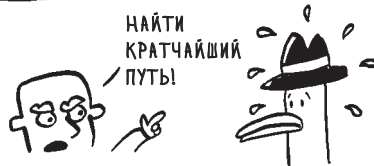
NP-полные задачи всегда являются задачами разрешимости. Задача разрешимости подразумевает ответ «да» или «нет». Задача о коммивояжере не относится к задачам разрешимости. Она требует найти кратчайший путь, а это задача оптимизации.

ПРИМЕЧАНИЕ

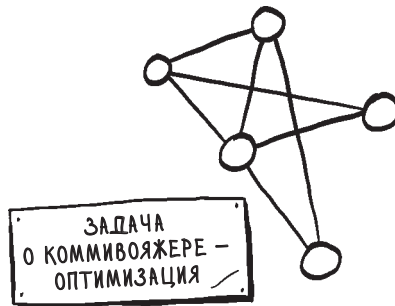
Я знаю, что во вводной части я говорил об NP-**трудных** задачах, а сейчас говорю о NP-**полных** задачах. Вскоре вы узнаете, чем отличаются эти два понятия.



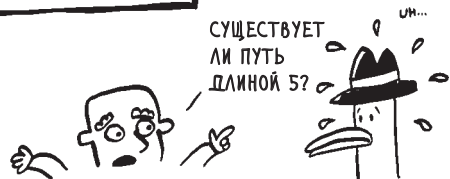
Найти кратчайший путь!



А вот как выглядит задача о коммивояжере, сформулированная в виде задачи разрешимости:



Существует ли путь длиной 5?



Вопрос формулируется так, чтобы на него можно было дать ответ «да»/«нет»: существует ли путь длиной 5? Я решил начать с задач разрешимости, потому что все NP-полные задачи являются таковыми. Итак, *все задачи, которые будут рассматриваться в приложении, будут задачами разрешимости*. А значит, при упоминании «задачи о коммивояжере» будет иметься в виду задача о коммивояжере в форме *разрешимости*.

А теперь разберемся, что же на самом деле означает NP-полнота. Для начала необходимо познакомиться с задачами выполнимости (SAT, Satisfiability Problem).

Задачи выполнимости



Джерри, Джордж, Элен и Крамер заказывают пиццу.

«Давайте возьмем с пепперони!» — говорит Элен.

«Пепперони — хорошо. Колбаса — хорошо. Можно взять с пепперони или с колбасой», — добавляет Джерри.

«А мне пиццу с оливками, чтобы фигуру не портить, — говорит Крамер. — Много оливок. Или лука».

«Согласен на любую пиццу, только без лука, — возражает Джордж. — Терпеть не могу лук, Джерри!»

«Ясно. Давайте разберемся. Еще раз, какие топпинги нам нужны?» — спрашивает Джерри.

А вы сможете ему помочь? Еще раз перечислим требования всех участников:

- Пепперони (Элен).
- Пепперони или колбаса (Джерри).
- Оливки или лук (Крамер).
- Без лука (Джордж).

Прежде чем читать дальше, попробуйте самостоятельно определить, какие топпинги должны быть в пицце.

Получилось? Пицца с пепперони и оливками удовлетворяет всем требованиям. Перед вами пример задачи SAT. На псевдокоде ее можно записать следующим образом. Сначала определяются четыре логические переменные:

```
pepperoni = ?
sausage = ?
olives = ?
onions = ?
```

Затем записывается логическая формула:

```
(pepperoni) and (pepperoni or sausage) and (olives or
onions) and (not onions)
```

Эта формула содержит требования каждого участника в форме булевой логики. Задача SAT задает вопрос: можно ли присвоить этим переменным такие значения, чтобы следующее выражение давало результат `true`?

Задача SAT хорошо известна, потому что это первая NP-полная задача, описанная в 1971 году (хотя я не думаю, что авторы решали ее на примере пиццы). До этого концепции NP-полноты не существовало. Давайте разберемся, как работает задача SAT. Все начинается с логической формулы:

```
if (pepperoni) and (olives or onions):
    print("pizza")
```

После этого задается вопрос: можно ли присвоить значения переменным так, чтобы этот код вывел сообщение `pizza`?

Это довольно простой пример, который легко решить вручную. Если pepperoni и onions истинны, то код выведет сообщение pizza. Следовательно, ответ будет положительным.

А вот пример задачи, в которой ответ будет отрицательным:

```
if (olives or onions) and (not olives) and (not onions):
    print("pizza")
```

Невозможно задать значения переменных, чтобы этот код вывел сообщение pizza!



Задача SAT всегда ищет ответ «да» или «нет», поэтому это задача *разрешимости*.

Задача SAT в действительности довольно сложная. Приведу более серьезный пример просто, чтобы вы это поняли:

```
if (pepperoni or not olives) and (onions or not
pepperoni) and (not olives or not pepperoni):
    print("pizza")
```

Решать ее не нужно. Это всего лишь пример, который показывает, насколько трудной может быть такая задача. Количество переменных и условий может быть любым, и задачи очень быстро усложняются.

С n добавками возможны 2^n вариантов пиццы. Если составить их список и проверить каждый вариант, получится так называемая таблица истин-

ности. Справа приведена таблица истинности для комбинации pepperoni and (olives or onions).

ПЕПЕРОНИ	ОЛИВКИ	ЛУК	ОТВЕТ
F	F	T	НЕТ
F	T	F	НЕТ
F	T	T	НЕТ
T	F	F	НЕТ
T	F	T	ДА
T	T	F	ДА
T	T	T	ДА
F	F	F	НЕТ

ПЕПЕРОНИ И
(ОЛИВКИ ИЛИ ЛУК)

Иногда необходимо перебирать все варианты, как в задаче о покрытии множества и задаче о коммивояжере. Как выясняется, задача SAT так же сложна, как и эти две задачи. Она характеризуется временем выполнения $O(2^n)$.

Трудно решить, легко проверить

Мы часто видим задачи, в которых найти решения намного сложнее, чем проверить его. Допустим, я предлагаю вам придумать предложение, которое является палиндромом (то есть читается одинаково слева направо и справа налево) и при этом включает слова *cat* и *car*. Как вы думаете, сколько времени вам понадобится, чтобы решить эту задачу?

Теперь представьте, как я говорю, что мне известно такое предложение. Вот оно: «*Was it a car or a cat I saw?*»

Проверить это утверждение можно гораздо быстрее, чем составить собственное предложение. Задачу трудно решить, но потенциальное решение быстро проверяется.

Задачу SAT решить так же сложно, как и задачу о покрытии множества или задачу о коммивояжере, но, в отличие от этих задач, ее решение проверяется быстро. Например, для приведенного выше запроса (pepperoni or not olives) and (onions or not pepperoni) and (not olives or not pepperoni) решение выглядит так:

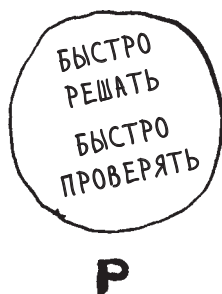
```
pepperoni = False
olives = False
onions = False
```

Легко убедиться, что с этими значениями логическая формула дает истинный результат. Проверить эти значения намного быстрее, чем решить задачу самостоятельно!



Решения задачи SAT *быстро проверяются*, поэтому она относится к классу NP.

В класс NP входят задачи, решения которых могут быть *проверены* за полиномиальное время. NP-задачи могут быть простыми или трудными в решении, но их решения легко проверить. Этим они отличаются от класса P.



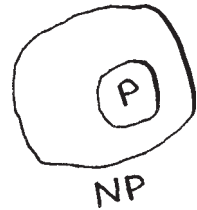
Задача относится к классу P, если ее можно *решить и проверить* за полиномиальное время. Полиномиальное время означает, что оценка «О-большое» не больше полиномиальной. В книге я не даю определение полиномиальности, но ниже приведены два примера полиномов.

$$n^3 \quad n^2 + n$$

И еще пара примеров, которые не являются полиномиальными.

$$n! \quad 2^n$$

P является подмножеством NP. Таким образом, NP содержит все задачи из P, а также другие задачи.



P и NP

Вы наверняка знаете известную задачу равенства классов P и NP. Мы только что видели, что задачи P быстро проверяются и решаются. Задачи NP быстро проверяются, но могут решаться быстро или не быстро. Задача равенства классов P и NP фактически спрашивает, будет ли каждая задача, которая быстро проверяется, так же быстро решаться. Если эта гипотеза истинна, то класс P не будет подмножеством NP; он будет равен NP.

NP-трудность — следующий термин, которому мы дадим определение, но сначала необходимо (кратко) обсудить, что такое сведéние.

Сведéние

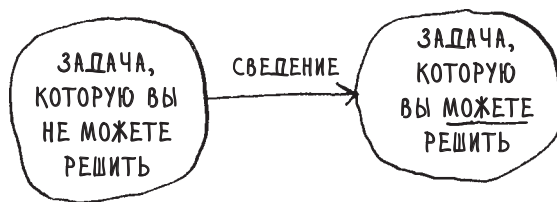
Что вы делаете, столкнувшись со сложной задачей? Стараетесь свести ее к задаче, которую вы можете решить! В реальных условиях при решении трудной задачи очень часто ее стараются изменить.

Посмотрим, как это сделать. Как перемножить два двоичных числа? Например, таких:

$$101 * 110$$

Если вы думаете как я, скорее всего, вам не захочется выполнять умножение в двоичной системе. Можно просто определить, что 101 соответствует десятичному числу 5, а 110 — десятичному числу 6, а затем умножить 5 на 6.

Такой подход называется «сведение». Задача, которую вы не умеете решать, сводится к задаче, которую вы решать умеете. В информатике он встречается очень часто.



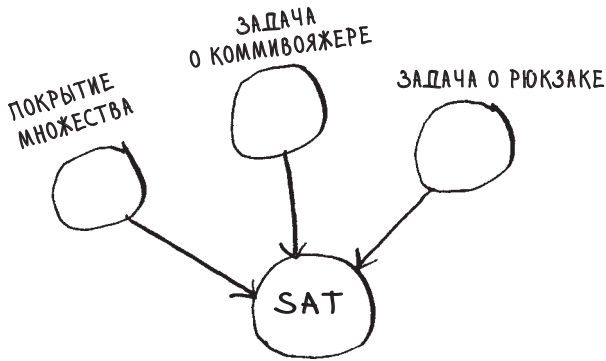
NP-трудность

Мы уже рассмотрели три примера NP-трудных задач:

- задачу о покрытии множества;
- задачу о коммивояжере;
- задачу SAT.

(Напомню, что речь идет о формах *разрешимости* этих задач — все задачи, рассматриваемые в этом приложении, относятся к категории разрешимых.) Три ранее рассмотренные задачи являются NP-трудными. Задача называется NP-трудной, *если какая-либо задача из класса NP может быть сведена к этой задаче*. Это определение NP-трудности.

Все NP-задачи также могут быть сведены к любой NP-трудной задаче. Например, все задачи NP можно свести к SAT.



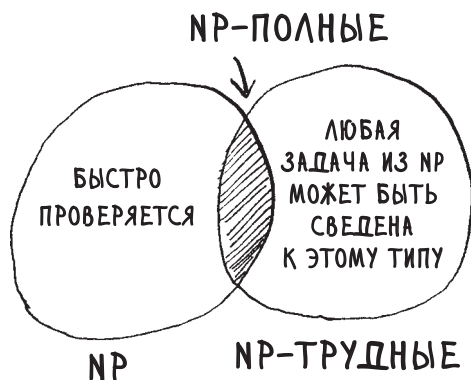
Еще одно требование: сведение всех этих задач должно выполняться *за полиномиальное время*. Уточнение «за полиномиальное время» важно, потому что фаза сведения не должна стать «узким местом». Любая NP-задача может быть сведена к SAT за полиномиальное время, поэтому она является NP-трудной. Так как любая задача из NP может быть сведена к любой NP-трудной задаче, решение с полиномиальным временем для любой NP-трудной задачи дает решение с полиномиальным временем для каждой задачи в NP!

NP-полнота

Мы познакомились с двумя утверждениями:

- Решения задач из класса NP быстро проверяются, и при этом такие задачи могут быстро решаться (но не обязательно).
- Задачи, относящиеся к NP-трудным, по крайней мере не уступают по сложности самым сложным задачам из NP, а любая задача из NP может быть сведена к задаче в классе NP-трудных.

Приведу последнее утверждение: задача является NP-полной, если она одновременно является NP-задачей и NP-трудной задачей.



NP-полные задачи:

- сложно решать (по крайней мере пока; если кто-нибудь докажет, что $P = NP$, ситуация изменится);
- легко проверить.

И любую задачу из NP можно свести к задаче, которая является NP-полной.

Итак, в этом приложении мы дали определение следующим понятиям:

- задачи разрешимости;
- задача SAT;
- P и NP;
- сведение;
- NP-трудность;
- NP-полнота.

Когда в следующий раз речь пойдет об NP-полных задачах, надеюсь, вы будете уверенно чувствовать себя, зная эти термины.

Шпаргалка

- Задача относится к классу P , если поиск и проверка ее решения выполняются быстро.
- Задача относится к классу NP , если ее решение можно быстро проверить, но не обязательно быстро решить.
- Если можно найти быстрый (с полиномиальным временем) алгоритм для каждой задачи в NP , значит, $P = NP$.
- Задача NP является NP -трудной, если ее можно свести к этой задаче.
- Задача является NP -полной, если она одновременно является и NP -задачей, и NP -трудной задачей.

В

Ответы к упражнениям



Глава 1

- 1.1** Имеется отсортированный список из 128 имен, и вы ищете в нем значение методом бинарного поиска. Какое максимальное количество проверок для этого может потребоваться?

Ответ: 7.

- 1.2** Предположим, размер списка увеличился вдвое. Как изменится максимальное количество проверок?

Ответ: 8.

- 1.3** Известна фамилия, нужно найти номер в телефонной книге.

Ответ: $O(\log n)$.

- 1.4** Известен номер, нужно найти фамилию в телефонной книге. (Подсказка: вам придется провести поиск по всей книге!)

Ответ: $O(n)$.

- 1.5** Нужно прочитать номера всех людей в телефонной книге.

Ответ: $O(n)$.

- 1.6** Нужно прочитать телефоны всех людей, фамилии которых начинаются с буквы «А». (Вопрос с подвохом! В нем задействованы концепции, которые более подробно рассматриваются в главе 4. Прочитайте ответ — скорее всего, он вас удивит!)

Ответ: $O(n)$. Возможно, кто-то подумает: «Я делаю это только для одной из 26 букв, а значит, время выполнения должно быть равно $O(n/26)$ ». Запомните простое правило: в «О-большое» игнорируются числа, задействованные в операциях сложения, вычитания, умножения или деления. Ни одно из следующих значений не является правильной записью «О-большое»: $O(n + 26)$, $O(n - 26)$, $O(n * 26)$, $O(n/26)$. Все они эквивалентны $O(n)$! Почему? Если вам интересно, найдите раздел «Снова об “О-большом”» в главе 4 и прочитайте о константах в этой записи (константа — это просто число; в этом вопросе 26 является константой).

Глава 2

2.1 Допустим, вы строите приложение для управления финансами.

- 1. ПРОДУКТЫ**
- 2. КИНО**
- 3. ВЕЛОСИПЕДНЫЙ КЛУБ**

Ежедневно вы записываете все свои траты. В конце месяца вы анализируете расходы и вычисляете, сколько денег было потрачено. При работе с данными выполняется множество операций вставки и относительно немного операций чтения. Какую структуру использовать — массив или список?

Ответ: В данном случае траты добавляются в список ежедневно, а чтение всех данных происходит один раз в месяц. Для массивов характерно быстрое чтение и медленная вставка, а для связанных списков — медленное чтение и быстрая вставка. Так как вставка будет выполняться намного чаще, чем чтение, есть смысл воспользоваться связанным списком. Кроме того, чтение в связанных списках происходит медленно только при обращении к случайным элементам списка. Так как читаться будут все элементы списка, связанный список также неплохо справится с чтением. Итак, связанный список станет хорошим решением этой задачи.

- 2.2** Допустим, вы пишете приложение для приема заказов от посетителей ресторана. Приложение должно хранить список заказов. Официанты добавляют заказы в список, а повара читают заказы из списка и выполняют их. Заказы образуют очередь: официанты добавляют заказы в конец очереди, а повар берет первый заказ из очереди и начинает готовить.



Какую структуру данных вы использовали бы для реализации этой очереди — массив или связанный список? (Подсказка: связанные списки хорошо подходят для вставки/удаления, а массивы — для произвольного доступа к элементам. Что из этого понадобится в данном случае?)

Ответ: Связанный список. Вставка происходит очень часто (официанты добавляют заказы), а связанные списки эффективно выполняют эту операцию. Ни поиск, ни произвольный доступ (сильные стороны массивов) вам не понадобятся, потому что повар всегда извлекает из очереди первый заказ.

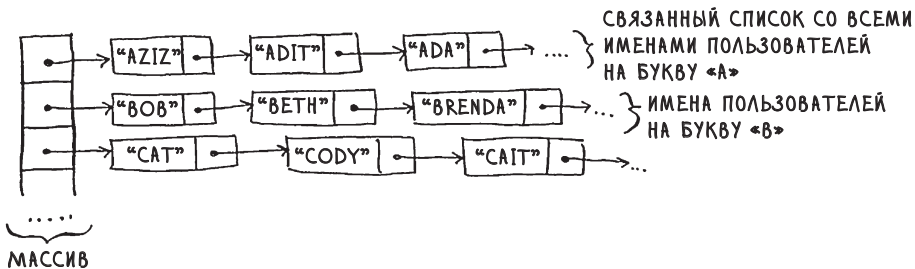
- 2.3** Проведем мысленный эксперимент. Допустим, Facebook хранит список имен пользователей. Когда кто-то пытается зайти на сайт Facebook, система пытается найти имя пользователя. Если имя входит в список имен зарегистрированных пользователей, то вход разрешается. Пользователи приходят на Facebook достаточно часто, поэтому поиск по списку имен пользователей будет выполняться часто. Будем считать, что Facebook использует бинарный поиск для поиска в списке. Бинарному поиску необходим произвольный доступ — алгоритм должен мгновенно обратиться к среднему элементу текущей части списка. Зная это обстоятельство, как бы вы реализовали список пользователей — в виде массива или связанного списка?

Ответ: В виде отсортированного массива. Массивы обеспечивают произвольный доступ — вы можете мгновенно получить элемент из середины массива. Со связанными списками это невозможно. Чтобы получить элемент из середины связанного списка, вам придется начать с первого элемента и переходить по ссылкам до нужного элемента.

- 2.4** Пользователи также довольно часто создают новые учетные записи на Facebook. Предположим, вы решили использовать массив для хранения списка пользователей. Какими недостатками обладает массив для выполнения вставки? Допустим, вы используете бинарный поиск для нахождения учетных данных. Что произойдет при добавлении новых пользователей в массив?

Ответ: Вставка в массив выполняется медленно. Кроме того, если вы используете бинарный поиск для нахождения имен пользователей, массив необходимо отсортировать. Предположим, пользователь по имени *Adit B* регистрируется на Facebook. Его имя будет вставлено в конец массива. Следовательно, массив нужно будет сортировать при каждой вставке нового имени!

- 2.5** В действительности Facebook не использует ни массив, ни связанный список для хранения информации о пользователях. Рассмотрим гибридную структуру данных: массив связанных списков. Имеется массив из 26 элементов. Каждый элемент содержит ссылку на связанный список. Например, первый элемент массива указывает на связанный список всех имен пользователей, начинающихся на букву «А». Второй элемент указывает на связанный список всех имен пользователей, начинающихся на букву «В», и т. д.



Предположим, пользователь с именем «Adit B» регистрируется в Facebook и вы хотите добавить его в список. Вы обращаетесь к элементу 1 массива, находите связанный список элемента 1 и добавляете

«Adit B» в конец списка. Теперь предположим, что зарегистрировать нужно пользователя «Zakhir H». Вы обращаетесь к элементу 26, который содержит связанный список всех имен, начинающихся с «Z», и проверяете, присутствует ли «Zakhir H» в этом списке.

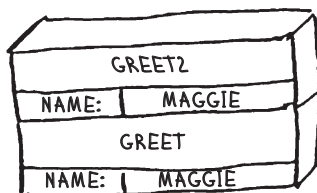
Теперь сравните эту гибридную структуру данных с массивами и связанными списками. Будет она быстрее или медленнее каждой исходной структуры при поиске и вставке? Приводить время выполнения «О-большое» не нужно, просто выберите одно из двух: быстрее или медленнее.

Ответ: Поиск — медленнее, чем для массивов, и быстрее, чем для связанных списков. Вставка — быстрее, чем для массивов, и с такой же скоростью для связанных списков. Итак, гибридная структура уступает массиву по скорости поиска, но по крайней мере не хуже связанных списков для всего остального. Далее в книге будет рассмотрена другая гибридная структура данных, называемая хеш-таблицей. Она даст некоторое представление о том, как строить сложные структуры данных из более простых.

Что же в действительности использует сервис Facebook? Вероятно, десяток разных баз данных, за которыми стоят разные структуры данных: хеш-таблицы, В-деревья и т. д. Массивы и связанные списки становятся структурными элементами для построения более сложных структур данных.

Глава 3

3.1 Предположим, имеется стек вызовов следующего вида:



Что можно сказать о текущем состоянии программы на основании этого стека вызовов?

Ответ: Некоторые наблюдения, о которых вы могли бы упомянуть:

- сначала вызывается функция `greet` для переменной `name = maggie`;
- затем функция `greet` вызывает функцию `greet2` для переменной `name = maggie`;
- на этой стадии функция `greet` находится в незавершенном, приостановленном состоянии;
- текущим вызовом функции является вызов `greet2`;
- после завершения этого вызова функция `greet` продолжит выполнение.

3.2 Предположим, вы случайно написали рекурсивную функцию, которая бесконечно вызывает саму себя. Как вы уже видели, компьютер выделяет память в стеке при каждом вызове функции. А что произойдет со стеком при бесконечном выполнении рекурсии?

Ответ: Стек будет расти бесконечно. Каждой программе выделяется ограниченный объем памяти в стеке. Когда все пространство будет исчерпано (а рано или поздно это произойдет), программа завершится с ошибкой переполнения стека.

Глава 4

4.1 Напишите код для функции `sum` (см. выше).

Ответ:

```
def sum(list):
    if list == []:
        return 0
    return list[0] + sum(list[1:])
```

4.2 Напишите рекурсивную функцию для подсчета элементов в списке.

Ответ:

```
def count(list):
    if list == []:
        return 0
    return 1 + count(list[1:])
```

4.3 Найдите наибольшее число в списке.

Ответ:

```
def max(list):
    if len(list) == 2:
        return list[0] if list[0] > list[1] else list[1]
    sub_max = max(list[1:])
    return list[0] if list[0] > sub_max else sub_max
```

4.4 Помните бинарный поиск из главы 1? Он тоже относится к классу алгоритмов «разделяй и властвуй». Сможете ли вы определить базовый и рекурсивный случаи для бинарного поиска?

Ответ: Базовым случаем для бинарного поиска является массив, содержащий всего один элемент. Если искомый элемент совпадает с элементом массива — вы нашли его! В противном случае элемент в массиве отсутствует.

В рекурсивном случае для бинарного поиска массив делится пополам, одна половина отбрасывается, а для другой половины проводится бинарный поиск.

Запишите «О-большое» для каждой из следующих операций.

4.5 Вывод значения каждого элемента массива.

Ответ: $O(n)$.

4.6 Удвоение значения каждого элемента массива.

Ответ: $O(n)$.

4.7 Удвоение значения только первого элемента массива.

Ответ: $O(1)$.

4.8 Создание таблицы умножения для всех элементов массива. Например, если массив состоит из элементов [2, 3, 7, 8, 10], сначала каждый элемент умножается на 2, затем каждый элемент умножается на 3, затем на 7 и т. д.

Ответ: $O(n^2)$.

Глава 5

Какие из следующих функций являются последовательными?

5.1 $f(x) = 1$ ◀..... **Возвращает "1" для любых входных значений**

Ответ: Функция последовательна.

5.2 $f(x) = \text{rand}()$ ◀..... **Возвращает случайное число**

Ответ: Функция непоследовательна.

5.3 $f(x) = \text{next_empty_slot}()$ ◀..... **Возвращает индекс следующего пустого элемента в хеш-таблице**

Ответ: Функция непоследовательна.

5.4 $f(x) = \text{len}(x)$ ◀..... **Возвращает длину полученной строки**

Ответ: Функция последовательна.

Предположим, имеются четыре хеш-функции, которые получают строки.

1. Первая функция возвращает «1» для любого входного значения.
2. Вторая функция возвращает длину строки в качестве индекса.
3. Третья функция возвращает первый символ строки в качестве индекса. Таким образом, все строки, начинающиеся с «a», хешируются в одну позицию, все строки, начинающиеся с «b», — в другую, и т. д.
4. Четвертая функция ставит в соответствие каждой букве простое число: a = 2, b = 3, c = 5, d = 7, e = 11 и т. д. Для строки хеш-функцией становится остаток от деления суммы всех значений на размер хеша. Например, если размер хеша равен 10, то для строки «bag» будет вычислен индекс $3 + 2 + 17 \% 10 = 22 \% 10 = 2$.

В каком из этих примеров хеш-функции будут обеспечивать хорошее распределение? Считайте, что хеш-таблица содержит 10 элементов.

5.5 Телефонная книга, в которой ключами являются имена, а значениями — номера телефонов. Задан следующий список имен: Esther, Ben, Bob, Dan.

Ответ: Хеш-функции C и D обеспечивают хорошее распределение.

5.6 Связь размера батарейки с напряжением. Размеры батареек: А, АА, ААА, АААА.

Ответ: Хеш-функции В и D обеспечивают хорошее распределение.

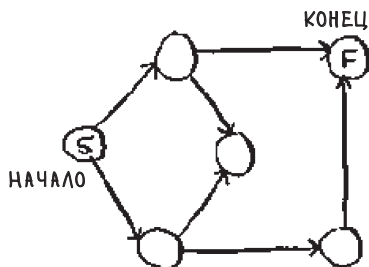
5.7 Связь названий книг с именами авторов. Названия книг: «Maus», «Fun Home», «Watchmen».

Ответ: Хеш-функции В, С и D обеспечивают хорошее распределение.

Глава 6

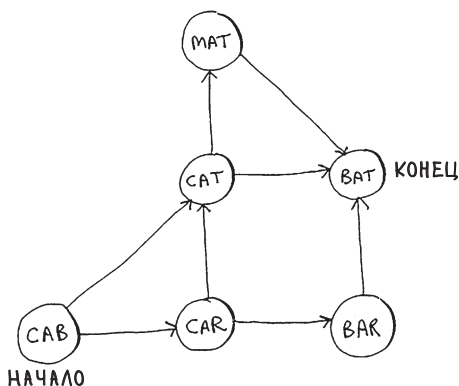
Примените алгоритм поиска в ширину к каждому из этих графов, чтобы найти решение.

6.1 Найдите длину кратчайшего пути от начального до конечного узла.



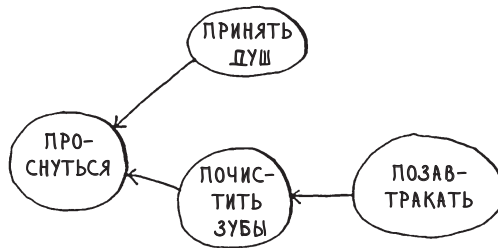
Ответ: Длина кратчайшего пути равна 2.

6.2 Найдите длину кратчайшего пути от «cab» к «bat».



Ответ: Длина кратчайшего пути равна 2.

6.3 Перед вами небольшой граф моего утреннего распорядка.

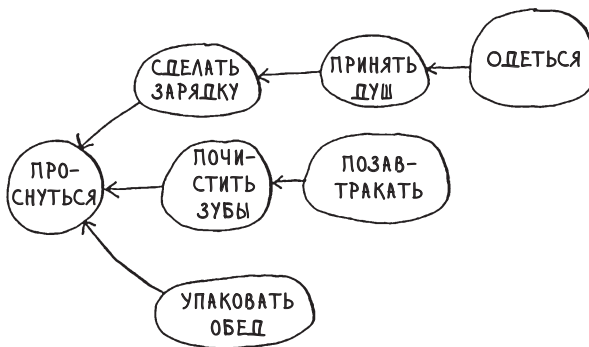


Для каждого из следующих трех списков укажите, действителен он или недействителен.

А	В	С
1. ПРОСНУТЬСЯ	1. ПРОСНУТЬСЯ	1. ПРИНЯТЬ ДУШ
2. ПРИНЯТЬ ДУШ	2. ПОЧИСТИТЬ ЗУБЫ	2. ПРОСНУТЬСЯ
3. ПОЗАВТРАКАТЬ	3. ПОЗАВТРАКАТЬ	3. ПОЧИСТИТЬ ЗУБЫ
4. ПОЧИСТИТЬ ЗУБЫ	4. ПРИНЯТЬ ДУШ	4. ПОЗАВТРАКАТЬ

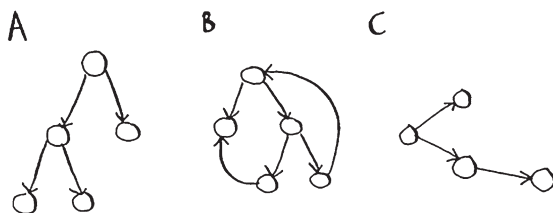
Ответы: А — недействителен; В — действителен; С — недействителен.

6.4 Немного увеличим исходный граф. Постройте действительный список для этого графа.



Ответ: 1 — Проснуться; 2 — Сделать зарядку; 3 — Принять душ; 4 — Почистить зубы; 5 — Одеться; 6 — Упаковать обед; 7 — Позавтракать.

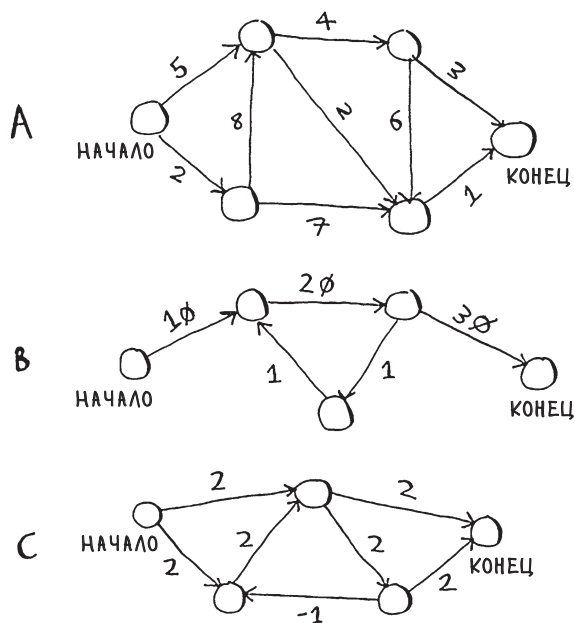
6.5 Какие из следующих графов также являются деревьями?



Ответы: А — дерево; В — не дерево; С — дерево. В последнем примере дерево просто повернуто набок. Деревья составляют подкатегорию графов, поэтому любое дерево является графом, но граф не обязательно является деревом.

Глава 9

9.1 Каков вес кратчайшего пути от начала до конца в каждом из следующих графов?



Ответы: А — 8; В — 60; С — каверзный вопрос (кратчайший путь не существует из-за наличия цикла с отрицательным весом).

Глава 10

10.1 Вы работаете в фирме по производству мебели и поставляете мебель по всей стране. Коробки с мебелью размещаются в грузовике. Все коробки имеют разный размер, и вы стараетесь наиболее эффективно использовать доступное пространство. Как выбрать коробки для того, чтобы загрузка имела максимальную эффективность? Предложите жадную стратегию. Будет ли полученное решение оптимальным?

Ответ: Жадная стратегия заключается в том, чтобы выбрать самую большую коробку, помещающуюся в оставшемся пространстве, и повторять это до тех пор, пока еще можно выбрать хотя бы одну коробку. Нет, такое решение оптимальным не будет.

10.2 Вы едете в Европу, и у вас есть 7 дней на знакомство с достопримечательностями. Вы присваиваете каждой достопримечательности стоимость в баллах (насколько вы хотите ее увидеть) и оцениваете продолжительность поездки. Как обеспечить максимальную стоимость (увидеть все самое важное) во время поездки? Предложите жадную стратегию. Будет ли полученное решение оптимальным?

Ответ: Выбирайте достопримечательность с наибольшей стоимостью в баллах, которую вы успеете посетить в оставшееся время. Остановитесь, когда таких достопримечательностей не останется. Нет, такое решение оптимальным не будет.

Глава 11

11.1 Предположим, к предметам добавился еще один: механическая клавиатура. Она весит 1 фунт и стоит \$1000. Стоит ли ее брать?

Ответ: Да. Вы сможете положить в рюкзак механическую клавиатуру, iPhone и гитару общей стоимостью \$4500.

11.2 Предположим, что вы собираетесь в турпоход. Емкость вашего рюкзака составляет 6 фунтов, и вы можете взять предметы из следующего списка. У каждого предмета имеется стоимость; чем она выше, тем важнее предмет:

- вода, 3 фунта, 10;
- книга, 1 фунт, 3;

- еда, 2 фунта, 9;
- куртка, 2 фунта, 5;
- камера, 1 фунт, 6.

Как выглядит оптимальный набор предметов для похода?

Ответ: Возьмите воду, еду и камеру.

11.3 Нарисуйте и заполните таблицу для вычисления самой длинной общей подстроки между строками *blue* и *clues*.

Ответ:

	C	L	U	E	S
B	0	0	0	0	0
L	0	1	0	0	0
U	0	0	2	0	0
E	0	0	0	3	0

Глава 12

12.1 В примере с Netflix сходство между двумя пользователями оценивалось по формуле расстояния. Но не все пользователи оценивают фильмы одинаково. Допустим, есть два пользователя, Йоги и Пинки, вкусы которых совпадают. Но Йоги ставит 5 баллов любому фильму, который ему понравился, а Пинки более разборчива и ставит «пятерки» только лучшим фильмам. Вроде бы вкусы одинаковые, но по метрике расстояния они не являются соседями. Как учесть различия в стратегиях выставления оценок?

Ответ: Можно воспользоваться нормализацией: вы вычисляете среднюю оценку для каждого человека и используете ее для масштабирования оценок. Например, вы определили, что средняя оценка Пинки равна 3, а средняя оценка Йоги — 3,5. Соответственно оценки Пинки немного увеличиваются так, чтобы ее средняя оценка тоже была равна 3,5. После этого оценки можно сравнивать по единой шкале.

12.2 Предположим, Netflix определяет группу «авторитетов». Скажем, Квентин Тарантино и Уэс Андерсон относятся к числу авторитетов Netflix, поэтому их оценки оказывают более сильное влияние, чем оценки рядовых пользователей. Как изменить систему рекомендаций, чтобы она учитывала повышенную ценность оценок авторитетов?

Ответ: При применении алгоритма k ближайших соседей можно увеличить вес оценок авторитетов. Предположим, у вас трое соседей: Джо, Дэйв и Уэс Андерсон (авторитет). Они поставили фильму «Гольф-клуб» оценки 3, 4 и 5 соответственно. Вместо того чтобы вычислять среднее арифметическое их оценок ($3 + 4 + 5 / 3 = 4$ звезды), вы просто повышаете вес оценки Уэса Андерсона: $3 + 4 + 5 + 5 + 5 / 5 = 4,4$ звезды.

12.3 У сервиса Netflix миллионы пользователей. В приведенном ранее примере рекомендательная система строилась для пяти ближайших соседей. Пять — это слишком мало? Слишком много?

Ответ: Слишком мало. Если ограничиться малым числом соседей, существует высокая вероятность того, что результаты будут искажены. Существует хорошее эмпирическое правило: для N пользователей следует рассматривать \sqrt{N} соседей.

Адितья Бхаргава

Грокаем алгоритмы

2-е издание

Перевел с английского Е. Матвеев

Руководитель дивизиона	<i>Ю. Сергиенко</i>
Ведущий редактор	<i>Е. Строганова</i>
Литературный редактор	<i>М. Трусковская</i>
Художественный редактор	<i>В. Мостипан</i>
Корректоры	<i>Н. Викторова, Т. Никифорова</i>
Верстка	<i>Л. Егорова</i>

Изготовлено в России. Изготовитель: ООО «Прогресс книга».
Место нахождения и фактический адрес: 194044, Россия, г. Санкт-Петербург,
Б. Сампсониевский пр., д. 29А, пом. 52. Тел.: +78127037373.

Дата изготовления: 08.2024. Наименование: книжная продукция. Срок годности: не ограничен.

Налоговая льгота — общероссийский классификатор продукции ОК 034-2014, 58.11.12.000 — Книги печатные профессиональные, технические и научные.

Импортер в Беларусь: ООО «ПИТЕР М», 220020, РБ, г. Минск, ул. Тимирязева, д. 121/3, к. 214, тел./факс: 208 80 01.

Подписано в печать 19.06.24. Формат 70×100/16. Бумага офсетная. Усл. п. л. 28,380. Тираж 10000. Заказ 0000.

Эрик Норманд

ГРОКАЕМ ФУНКЦИОНАЛЬНОЕ МЫШЛЕНИЕ



Кодовые базы разрастаются, становясь все сложнее и запутаннее, что не может не пугать разработчиков. Как обнаружить код, изменяющий состояние вашей системы? Как сделать код таким, чтобы он не увеличивал сложность и запутанность кодовой базы?

Большую часть «действий», изменяющих состояние, можно превратить в «вычисления», чтобы ваш код стал проще и логичнее.

Вы научитесь бороться со сложными ошибками синхронизации, которые неизбежно проникают в асинхронный и многопоточный код, узнаете, как komponуемые абстракции предотвращают дублирование кода, и откроете для себя новые уровни его выразительности.

Книга предназначена для разработчиков среднего и высокого уровня, создающих сложный код. Примеры, иллюстрации, вопросы для самопроверки и практические задания помогут надежно закрепить новые знания.

КУПИТЬ

Ришал Харбанс

ГРОКАЕМ АЛГОРИТМЫ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА



Искусственный интеллект — часть нашей повседневной жизни. Мы встречаемся с его проявлениями, когда занимаемся шопингом в интернет-магазинах, получаем рекомендации «вам может понравиться этот фильм», узнаем медицинские диагнозы...

Чтобы уверенно ориентироваться в новом мире, необходимо понимать алгоритмы, лежащие в основе ИИ.

«Грокаем алгоритмы искусственного интеллекта» объясняет фундаментальные концепции ИИ с помощью иллюстраций и примеров из жизни. Все, что вам понадобится, — это знание алгебры на уровне старших классов школы, и вы с легкостью будете решать задачи, позволяющие обнаружить банковских мошенников, создавать шедевры живописи и управлять движением беспилотных автомобилей.

КУПИТЬ