

Теоретическая информатика для биологов

А.А.Миронов

2019

Предисловие

Предлагаемая Вашему вниманию книга основана на семестровом курсе лекций по теоретической информатике, который читается на факультете Биоинженерии и биоинформатики Московского Государственного университета им М. В. Ломоносова. Основной целью данного курса является стремление дать самые общие представления о теоретической информатике. Это необходимо, чтобы читать и понимать алгоритмические статьи по биоинформатике и не падать в обморок от пары формул, псевдокода или теорем. Поэтому в этой книге приведено достаточно много алгоритмов на псевдокоде. Есть леммы и теоремы — достаточно простые, поскольку доказательство, как правило, умещается в нескольких строчках. С другой стороны, иногда даются некоторые абстрактные понятия, такие, как полукольцо, язык, конечный автомат. Это представляется важным для того, чтобы у читателя возникла привычка если не оперировать абстракциями, то по крайней мере наблюдать за тем, как это делают другие. При выборе тем и рассмотренных алгоритмов автор опирался на собственный опыт и представление о прекрасном. Возможно другой преподаватель при составлении плана лекций по теоретической информатике для биологов несколько изменил бы содержание, но, думаю, примерно 70% материала сохранилось бы, поскольку они являются базовыми и используются в биоинформатике. В первой главе приведены некоторые основные понятия информатики и булевой алгебры, разбираются описываются способы записи алгоритмов. Замечу, что поскольку для записи псевдокода нет стандарта, то в книге использован синтаксис наиболее похожий на язык C или Java. Вторая глава посвящена поиску, сортировке и некоторым основным структурам данных — списки, деревья, хеш-таблицы. В третьей главе обсуждаются алгоритмы поиска образца в тексте. В конце главы даны основные представления о теории сложности А.Н. Колмогорова. Четвертая глава является продолжением третьей и посвящена преобразованию Барроуза-Уиллера. Это довольно специальный раздел, но это преобразование лежит в основе задач картирования прочтений на геном, связанных с транскриптомикой и другими пост-геномными методами. Пятая глава посвящена алгоритмам на графах. Приведено несколько примеров того, как биологические задачи могут формулироваться на языке графов. Разобраны наиболее важные алгоритмы в этой области. Наконец, шестая глава является наиболее абстрактной. В ней даются представления о теории вычислительной сложности.

В книге есть некоторое количество задач, достаточно простых, чтобы не давать их решений. Книга рассчитана на студентов и аспирантов, специализирующихся в биоинформатике и молекулярной биологии. Автор приносит глубокую благодарность своим бывшим студентам Ольге Ганчаровой и Надежде Тухтумбаевой, которые взяли на себя труд записать и оцифровать конспект лекций. Именно эти конспекты легли в основу предлагаемой книги. Позже конспекты были существенно поправлены усилиями Евгении Елизаровой и Марии Гурылевой, за что им отдельное спасибо.

Оглавление

1	Некоторые предварительные сведения	7
1.1	Компьютеры, программы и алгоритмы	8
1.1.1	Время работы алгоритма	9
1.2	Типы данных	10
1.2.1	Основные типы данных	10
1.2.2	Массивы	11
1.2.3	Указатели	12
1.2.4	Структуры	12
1.3	Способы описания алгоритмов	13
1.4	Элементы булевой алгебры	16
1.4.1	Основные операции булевой алгебры	16
1.4.2	Другие часто используемые функции	17
1.5	Машина Тьюринга	19
1.5.1	Структура Машины Тьюринга	19
1.5.2	Краткое описание машины Тьюринга	20
1.6	Примеры машин Тьюринга	21
1.6.1	Скопировать строку	21
1.6.2	Вычислить длину строки	23
2	Поиск, сортировка и структуры данных	27
2.1	Поиск элемента в массиве	27
2.1.1	Поиск в несортированном массиве	27
2.1.2	Поиск в сортированном массиве	30
2.1.3	Анализ алгоритма бинарного поиска	32
2.2	Сортировка массивов	34
2.2.1	Анализ алгоритма Merge_Sort	36
2.2.2	Быстрая сортировка (QSort)	37
2.3	Списки	40
2.3.1	Односвязный список	40
2.3.2	Двусвязный список	43
2.3.3	Пример применения связного списка. Организация дан- ных в памяти компьютера	45
2.3.4	Очередь и стек	46
2.4	Бинарное дерево поиска	49

2.4.1	Поиск в бинарном дереве	50
2.4.2	Добавление элемента в бинарное дерево поиска	50
2.4.3	Поиск минимального (максимального) элемента в дереве	52
2.4.4	Поиск следующего по возрастанию элемента	52
2.4.5	Удаление элемента из бинарного дерева поиска	54
2.4.6	Преобразование деревьев поиска	56
2.4.7	Красно-черные деревья	57
2.5	Хеш-таблицы	61
2.6	Заключительные замечания	64
3	Алгоритмы для строк	67
3.1	Определения, постановка задачи. Наивный алгоритм	67
3.2	Алгоритм Рабина-Карпа	70
3.3	Алгоритм Кнута-Морриса-Пратта	72
3.3.1	Префикс-функция $sp_i(P)$	72
3.3.2	Алгоритм поиска образца	73
3.3.3	Препроцессинг Кнута-Морриса-Пратта	75
3.4	Конечные автоматы	77
3.4.1	Поиск образца с помощью конечного автомата	79
3.4.2	Поиск регулярных выражений. Недетерминированные конечные автоматы.	82
3.5	Поиск многих образцов	85
3.6	Поиск многих образцов в тексте	90
3.6.1	Суффиксные массивы	90
3.6.2	Суффиксное дерево	92
3.6.3	Алгоритмы построения суффиксного дерева	93
3.6.4	Сравнение суффиксных деревьев и суффиксных мас- сивов	95
3.6.5	Примеры применения суффиксных деревьев	96
3.7	Сложность текста	98
3.8	Посимвольное сжатие (код Хаффмана)	100
3.9	Сжатие RLE	102
3.10	Сжатие информации по Лемпелю-Зиву	102
3.11	Сжатие информации по Лемпелю-Зиву-Вельчу	103
3.11.1	Сжатие	104
3.11.2	Декодирование	105
4	Преобразование Барроуза-Уиллера	107
4.1	Преобразование Барроуза Уиллера	107
4.1.1	Концепция	107
4.1.2	Быстрое восстановление исходного текста	109
4.1.3	Поиск образца с помощью преобразования BWT	110
4.1.4	Оптимизация памяти	112
4.1.5	Сжатие суффиксного массива	116
4.1.6	Итоги	117

5	Алгоритмы для графов	119
5.1	Способы представления графов	120
5.2	Обход графа в ширину	123
5.3	Обход графа в глубину	127
5.3.1	Типы ребер графа	129
5.3.2	Построение покрывающего дерева	130
5.3.3	Поиск связных компонент	130
5.3.4	Поиск циклов	130
5.3.5	Некоторые приложения к биоинформатике	131
5.4	Задача Эйлера	133
5.4.1	Алгоритм поиска Эйлера цикла в графе	134
5.4.2	Биоинформатика и задача Эйлера: секвенирование геномов	136
5.5	Топологическая сортировка	139
5.5.1	Алгоритм топологической сортировки	140
5.6	Поиск оптимального пути в графе	141
5.6.1	Динамическое программирование для поиска оптимального пути	142
5.6.2	Оценка времени работы динамического программирования	145
5.6.3	Биоинформатические применения поиска оптимального пути в графе. Выравнивание	145
5.6.4	Полукольцо. Динамическое программирование над полукольцом	148
5.6.5	Поиск минимального пути в графе, содержащем циклы. Алгоритм Дейкстры	149
5.6.6	Алгоритм Беллмана-Форда	156
5.6.7	Сводка основных алгоритмов, ищущих оптимальный путь на графе	157
5.7	Сети и потоки в них	157
5.7.1	Основные понятия	157
5.7.2	Метод Форда-Фалкерсона	159
5.7.3	Разрезы в сетях	160
5.7.4	Теорема о максимальном потоке и минимальном разрезе	163
5.7.5	Случай многих источников и стоков	163
5.7.6	Паросочетания в двудольном графе	163
5.8	Понятие о линейном программировании	164
5.8.1	Геометрическая интерпретация	166
5.8.2	Немного истории	167
5.8.3	Симплекс метод	167
6	Свойства задач и NP-полные задачи.	169
6.1	Нижняя оценка времени сортировки массива	169
6.2	Допускающий и распознающий алгоритмы	171
6.3	Типы задач	172
6.4	Проверяющие алгоритмы. Класс языков NP	173

6.5	NP-полные задачи	174
6.5.1	Иерархия типов задач	176
6.6	Некоторые NP-полные задачи	176
6.6.1	Задачи выполнимости SAT и 3CNF.	179
6.6.2	Задача о клике (CLIQUE).	182
6.6.3	Решение NP-полных задач	183
6.6.4	Переборные алгоритмы	184
6.7	Стохастические алгоритмы	185
6.7.1	Искусственный отжиг	186
6.7.2	Генетические алгоритмы	188

Глава 1

Некоторые предварительные сведения

Предлагаемая книга дает краткое и достаточно поверхностное введение в теоретическую информатику. Эта наука по-английски называется Computer Science. Зачем информатика биологу? Существенной частью практически любого современного исследования в области молекулярной биологии является биоинформатическое исследование, которое либо предваряет экспериментальную работу, либо ее завершает, хотя чаще всего биоинформатика применяется и до и после экспериментальной работы. Это связано в большой степени с тем, что современные работы в области молекулярной биологии опираются на большое количество разнообразных данных, таких как геномы, протеомы, экспрессионные данные и т.п. Анализ этих данных немислим без современных методов компьютерного анализа. Разумеется, автор не предполагает, что биологи должны сами создавать алгоритмы анализа этих данных, писать программы и т.п. — для этого есть специалисты. Однако биологу придется взаимодействовать со специалистами в области информатики и программистами. Для этого необходимо иметь хотя бы общее представление о методах и подходах теоретической информатики. Кроме того, грамотный биолог читает научную литературу. А в ней часто встречаются работы, описывающие новые методы анализа. Можно привести ряд биологических журналов, в которых публикуются работы, посвященные методам компьютерного анализа данных, зачастую эти работы содержат также экспериментальную часть. Чтобы понимать эти статьи необходимо знакомство с языком информатики и основными идеями и методами. Что касается литературы. Во-первых, Google в помощь. Во-вторых есть замечательная толстая книжка Т.Кормен, Ч.Лейзерсон, Р.Ривест, К.Штайн — Алгоритмы. Построение и анализ. Это книга намного больше предлагаемого курса, но она совершенно замечательная. Ну, и еще рекомендую сайт <http://rain.ifmo.ru/cat/view.php/aboutus/vis-authors>, на котором есть весьма наглядная визуализация различных алгоритмов, надо толь-

ко браузеру разрешить работать с апплетами. С алгоритмами и/или программами вы встречаетесь постоянно, даже не замечая этого. Например, вы назначили встречу с Васей Пупкиным на факультете Биоинженерии и биоинформатики. Вы звоните по телефону и говорите примерно следующий текст. “Доедешь до метро Университет, выйдешь по указателям к Университету, перейдешь Ломоносовский проспект. Если нет дождя, то пойдешь вдоль Ломоносовского проспекта. Пройдешь примерно два километра, далее мимо Биологического факультета, увидишь наш корпус. Если же идет дождь, то садись на любой автобус, кроме номеров 1, 47, 119, или на троллейбус 34. На транспорте проедешь 3 остановки и выйдешь на ул. Менделеева. На входе покажешь что-нибудь похожее на студ. билет. Если пропустят, то поднимайся на 4-й этаж, иначе звони мне, я проведу по своей карточке”. В таком объяснении много неточностей и недоговорченностей. Что здесь важно? Во-первых, вы предполагаете, что Вася Пупкин ориентируется в метро, умеет сесть в автобус, умеет считать до четырех, наконец, умеет ходить и читать. Теперь представьте, что вы объясняете, как доехать до факультета, вашей подруге из Швеции. Если она только что приехала в Москву, то ей надо еще объяснить, как дойти до метро, как туда войти, где, как и за сколько купить карточку, где и как сделать пересадку и т.п. Итак, при описании алгоритма предполагается, что существует исполнитель, и что исполнитель кое-что умеет. Определений для термина алгоритм существует довольно много. Вот лишь одно из них. Алгоритм — это точный однозначно понимаемый набор инструкций, описывающий порядок действий некоторого исполнителя для достижения конкретного результата за конечное время.

1.1 Компьютеры, программы и алгоритмы

Здесь мы будем рассматривать алгоритмы для компьютеров. Хотя компьютеры бывают разные (в мире есть не только пентиумы), все они обладают рядом общих свойств:

- Компьютеры имеют линейно организованную память, в которой можно запоминать промежуточные результаты. К элементам памяти можно обращаться по *адресам* — целым числам.
- Компьютеры имеют процессор (исполнитель), который умеет выполнять арифметические операции, операции сравнения, операции записи в память и извлечения данных из памяти.
- Процессор один. Вы слышали, конечно, о многоядерных процессорах и многопроцессорных кластерах. Наверняка у кого-то из вас есть дома четырехядерный компьютер. Количество процессоров обычно заранее известно и ограничено, поэтому нет принципиальной разницы один процессор или несколько. Многопроцессорными вычислениями называются такие, где количество процессоров пропорционально размеру

задачи, например, при молекулярном моделировании каждому атому отводится отдельный процессор. Многопроцессорные вычисления выходят за рамки этого курса.

- Для современных компьютеров процессор, строго говоря, не всегда исполняет команды последовательно одну за другой. Здесь мы предполагаем, что инструкции выполняются последовательно.

Алгоритм должен за конечное время (пусть и достаточно большое) найти решение.

Следствие 1. *Для решения любой задачи используется конечная память (докажите).*

Хотя в любом реальном компьютере размер памяти вполне определен, мы будем считать, что память в компьютере не ограничена. Важной особенностью алгоритма является то, что он имеет входные данные и производит результат. Для одного и того же алгоритма можно использовать разные входные данные, и он будет производить разные результаты. Даже алгоритм, вычисляющий значение числа π , имеет разные входные данные в зависимости от требуемой точности. Один и тот же алгоритм может быть

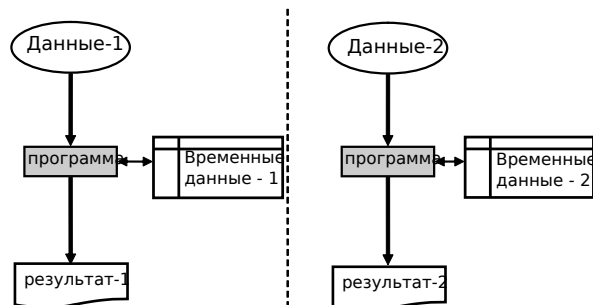


Рис. 1.1: Программа (алгоритм) принимает входные данные, в процессе обработки порождает временные данные и выдает результат.

реализован на компьютере в виде разных программ — это зависит от платформы компьютера, языка программирования, стиля программиста и т.п.

1.1.1 Время работы алгоритма

Важнейшим вопросом теоретической информатики является время работы того или иного алгоритма. Размер входной задачи (например, длина анализируемого слова или текста) может варьировать. Поскольку алгоритм

принимает разные входные данные, то и время его работы может быть разным. Будем считать, что время работы алгоритма пропорционально числу выполняемых операций (сложений, сравнений и т.п.). На самом деле это не совсем правильно, поскольку в реальных компьютерах время исполнения разных команд может отличаться в десятки раз. Обычно принято оценивать время работы алгоритма как функцию от размера задачи, причем важно его асимптотическое поведение (при больших размерах задачи). Но для современных архитектур компьютеров важными становятся и другие аспекты и для алгоритмов вводят и другие характеристики. Поскольку реальное время зависит от деталей устройства процессора и от деталей перевода кода в инструкции процессора, то время работы оценивается с точностью до порядка. Говорят, что величина T есть $O(f(L))$ (читается: O большое от $f(L)$), если существует такая константа C , что для любого допустимого L $T(L) \leq C \cdot f(L)$. Время работы алгоритма может зависеть не только от размера задачи, но и от самих данных. Поэтому для времени работы алгоритма существуют оценки в худшем и в среднем. При оценке времени в худшем рассматриваются такие данные (заданного размера), которые дадут самое большое количество операций. При оценке времени в среднем исходным данным заданного размера приписывается вероятность их реализации и оценивается математическое ожидание времени работы алгоритма. Говорят, что один алгоритм A_1 эффективнее алгоритма A_2 , если $T_1(N)/T_2(N) \rightarrow 0$ при $N \rightarrow \infty$, где N — размер задачи. Алгоритмы используют память для хранения промежуточных данных. Количество используемой памяти также является важной характеристикой алгоритма.

1.2 Типы данных

В памяти компьютера хранятся данные, которые бывают разных типов.

1.2.1 Основные типы данных

Примитивные числовые типы данных — это байт, целое число со знаком, целое число без знака, число с плавающей точкой. Также к примитивным типам данных относятся `boolean` и `char`.

Байт (*byte*). Это целочисленный беззнаковый тип данных. Принимает значения от -128 до $+127$, всего 256 значений, которые можно записать с помощью двух шестнадцатеричных цифр (0123456789abcdef).

Целое число со знаком. Тип включает переменные, способные принимать целочисленные значения, положительные или отрицательные, и значение нуль. К этому типу относятся переменные вида `int` (занимает 4 байта, соответственно может принимать значения в пределах ± 2147483647), `long` и `short`.

Целое число без знака (*unsigned short*, *unsigned long*, *unsigned int*). переменные, аналогичные целочисленным переменным со знаком, но

способные принимать только положительные значения и значение нуль. Беззнаковый `int` может принимать значения от 0 до $4 \cdot 10^9$.

Число с плавающей точкой (*double*, *float*) представление дробных чисел, когда число хранится в форме мантиссы и показателя степени. К примеру, число 0.02 представляется в виде двух чисел: мантиссы 2 и показателя степени -2 . Как и целые числа, этот тип данных может быть реализован с привлечением разных объемов памяти, что влияет на максимальное принимаемое значение. К примеру, в языках C и Java есть два типа чисел с плавающей точкой — *float* (занимает 4 байта) и *double* (занимает 8 байтов). *double* — 64-разрядный тип данных, при этом 54 разряда выделяется под мантиссу, 10 — под показатель степени. К точности чисел с плавающей точкой нужно относиться с осторожностью, особенно при сравнении двух чисел. Известно, что если поделить единицу на число 9999999999, а потом умножить на него же, и спросить Java, Python, или C, равен ли результат единице, получим ответ *false*, поскольку при округлении и переводе из двоичной в десятичную систему и обратно происходит снижение точности. В представлении *double* мантисса меняется в пределах от $+(2^{54})/2$ до $-(2^{54})/2$, что соответствует примерно 15-ти десятичным знакам. Однако, существуют числа с длиной значащей части больше 15 знаков и даже иррациональные числа. Записав любое такое число в переменную типа *double*, мы неизбежно теряем точность, поскольку 16-ое и последующие значащие числа будут отброшены. Так что, если мы хотим сравнить два числа с плавающей точкой (*a* и *b*), лучше делать это с умом: взять какое-либо малое число (уровень значимости) и сравнивать разность $a - b$ с этим числом.

Булева переменная *boolean* может принимать значения “истина” или “ложь”.

Для реализации типа достаточно одного бита, и часто для экономии памяти удобно представлять данные в виде переменных типа *boolean*. Если у вас есть только два значения ($\{1,0\}$ или $\{+, -\}$) разумнее всего представлять их в виде *boolean*. Однако в архитектуре современных компьютеров нет прямого доступа к отдельным битам, а элементарной единицей является байт. Поэтому в реальности для значений *boolean* используется байт.

1.2.2 Массивы

Часто используются массивы — это пронумерованные множества данных. К элементу массива можно обратиться с помощью индекса. Массивы лежат в памяти подряд. Мы не будем вдаваться в достаточно сложную архитектуру памяти реальных компьютеров. Следуя соглашениям в языках C, Java, Python мы будем считать, что первый элемент имеет индекс 0. Надо заметить, что в некоторых языках программирования, например R или SQL, нумерация идет с 1. Если у нас есть массив *a*, то элементом массива с но-

мером i будет $a[i]$. Длину массива будем обозначать как `a.length`. Частным случаем массива является строка — массив из символов.

1.2.3 Указатели

Особое место в современном программировании занимают *указатели* — это переменные, хранящие в памяти адреса в памяти, которые указывают нам место, где хранятся данные. При этом сами данные могут быть размещены в самых разных местах памяти. На рис. 1.2 а) показан фрагмент памяти компьютера. Верхняя строка — адреса (они идут подряд). Нижняя строка — содержимое ячеек памяти (данные). Данные могут иметь самый разный характер — числа, символы (тоже представлены числами). В частности в данных могут лежать адреса — указатели на другие ячейки памяти. Использование указателей дает множество преимуществ. Например, массив строк (а каждая строка в свою очередь является массивом) может быть организован как массив указателей на строки. В этом случае перестановка двух строк не требует перетаскивания всех символов строки с места на место (что потребует большого количества операций). Достаточно только переставить два указателя. Кроме того, использование указателей позволяет не плодить сущностей. Например, полный телефонный номер содержит код страны, код места или сети и собственно телефонный номер. Вместо кода места или сети удобно хранить указатель на этот код. При изменении кода мы можем просто заменить содержание одной ячейки памяти — и все номера, содержащие этот код, автоматически его поменяют. Заметим, что косвенная адресация (через указатели) может замедлять работу программы. Это связано с особенностями архитектуры процессоров и кэшированием памяти. Но мы не будем обращать на это внимания.

1.2.4 Структуры

Часто некоторые данные образуют комплекс. Например, фамилия, имя, дата рождения и номер зачетной книжки образуют неразделимую запись и эти данные должны жить вместе. В этом случае говорят о *структуре*. При этом отдельные элементы называются полями. Обратите внимание, что дата также является структурой, которая состоит из числа, месяца, года. Мы будем описывать структуру следующим образом:

1. `Struct Date{year, month, day;}`
2. `Struct Stud{FamilyName, FirstName, BirthDay, Id}`

К полям структуры мы будем обращаться через символ `.` (точка). Например, если у нас есть студент `Purkine`, то номер его зачетки будем обозначать как `Purkine.Id`. Поля структуры в памяти лежат подряд. Отметим важную деталь. При ссылке на другую структуру или строку мы всегда будем понимать, что в структуре лежит указатель на нее, а сама структура может при этом лежать в другом месте памяти. Наши соглашения следуют

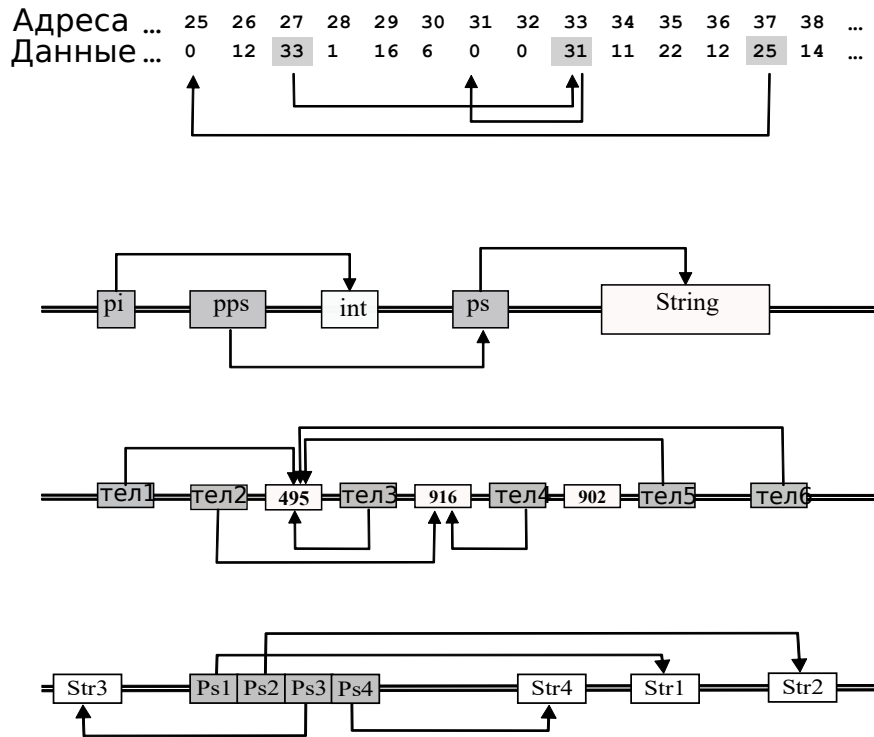


Рис. 1.2: Указатели (показаны серым) содержат адреса в памяти, по которым расположены сами объекты. а) фрагмент памяти. Верхняя строка — адреса (они идут подряд), нижняя строка — содержимое ячеек памяти. б) pps — указатель на указатель. в) иллюстрирует ситуацию с кодами сети. Четыре телефона ссылаются на код 495. При смене кода на 499 достаточно поменять значение в одной ячейке памяти — и сразу автоматически Тел1, Тел.3, Тел.5 и Тел.6 приобретут новые коды. г) Массив строк — указатели лежат в памяти подряд, Сами строки разбросаны по памяти.

синтаксису языка Java и Python. На языке C, например, нам следовало бы писать `Purkine->FamilyName`.

Упражнение 1. Приведите примеры структур.

1.3 Способы описания алгоритмов

Для представления алгоритма обычно используются различные способы:

- в виде текстовых описаний действий (текстуальная форма); Эта форма не достаточно строгая, хотя часто она более понятна

- в виде блок-схем (графическая форма); Эта форма достаточно прозрачна и на нее даже есть стандарт. Она наиболее удобна для представления потоков данных, но описание алгоритмов с помощью блок-схем громоздко, особенно, если в алгоритме есть циклы и ветвления. Ее можно запатентовать.
- в виде псевдокода (способ описания алгоритмов, использующий ключевые слова языков программирования, но опускающий подробности и специфический синтаксис). Этот способ весьма удобен для тех, кто знаком с программированием. На псевдокод нет стандарта — в разных источниках встречаются разные модификации. Мы будем использовать версию, наиболее похожую на язык C или Java. Текст на псевдокоде не есть непосредственно программа, поскольку обычно опускается целый ряд важных деталей

Пример. Дан массив из чисел. Надо его упорядочить, т.е. переставить числа так, что-бы они шли по возрастанию. Представим решение во всех трех формах.

Словесное описание (наивная сортировка). Ищем в массиве наименьший элемент. Меняем местами первый элемент и наименьший. Ищем наименьший элемент в части массива, начинающейся со второго элемента. Меняем местами второй элемент и найденный наименьший и т.д.

Блок-схема и пояснения к ней показаны на рис.1.3

Псевдокод:

```

1. NaiveSort (a[ ]){
2.     for ( i=0; i < a.length ; i++){
3.         min = a[i];
4.         jmin = i;
5.         for (j=i+1; j < a.length; j++){
6.             if (min > a[j]){
7.                 min = a[j];
8.                 jmin = j;
9.             }
10.        }
11.        a[jmin] = a[i];
12.        a[i] = min;
13.    }
14. }
```

Несколько особенностей и встречающихся вариаций. В этом примере есть цикл, условный оператор, переменные, присваивания. В разных версиях

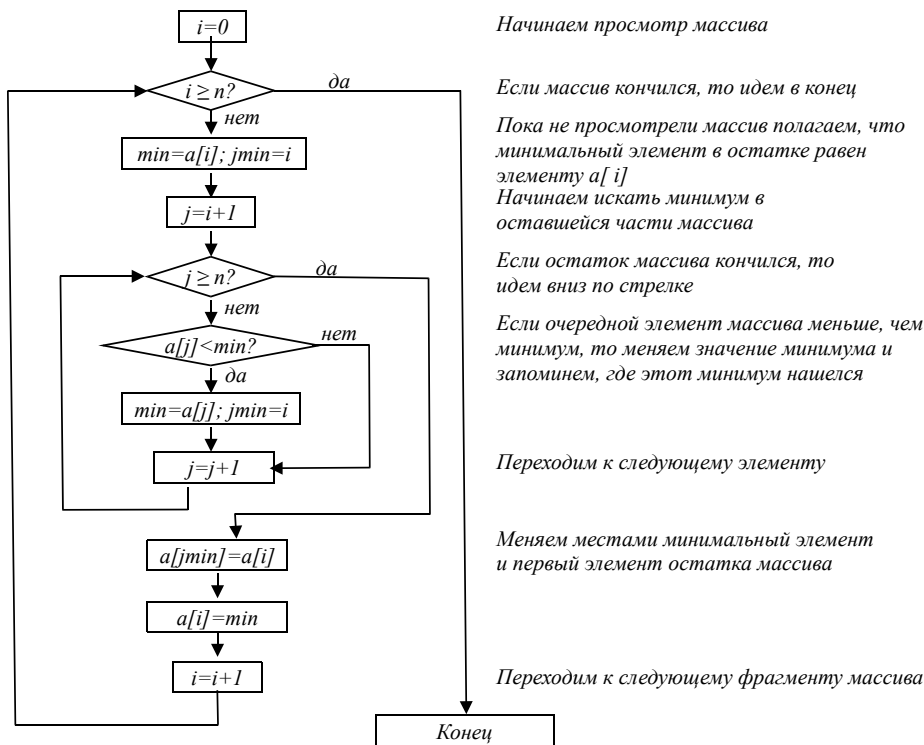


Рис. 1.3: Пример блок-схемы для наивной сортировки массива.

псевдокода встречаются варианты для обозначения этих действий, например:

Цикл	for i = 0 until n
Условный оператор	if min > a[i] then
Присваивание	a[i] := min или a[i] ← min

Упражнение 2. Метод сортировки "всплывающий пузырек" заключается в следующем. Идем вдоль массива. Если встретим инверсию (ситуацию, когда два соседних элемента имеют неправильный порядок), меняем их местами и проверяем, не нарушился ли порядок для предыдущей пары элементов. Если порядок предыдущей пары неправильный, то также эти два элемента меняем местами. И так пока не исчезнет инверсия. После чего продолжаем поиск следующей инверсии. Представьте этот алгоритм в виде блок-схемы псевдокода.

Упражнение 3. Оцените время работы алгоритма наивной сортировки.

1.4 Элементы булевой алгебры

Булева алгебра (названа в честь английского математика XIX века Джорджа Буля) рассматривает величины, принимающие только два значения — 1 (true; истина) или 0 (false; ложь). С такими величинами можно производить различные операции — так же, как мы оперируем с утверждениями при рассуждениях. Поэтому булева алгебра является основой формальной логики. Итак, дано множество из двух элементов

$$B = \{0, 1\} \text{ или } B = \{false, true\}$$

В булевой алгебре рассматриваются функции в многомерном множестве B^n вида:

$$B^n \rightarrow B$$

Это функции, принимающие на вход n переменных, а выдающих в качестве значения только false(0) или true(1). В отличие от функций действительных переменных, которые рассматривались на мат.анализе, количество разнообразных булевых функций заданной размерности — конечно. Сколько всего разных булевых функций от n переменных? Например, для отображения $B \rightarrow B$ могут быть такие варианты:

0 → 0	0 → 1	0 → 1	0 → 0
1 → 0	1 → 1	1 → 0	1 → 1
false	true	$\neg x$	x

то есть существует 4 типа функции $B \rightarrow B$:

$$\{B \rightarrow B\} = \{true, false, \neg x, x\}$$

Если речь идет об отображении $B^2 \rightarrow B$ ($n = 2$; двухмерность: 2 значения на входе, 1 — на выходе), то число вариантов функции увеличивается и становится равным уже 2^4 . В общем виде:

$$|\{B^n \rightarrow B\}| = 2^{2^n} \tag{1.1}$$

Упражнение 4. Докажите равенство 1.1.

1.4.1 Основные операции булевой алгебры

Выделяют несколько основных операций (функций) над булевыми переменными: логическое “или” (дизъюнкция, OR, логическое сложение, \vee , |), логическое “и” (конъюнкция, AND, логическое умножение, \wedge , &), логическое отрицание “не” (NOT, \neg , !). В программировании вместо математических символов для операций используют обозначения:

текст	программа
\vee	
\wedge	&
\neg	!

Значения функций двух переменных (например, OR или AND) или одной переменной (например, NOT) удобно записывать в виде матрицы: OR (\vee) — логическое «или»:

or	0	1
0	0	1
1	1	1

AND (\wedge) — логическое «и»:

and	0	1
0	0	0
1	0	1

NOT (\neg) — отрицание:

not	0
0	0
1	1

Можно записывать значение булевой функции (любого числа переменных) в виде таблицы. Например,

X,Y	and
0,0	0
0,1	0
1,0	0
1,1	1

Задача: чему равно значение $f(x, y) = (x \wedge y) \vee (x \wedge \neg y)$? Решение: запишем все возможные значения в виде таблицы:

X,Y	$f(x, y)$
0,0	0
0,1	0
1,0	1
1,1	1

Отсюда $f(x, y) = x$.

1.4.2 Другие часто используемые функции

XOR (\wedge) — «исключающее или»:

\wedge	0	1
0	0	1
1	1	0

EQ (\equiv) — совпадение:

\equiv	0	1
0	1	0
1	0	1

Следствие (импликация, $y \Rightarrow x$ — “из y следует x ”):

$y \Rightarrow x$	0	1
0	1	1
1	0	1

Задача: Построим булеву функцию от n переменных x_1, x_2, \dots, x_n следующим образом: для начала напишем ВСЕ комбинации вида

$$[\neg]x_1 \wedge [\neg]x_2 \wedge \dots \wedge [\neg]x_n,$$

где $[\neg]$ означает, что в каких-то комбинациях берем отрицание, а в каких-то — нет. Затем все полученные комбинации объединяем знаком \vee . Чему равно значение полученной функции?

Решение: Итак, мы рассматриваем функцию такого вида:

$$f(x_1, x_2, \dots, x_n) = \begin{aligned} &(x_1 \wedge x_2 \wedge \dots \wedge x_n) \vee \\ &(\neg x_1 \wedge x_2 \wedge \dots \wedge x_n) \vee \\ &(x_1 \wedge \neg x_2 \wedge \dots \wedge x_n) \vee \\ &\dots \\ &(\neg x_1 \wedge \neg x_2 \wedge \dots \wedge \neg x_n) \end{aligned}$$

Идея решения состоит в том, что при любом наборе значений x_1, x_2, \dots, x_n хотя бы в одной скобке будет 1 (true) и, следовательно, $f(x_1, x_2, \dots, x_n) = true$. Более формальное доказательство возможно по индукции.

Важные тождества

$$\begin{aligned} \neg(x \vee y) &= \neg x \wedge \neg y \\ \neg(x \wedge y) &= \neg x \vee \neg y \\ x \wedge true &= x \\ x \vee true &= true \\ x \wedge false &= false \\ x \vee false &= x \\ x \equiv y &= (x \wedge y) \vee (\neg x \wedge \neg y) \\ x \hat{=} y &= (x \vee y) \wedge (\neg x \vee \neg y) \\ x \hat{=} y &= \neg(x \equiv y) \end{aligned}$$

Теорема 1. Любую логическую функцию n переменных можно представить с помощью основных операций (\wedge , \vee , \neg).

Доказательство. Доказательство данного утверждения предлагается проводить по индукции:

- для $n = 1$ есть 4 типа функции от одной переменной: $true, false, x, \neg x$. Они очевидно выражаются через базовые операции

$$\begin{aligned} true &= x \vee \neg x; \\ false &= x \wedge \neg x; \\ (x) &= x; \\ (\neg x) &= \neg x \end{aligned}$$

- пусть с помощью основных операций можно представить функцию $n - 1$ переменных. Тогда для любого n :

$$f(x_1, \dots, x_{n-1}, x_n) = (f(x_1, \dots, x_{n-1}, true) \wedge x_n) \vee (f(x_1, \dots, x_{n-1}, false) \wedge \neg x_n)$$

Заметим, что $f(x_1, \dots, x_{n-1}, true)$ и $f(x_1, \dots, x_{n-1}, false)$ на самом деле являются двумя функциями $n - 1$ переменной, и по предположению индукции могут быть представлены в виде комбинации базовых операций. Следовательно, всякую булеву функцию любого числа переменных можно разложить на функции *and*, *or*, *not*. Индукция завершена.

□

Утверждение теоремы можно еще сузить. Для представления любой логической функции любого числа переменных достаточно операций \vee, \neg , так как \wedge представить можно как: $(x \wedge y) = \neg(\neg x \vee \neg y)$. На самом деле булева алгебра может рассматриваться с одной стороны, как алгебра высказываний, а с другой, как алгебра множеств. Например утверждение {элемент x принадлежит множеству A } \vee {элемент x принадлежит множеству B } эквивалентно утверждению {элемент x принадлежит множеству $A \cup B$ }.

1.5 Машина Тьюринга

Машина Тьюринга (МТ) – математическая абстракция, представляющая вычислительную машину общего вида. Была предложена Аланом Тьюрингом в 1936 году для формализации понятия алгоритма

1.5.1 Структура Машины Тьюринга

Машина Тьюринга (МТ) состоит из двух частей – ленты и автомата.

Структура Машины Тьюринга и пояснения к ней показаны на рис.1.4. Лента используется для хранения информации. Она бесконечна в обе стороны и разбита на клетки, которые никак не нумеруются и не именуется. В каждой клетке может быть записан один символ или ничего не записано. Содержимое клетки может меняться – в нее можно записать другой символ или стереть находящийся там символ. Автомат – это активная часть МТ.

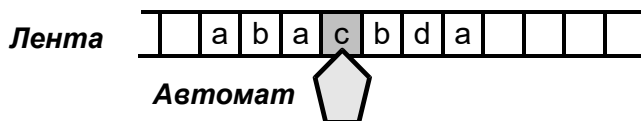


Рис. 1.4: Структура Машины Тьюринга.

В каждый момент он размещается под одной из клеток ленты и видит её содержимое; это видимая клетка, а находящийся в ней символ — видимый символ; содержимое же соседних и других клеток автомат не видит. Кроме того, в каждый момент автомат находится в одном из состояний, которые будем обозначать буквой q с дополнениями: q_1 , q_2 , q_{Left} и т.п. Находясь в некотором состоянии, автомат выполняет какую-то определённую операцию (например, перемещается направо по ленте, заменяя все символы b на a), находясь в другом состоянии — другую операцию.

1.5.2 Краткое описание машины Тьюринга

Пару из видимого символа (S) и текущего состояния автомата (q) будем называть конфигурацией и обозначать $\langle S, q \rangle$. Автомат может выполнять три элементарных действия: 1) записывать в видимую клетку новый символ (менять содержимое других клеток автомат не может); 2) сдвигаться на одну клетку влево или вправо (“перепрыгивать” сразу через несколько клеток автомат не может); 3) переходить в новое состояние. Ничего другого делать автомат не умеет, поэтому все более сложные операции так или иначе должны быть сведены к этим трём элементарным действиям. Автомат принимает некий набор состояний (конечный). Есть выделенное состояние $START$ и несколько (или одно) состояние $STOP$. МТ работает тактами (по шагам), которые выполняются один за другим. На каждом такте автомат МТ выполняет три следующих действия, причем обязательно в указанном порядке: 1) записывает некоторый символ S' в видимую клетку (в частности, может быть записан тот же символ, что и был в ней, тогда содержимое этой клетки не меняется); 2) сдвигается на одну клетку влево (обозначение — L , от left), либо на одну клетку вправо (обозначение — R , от right), либо остается неподвижным (обозначение — N). 3) переходит в некоторое состояние q' (в частности, может остаться в прежнем состоянии). Пример таблицы состояний. Машина Тьюринга является стандартной машиной, на которой можно показать работу всех алгоритмов, выполняемых на одном процессоре.

		Прочитанный символ				
		a1	a2	...	a _i	...
Текущее состояние автомата	q1	ak11 qm11 →	ak12 qm12 →			
	q2	ak21 qm21 →	ak qm →			
	...					
	q _i	ak _i 1 qm _i 1 →	ak _i 2 qm _i 2 →			
					q _m →	

ak – Символ, который будет напечатан

qm – Новое состояние

→ – то, куда автомат перейдет на следующем шаге

Рис. 1.5: Таблица состояний Машины Тьюринга.

1.6 Примеры машин Тьюринга

1.6.1 Скопировать строку

Дано: строка в некотором алфавите $A = \{0, 1, \emptyset, L, R\}$, где \emptyset – пустой символ, которым заполнены все пустые позиции ленты, L, R –обозначают правую и левую границы текста. С начала читающая головка автомата стоит в начале строки на символе L . Задача – написать после заданной строки ее копию. В целом автомат должен работать так: читаем символ и запоминаем его – как? – просто создаем два состояния в соответствии с основными буквами алфавита. Помечаем прочитанный символ – как? – записываем на его место символ $0'$ или $1'$. Мы расширили допустимый алфавит добавив в него дополнительные символы. Потом надо пойти вправо до тех пор, пока на ленте чего-то записано, потом, в зависимости от состояния (того, что мы запомнили), записать символ – 0 или 1. Здесь мы перейдем в состояние поиска очередного символа для чтения – идем влево до символа R – чтобы пропустить вновь записанные символы. Потом идем влево, пока не встретим прочитанный символ – это является границей прочитанных символов. Шаг вправо – и вот очередной непрочитанный символ. Если этот символ оказался R , то ясно, что мы уже все скопировали. Осталось символы, помеченный как прочитанные вернуть в исходное состояние. Теперь это надо формализовать. Итак, матрица состояний и переходов будет такая, как показано на таблица 1.1. Как работает эта машина показано на рис.1.6. Вначале стоим на начале строки, состояние start. Читаем символ L и переходим в состояние qx , которое ответственно за поиск очередного не скопированного символа. Движемся вперед, пока не найдем символа без штриха. Как только мы увидели символ без штриха, запоминаем его – переходим в соответствующее состояние – и несем его направо до первого свободного места, куда его и записываем. Освободившись от очередного символа идем налево, пока не встретим символа со штрихом. Первый символ со штрихом

Таблица 1.1: Таблица переходов для машины Тьюринга, копирующей строки

		L	0	1	0'	1'	R	\emptyset
Start	переходим вправо и читаем первый символ	L, qx →	err	err	err	err	err	err
qx	Готовы прочитать символ	err	0', q0 →	1', q1 →	err	err	R, fin ←	err
q0	Помним 0	err	0, q0 →	1, q0 →	err	err	R, q0 →	0, qL ←
q1	Помним 1	err	0, q1 →	1, q1 →	err	err	R, q1 →	1, qL ←
qL	Идем влево до следующего непрочитанного символа	err	0, qL ←	1, qL ←	0', qx →	1', qx →	R, qL ←	err
fin	Снимаем метки	stop	0, fin ←	1, fin ←	0, fin ←	1, fin ←	R, fin ←	err
err	ошибка	stop	stop	stop	stop	stop	stop	stop

– это последний символ, который мы перенесли. Переходим в состояние *qx* – поиска символа, который надо переместить. Он, конечно же следующий. Однако, если следующим символом является символ *R*, то это значит, что мы все скопировали, и надо снять все метки, чем и занимается состояние *fin*

Start	L	0	0	1	1	R	⊘	⊘	⊘	⊘
qx	L	0	0	1	1	R	⊘	⊘	⊘	⊘
q0	L	0'	0	1	1	R	0	⊘	⊘	⊘

q0	L	0'	0	1	1	R	0	⊘	⊘	⊘
qL	L	0'	0	1	1	R	0	⊘	⊘	⊘

qL	L	0'	0	1	1	R	0	⊘	⊘	⊘
qx	L	0'	0	1	1	R	0	⊘	⊘	⊘
qx	L	0'	0	1	1	R	0	⊘	⊘	⊘
q0	L	0'	0'	1	1	R	0	⊘	⊘	⊘

q0	L	0'	0'	1	1	R	0	0	⊘	⊘
qL	L	0'	0'	1	1	R	0	0	⊘	⊘

qL	L	0'	0'	1	1	R	0	0	⊘	⊘
qx	L	0'	0'	1	1	R	0	0	⊘	⊘
q1	L	0'	0'	1'	1	R	0	0	⊘	⊘
									⊘	⊘
q1	L	0'	0'	1'	1	R	0	0	1	⊘

qL	L	0'	0'	1'	1'	R	0	0	1	1
qx	L	0'	0'	1'	1'	R	0	0	1	1
fin	L	0'	0'	1'	1	R	0	0	1	1
fin	L	0'	0'	1	1	R	0	0	1	1

fin	L	0	0	1	1	R	0	0	1	1
stop	L	0	0	1	1	R	0	0	1	1

Рис. 1.6: Пример работы машины Тьюринга для копирования строк

1.6.2 Вычислить длину строки

Машина тьюринга также способна производить арифметические операции. Рассмотрим в качестве примера вычисление длины строки. Строка, как и в предыдущем примере представлена в алфавите $\{0,1\}$, плюс два символа, ограничивающих строку $\{L, R\}$, кроме того есть незаполненные ячейки \emptyset . Также, как и в предыдущем примере, введем два дополнительных символа $\{0', 1'\}$ для того, чтобы пометать уже обработанные символы. Длину строки будем вычислять в двоичной системе счисления. Значение длины строки поместим слева от строки. Машина будет работать так (1.2): Считываем очередной символ, переходим в состояние q_+ , которое означает, что теперь нам надо прибавить 1 к уже существующему счетчику. После этого двигаемся влево до начала строки. Здесь переходим в состояние прибавления q_{++} и прибавляем 1 к существующему счетчику. Увеличение счетчика

делаем так. Если в текущей ячейке 0, то заменяем его на 1, переходим в состояние q_0 , которое обнулит все ячейки с младшими разрядами. Чтобы понять, почему так, достаточно посмотреть, на то, как происходит прибавление 1 (рис.1.7). Прибавление сделано, можно идти вперед за следующим символом. Вот как это работает (рис.1.8). Сначала автомат стоит в начале

000	+1	→	001
001	+1	→	010
010	+1	→	011
011	+1	→	100
100	+1	→	101

Рис. 1.7: Прибавление 1 к разным числам в двоичной системе

строки. Далее он ищет вправо первую не посчитанную букву – букву без штриха. Как только он ее находит, он помечает эту букву, переходит в состояние готовности складывать ($q+$). Для того, чтобы начать складывать ему надо в этом состоянии пройти налево до начала строки. Пройдя начало строки, автомат меняет состояние на активное сложение ($q++$) и начинает прибавлять единицу к числу, которое записано перед строкой. Чтобы сделать сложение, надо найти самый правый ноль. Из этого нуля сделать единицу, а все, что идет за этой новой единицей и до начала строки – обнулить, чем и занимается состояние q_0 . После того, как мы учли последний символ строки, т.е. при поиске очередного символа уперлись символ конца строки R мы переводим автомат в завершающее состояние fin , которое должно стереть все метки. Почему мы пишем число символов слева от текста, а не справа? Дело в том, что мы начинаем прибавлять с младшего, самого правого разряда, а он как раз является первым символом перед началом строки. Если бы мы писали счетчик справа от текста, то мы получили бы число символов, записанное задом наперед со младшими разрядами слева и нам пришлось бы разворачивать строку.

Упражнение 5. Постройте машину Тьюринга для сложения двух двоичных чисел. Упражнение посложнее – построить машину Тьюринга для сложения десятичных чисел.

start	⊘	⊘	⊘	L	1	0	1	R	⊘
qx	⊘	⊘	⊘	L	1'	0	1	R	⊘
q+	⊘	⊘	⊘	L	1'	0	1	R	⊘
q++	⊘	⊘	1	L	1'	0	1	R	⊘
q0	⊘	⊘	1	L	1'	0	1	R	⊘
q++	⊘	⊘	1	L	1'	0'	1	R	⊘
...
q++	⊘	1	1	L	1'	0'	1	R	⊘
q0	⊘	1	0	L	1'	0'	1	R	⊘
qx	⊘	1	0	L	1'	0'	1	R	⊘
...
q+	⊘	1	0	L	1'	0'	1'	R	⊘
q++	⊘	1	1	L	1'	0'	1'	R	⊘
q0	⊘	1	1	L	1'	0'	1'	R	⊘
...
qx	⊘	1	1	L	1'	0'	1'	R	⊘
fin	⊘	1	1	L	1'	0'	1	R	⊘
...
stop	⊘	1	1	L	1	0	1	R	⊘

Рис. 1.8: Пример работы машины Тьюринга, подсчитывающей длину строки

Глава 2

Поиск, сортировка и структуры данных

2.1 Поиск элемента в массиве

2.1.1 Поиск в несортированном массиве

Рассмотрим простейшую задачу — найти телефон в телефонной книге. Для этого разберем сначала, что такое телефонная книга. Телефонная книга состоит из записей (вспомним, что такое структура), а каждая запись состоит из имени и номера телефона. При поиске в телефонной книге, как правило, ищут по имени. То поле, по которому производят поиск, называется ключом, а остальные поля называются данными. В некоторых случаях ключ и данные совпадают (например, в простом массиве чисел). Итак, ищем в книге Тюткина. Поиск устроен очень просто.

1. $i = 0$
2. Берем элемент номер i
3. Проверим, не Тюткин ли это
4. Если Тюткин, то **Ура!** (сообщаем индекс Тюткина)
5. Иначе переходим к следующему элементу: $i = i + 1$
6. Если массив не кончился, то идем к 2
7. **Увы...** (индекс -1)

Блок-схема алгоритма приведена на рис.2.1

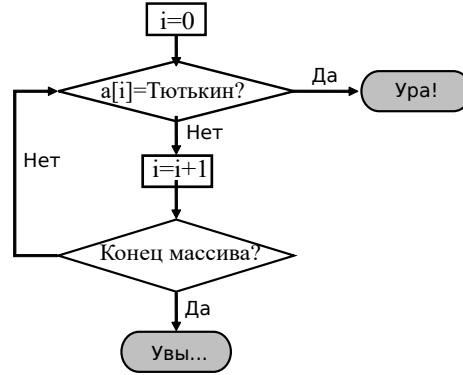


Рис. 2.1: блок-схема поиска в не сортированном массиве

Псевдокод:

```

1. SearchElement (Elm e, a[]){
2.   for (i =0; i<a.length; i++){
3.     if (a[i] == e)
4.       return ‘‘Ура!!!’’;
5.   }
6.   return (Увы...)
7. }

```

Время работы алгоритма в лучшем случае $T = O(1)$ — когда первый элемент содержит искомый ключ. Но это бывает крайне редко. Худшим случаем для алгоритма является отсутствие искомого элемента, или если искомый элемент находится в конце массива. Время работы в худшем случае есть $T = O(L)$. В среднем искомый элемент находится в середине массива. В первом случае нам в среднем надо просмотреть половину массива, во втором надо просмотреть весь массив. И в том и в другом случае время работы алгоритма $T = O(L)$, где L — длина массива. Так ли это? Проведем более строгий анализ. Пусть вероятность совпадения ключа с запросом поиска равна p , а вероятность несовпадения $q = 1 - p$. Тогда вероятность того, что цикл на строке 2 пройдет один раз, равна p . Вероятность того, что цикл пройдет две итерации равна pq , поскольку на первой итерации было несовпадение, а на второй — совпадение. Вероятность того, что цикл остановится на итерации i , равна $q^{i-1}p$. Математическое ожидание числа оборотов цикла равно:

$$E = p + 2pq + 3pq^2 + \dots + (L-1)pq^{L-2} + Lq^{L-1} = p(1 + 2q + 3q^2 + \dots + (L-1)q^{L-2}) + Lq^{L-1}$$

Последний член суммы не содержит множителя p , поскольку здесь нам уже все равно, совпал последний элемент или нет — цикл-то мы все равно

прокрутили. Вычислим сумму в скобках. Это вроде бы почти геометрическая прогрессия, только нам мешают множители $2, 3, \dots, (L - 1)$. Нетрудно заметить, что сумма в скобках есть производная по q от суммы:

$$\begin{aligned} S &= 1 + 2q + 3q^2 + \dots + (L - 1)q^{L-2} = \\ &= \frac{d(q + q^2 + \dots + q^{L-1})}{dq} = \frac{d(1 + q + q^2 + \dots + q^{L-1})}{dq} \\ &= \frac{d}{dq} \left(\frac{1 - q^L}{1 - q} \right) \end{aligned}$$

Добавление 1 в сумму не влияет на производную, поскольку производная константы (1) равна 0. Вспоминая математический анализ, легко вычисляем эту производную:

$$S = \frac{-Lq^{L-1} \cdot (1 - q) + 1 - q^L}{1 - q^2} = \frac{-Lq^{L-1}}{p} + \frac{1 - q^L}{p^2}$$

Подставляя в формулу для математического ожидания и вспоминая, что $q = 1 - p$ получаем:

$$E = pS + Lq^{L-1} = \frac{1 - q^L}{1 - q}$$

Из того же матанализа мы знаем, что для $p \ll 1/L$ $(1 - p)^L \approx 1 - Lp$. Поэтому получаем

$$E = \frac{1 - (1 - p)^L}{p} \approx \frac{1 - 1 + Lp}{p} = L = O(L)$$

Таким образом, мы получили оценку времени работы алгоритма — время работы есть $O(L)$. Заметим, что обычно p достаточно мало. Однако, если p велико (т.е. вероятность встречи ключа высока), то время работы будет достаточно малым — при больших p оценка будет $O(p/q)$. Иными словами, если в нашей записной книжке есть только Тютюкины и Пупкины, то поиск Тютюкина быстро приведет к успеху. Надо иметь возможность добавлять элементы в массив и удалять их из него. Добавление элемента происходит просто — надо в конец массива (**если есть место**) дописать элемент. Эта операция требует константного времени $T = O(1)$. Для вставки элемента в заданную позицию надо освободить место, для чего подвинуть все элемента хвоста массива. Это требует времени порядка $O(L)$. При удалении элемента надо часть массива справа от удаляемого элемента сдвинуть влево. Эта операция требует времени, пропорционального длине оставшейся части массива. В среднем — $O(L)$. Итак,

Простой массив	
Добавление	$O(1)$
Вставка	$O(L)$
Удаление	$O(L)$
Поиск	$O(L)$

2.1.2 Поиск в отсортированном массиве

Где нужно слово найти проще — в случайном списке слов или в словаре? А чем хорош словарь? Какое его основное свойство? Слова в словаре *отсортированы, упорядочены*. Сортировать можно только такие массивы, в которых содержатся сравнимые элементы. Объекты называются сравнимыми, если для $\forall x, y$ существует три ситуации: $x < y, x > y, x = y$, причем для операций $>, <$ выполняется условие транзитивности: если $a < b$ и $b < c$, то $a < c$. Заметим, что сравнивать можно не все элементы! Например, не существует естественного способа сравнивать точки на плоскости или в пространстве. Бывают случаи, когда между некоторыми элементами множества можно установить сравнение, а между другими — нельзя. Эти случаи мы рассмотрим при обсуждении алгоритмов на графах. Отсортировать массив в соответствии с правилом сравнения — значит, выстроить его элементы в таком порядке, что

$$\begin{array}{ll} \forall i > 0 & a[i-1] \leq a[i] \quad \text{массив упорядочен по возрастанию} \\ \forall i > 0 & a[i-1] \geq a[i] \quad \text{массив упорядочен по убыванию} \end{array}$$

Слова сравнивают лексикографически. Этот способ сравнения используется, в частности, при построении словарей. Суть его сводится к следующему. Отбрасываем совпадающие части слов сначала, потом сравниваем первую пару букв, которые не совпадают. Какая буква ближе к началу алфавита, то слово и меньше. Если одно слово является началом другого слова, то более короткое слово меньше. Пример: "пингвин" < "прорыв"; Преимущество отсортированного массива в скорости поиска. Вспомним, как ищется слово в словаре. Можно листать словарь сначала, пока не наткнешься на нужное слово. Однако, это не самый эффективный способ. Можно поступить по-другому:

1. Делим массив пополам.
2. Проверяем, попало ли слово на середину, и если да, то Ура!
3. Смотрим, следует ли искать слово в левой или в правой половине.
4. Ищем в соответствующей половине таким же образом.

Покажем работу алгоритма на примере поиска слова "Добро" в словаре на рисунке 2.2. После первого деления пополам находим слово "Клюв". Слово, которое ищем (ключ) меньше, чем "Клюв", поэтому продолжаем поиск в верхней половине словаря. Делим ее пополам и находим слово "Викинг". Это слово меньше ключа, следовательно ищем в части словаря между словами "Викинг" и "Клюв". Заметим, что на каждом этапе мы ищем в части словаря — полный словарь есть частный случай части словаря. Вот алгоритм поиска в части словаря (from - индекс начала, to - индекс конца + 1):

1. B_Search (key, a, from, to){
2. i = int((to+from)/2) ;

Автор	Автор	Автор	Автор
Ангар	Ангар	Ангар	Ангар
Бензин	Бензин	Бензин	Бензин
Викинг	Викинг	Викинг	Викинг
Добро	Добро	Добро	Добро
Квадрат	Квадрат	Квадрат	Квадрат
Клюв	Клюв	Клюв	Клюв
Лень	Лень	Лень	Лень
Пингвин	Пингвин	Пингвин	Пингвин
Порыв	Порыв	Порыв	Порыв
Рассвет	Рассвет	Рассвет	Рассвет
Сказка	Сказка	Сказка	Сказка
Яхта	Яхта	Яхта	Яхта

Рис. 2.2: Пример словаря

```

3.     if (key == a[i])      return(i);
4.     if (to - from < 2)   return(УВЫ...);
5.     if (key>a[i])        return B_Search(e,i+1,to) ; // рекурсия!!
6.     if (key<a[i])        return B_Search(e,from,i-1); // рекурсия!!
7.  }
8.  B_Search (key, a){
9.  return B_Search (key, a, 0, n)
10. }
```

Пояснение кода. Строки 1–7 – поиск в фрагменте массива; строки 8–10 – поиск во всем массиве – весь массив является частным случаем фрагмента массива. В строке 2 происходит целочисленное деление, т.е. остаток отбрасывается, например, $2/3 = 0$. В строках 5, 6 процедура `B_Search` вызывает саму себя. Такой вызов называется *рекурсией*. Замечание: Обязателен выход из алгоритма без рекурсивного вызова! Рекурсия — это аналог метода математической индукции, когда более длинная задача сводится к решению точно такой же, но более короткой. Рекурсивный вызов процедур широко используется в алгоритмах. Например, вычисление факториала можно представить в виде рекурсии:

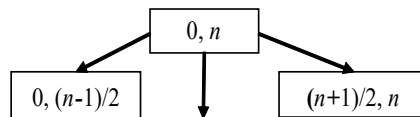
```

1. factorial(n){
2.   if(n==0) return 1;
3.   return n*factorial(n-1)
4. }
```

Упражнение 6. *Напишите алгоритм бинарного поиска без рекурсии в виде цикла.*

2.1.3 Анализ алгоритма бинарного поиска

Чтобы оценить время работы бинарного поиска элемента в отсортированном массиве, представим этапы работы алгоритма в виде узлов и листьев дерева (дерево поиска). Узлам соответствуют части (“половинки”) массива. Корень дерева — сам массив (пусть *from* будет 0 — первый элемент массива, *to* будет $n - 1$ — последний элемент массива). На первом шаге работы алгоритма листьями дерева становятся половинки (без среднего элемента) этого массива: до и после $n/2$:



Здесь центральная стрелка отвечает случаю, когда средний ($n/2$) элемент совпал с ключом, левая и правая стрелки соответствуют случаю, когда средний элемент больше или меньше ключа. В итоге возникает дерево, в котором каждому возможному ключу, соответствует свой путь (рис.2.3). Понятно,

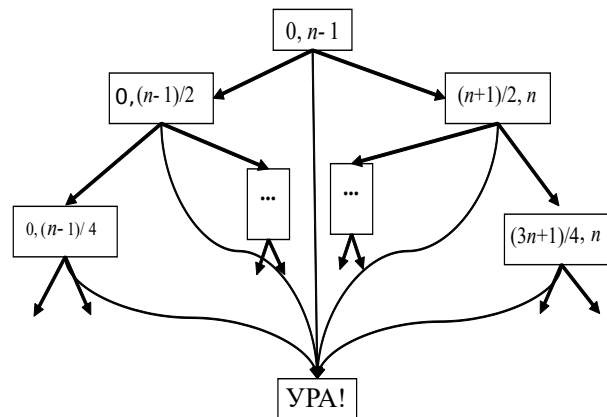


Рис. 2.3: дерево алгоритма поиска элемента в отсортированном массиве

но, что прохождение по этому дереву аналогично работе метода `B_Search`. Каждый следующий этап соответствует узлу, из которого выходят три ветви — три возможности:

1. закончить работу, найдя нужный элемент в середине соответствующей части массива (Ура!)
2. продолжить поиск в начальной половине данной части массива
3. продолжить поиск в конечной половине данной части массива

В худшем случае время работы `B_Search` равно количеству делений массива пополам, т.е. высоте дерева поиска. Какой случай является самым неблагоприятным? Когда мы начинаем из середины массива, а искомый элемент является самым первым или самым последним в этом массиве; в таком случае придется обойти все дерево, прежде чем мы найдем нужный элемент. Высота дерева h зависит от числа листьев в нем. Обратите внимание: наше дерево — "красивое", "правильное" и симметричное, из каждой его вершины выходят три дочерних. У этого дерева число листьев $\leq n$ (n — длина нашего массива). Для h и n имеем такую оценку:

$$2^{h-1} < n \leq 2^h$$

или

$$h - 1 < \log_2 n \leq h$$

Отсюда получаем, что время работы алгоритма $T(n) = O(h) = O(\log n)$. Время работы бинарного поиска — *логарифмическое*, тогда как время работы "наивного" поиска по массиву той же длины — *линейное*. Для массива длиной миллион элементов время работы алгоритма наивного поиска в худшем случае составит миллион операций, а бинарный поиск справится с этой задачей в худшем случае примерно за 20 операций. Однако для работы `B_Search` необходим отсортированный массив. Сортированный массив — хорошая конструкция для поиска, но если нам надо добавить или удалить элемент — возникают проблемы. Для удаления элемента нам надо все элементы с большими индексами подвинуть налево. Вот алгоритм удаления элемента номер k из массива:

```

1.   DelArray(a, k){
2.       for(i=k+1; i < a.length; i++){
3.           a[i-1]=a[i];
4.       }
5.   }
```

Добавление элемента еще хуже. Нам надо, чтобы после добавления элемента массив сохранил свойство упорядоченности. Кроме того, массив имеет заранее определенный размер. Поэтому добавление элемента не должно приводить к переполнению. Если массив полностью заполнен, то вставить элемент уже не удастся (это одно из самых серьезных ограничений в использовании массивов). Если же место для нового элемента есть, то надо, во-первых, найти место куда вставлять элемент (можно использовать бинарный поиск), затем освободить место для вставки (подвинуть весь хвост массива направо), и только после этого вставить элемент.

Упражнение 7. *Модифицируйте алгоритм бинарного поиска так, чтобы он находил либо индекс элемента, если он есть, либо индекс самого большого элемента, который меньше заданного (т.е. тот элемент, после которого надо вставлять новый).*

Упражнение 8. *Напишите алгоритм вставки элемента в сортированный массив.*

2.2 Сортировка массивов

Допустим, есть два отсортированных массива. Можно ли их слить так, чтобы в результате получился сортированный массив? Можно. При этом время работы такого слияния будет порядка суммарной длины входных массивов. Этот алгоритм был предложен Джоном фон Нейманом в 1945 году. Алгоритм называется сортировкой слиянием (Merge_Sort) В его основе лежат следующие соображения. Если:

- массив разбит на два отсортированных подмассива;
- и мы умеем сливать два отсортированных массива.

Тогда:

- разбиваем массив на две половинки;
- сортируем каждую половинку отдельно;
- сливаем половинки.

В основе этого алгоритма лежит парадигма "разделяй и властвуй", состоящая из 3-х этапов:

- *Разделение*: разделение задачи на несколько подзадач.
- *Покорение*: рекурсивное решение подзадач.
- *Комбинирование*: решение исходной задачи исходя из уже решенных подзадач.

Рекурсия достигает своего нижнего предела, когда длина сортируемого подмассива равна 1, а такой массив можно считать упорядоченным. После этого можно сливать массивы. Рассмотрим этот алгоритм, реализованный на псевдокоде:

```
1. MergeSort (A, i1, i2){
2.     if (i2 - i1 < 2) return A;
3.     i = int( (i1+i2)/2 );
4.     A1 = MergeSort (A, i1, i);
5.     A2 = MergeSort (A, i+1, i2);
6.     return Merge (A1, A2);
7. }
8. MergeSort (A){
9.     return MergeSort(A,0,n)
10. }
```

Пояснение кода. Здесь, также как и в бинарном поиске, две части — сортировка части массива (строки 1–7) и сортировка всего массива (строки 8–10). Строка 2 проверяет, что массив имеет длину 1 или 0. Такой массив уже отсортирован, поэтому мы его возвращаем. В противном случае ищем середину (строка 3) и сортируем две половинки массива (строки 4, 5), после чего их сливаем (строка 6). Теперь обратимся к алгоритму слияния массивов. Основную идею слияния двух отсортированных массивов можно объяснить на следующем примере. Пусть имеется две стопки карт, лежащих рубашками вниз так, что в любой момент видно только верхнюю карту в каждой из этих стопок. Пусть также, карты в каждой из этих стопок идут сверху вниз в неубывающем порядке (аналог отсортированных массивов). Как сделать из этих стопок одну? На каждом шаге мы берем меньшую из двух верхних карт и кладем ее (рубашкой вверх) в результирующую стопку. Когда одна из оставшихся стопок становится пустой, мы добавляем все оставшиеся карты второй стопки к результирующей стопке. Псевдокод такого слияния можно увидеть ниже:

```
1. Merge(a1, a2) {
2.     a[]=new array[a1.length+a2.length];
3.     int i1=0;
4.     int i2=0;
5.     while ((i1+i2)< a.length ){
6.         if(i1 < a1.length && i2 < a2.length){
7.             if(a1[i1]< a2[i2]){
8.                 a[i1+i2]=a1[i1];
9.                 i1++;}
10.            else {
11.                a[i1+i2]=a2[i2];
12.                i2++;}
13.            }
14.            if(i1 = a1.length) {
15.                a[i1+i2]=a2[i2];
16.                i2++;}
17.            if(i2 = a2.length) {
18.                a[i1+i2]=a1[i1];
19.                i1++;}
20.        }
21.    return a;
22. }
```

Пояснение кода. В строке 2 мы создаем вспомогательный массив, куда будем записывать результат слияния. Ясно, что длина этого массива равна сумме длин входных массивов. Строки 3 и 4 устанавливают указатели на начало массивов. В строке 6 мы проверяем, есть ли элементы в обоих массивах. Строка 7 — сравнение элементов и перенос соответствующих элементов (строки 8 и 11). При переносе элементов увеличиваются индексы массивов

(строки 9,12). Если же один из массивов исчерпан (строки 14, 17 проверяют это), то переносим элементы оставшегося массива (строки 15, 18). Процедура возвращает новый массив (строка 21).

Упражнение 9. *Продемонстрируйте работу слияния на примере массивов: [3, 7, 8, 10]; [2, 4, 6, 9].*

2.2.1 Анализ алгоритма Merge_Sort

Время работы алгоритма определяется временем работы процедуры Merge. Очевидно, что время работы этой процедуры есть $O(N)$. Действительно, время определяется временем работы цикла строки 5, а он прокручивается ровно N раз, где N — длина результирующего массива. Тогда время работы алгоритма можно представить в виде рекурсии:

$$T(N) = 2 \cdot T(N/2) + O(N) \quad (2.1)$$

Первое слагаемое — сортировка подмассивов (их два, длиной $N/2$), второе слагаемое — слияние

Теорема 1 (Теорема о рекурсии Merge_Sort). *Если*

$$T(N) = 2T(N/2) + O(N) \quad (2.2)$$

то существует C такое, что для любого N

$$T(N) \leq C \cdot N \cdot \log_2 N$$

Условие $T(N) = 2T(N/2) + O(N)$ означает по определению, что существует некоторая константа b , такая, что $T(N) \leq 2T(N/2) + bN$

Доказательство. Доказательство по индукции: База индукции: при $N = 3$ величина T равна чему-то, например, $a \cdot 3 \cdot \log_2 3$ и для любого $C \geq a$ не превышает оценку.

- В предположении индукции: $T(N) \leq C \cdot N \cdot \log_2 N$
- Докажем, что верно для $2N$:

$T(2N) = 2T(N) + O(2N)$	уравнение 1 для $2N$
$\leq 2T(N) + C \cdot 2N$	раскрываем $O(2N)$
$\leq C \cdot 2N \cdot \log_2 2N + C \cdot 2N$	по предположению индукции
$= 2 \cdot N \cdot C \cdot (\log_2 N + 1)$	вынесли $N \cdot C$ за скобки
$= 2N \cdot C \cdot \log_2 2N$	поскольку $1 = \log_2 2$

Таким образом, $T(2N) \leq C \cdot 2N \cdot \log_2 2N$ □

Время работы Merge_Sort $T(N) = O(N \log N)$, однако у сортировки слиянием имеется один недостаток: необходимо создание дополнительного массива, что занимает место в памяти.

Упражнение 10. *Попробуйте создать алгоритм сортировки Merge_Sort без использования дополнительных массивов (один массив все равно понадобится) и без рекурсии.*

2.2.2 Быстрая сортировка (QSort)

Алгоритм, предложенный Чарльзом Хоаром в 1962 году, был назван сортировкой разделением. Метод оказался настолько эффективным, что вскоре его стали называть алгоритмом быстрой сортировки (quick sort, QSort). Сортировка методом QSort основана на одном из свойств упорядоченного массива: Массив отсортирован тогда и только тогда, когда для любого элемента q все элементы массива слева не больше его, а все элемента справа — не меньше:

$$\forall k : (\forall i < ka[i] \leq a[k]) \& (\forall j > ka[j] \geq a[k])$$

Быстрая сортировка, как и сортировка слиянием, основана на принципе "разделяй и властвуй". Для этого берем первый элемент (опорный элемент) в массиве и попробуем сделать перестановки так, чтобы слева от него были элементы меньше, а справа — больше его:

1. Ищем с правого конца первый элемент, который меньше опорного и меняем их местами
2. Ищем слева первый элемент, который больше опорного и меняем их местами
3. Повторяем процедуру до тех пор, пока массив не окажется разбитым на два подмассива при этом опорный элемент занял свое окончательное положение.
4. После разбиения можно применить ту же процедуру к левой и правой частям массива

После первой итерации справа оказывается часть массива, которая больше q , после второй — слева оказывается часть массива, которая не превышает q , а зона массива, для которой неизвестно отношение к q — сокращается. Процедуру повторяем до тех пор, пока эта зона не сожмется до одного элемента — q . После того, как мы разделили массив на две (вообще-то говоря не равные) части можно к этим частям применить ту же процедуру. Теперь рассмотрим реализацию алгоритма на псевдокоде:

```

1.   QSort (A, from, to){
2.       int q;
3.       if (from < to){
4.           q = Partition (A, from, to);
5.           QSort (A, from, q-1);
6.           QSort (A, q+1, to);
7.       }
8.   }
```

Пояснения. Если массив содержит больше одного элемента (строка 3), то его сортируем: перетасовываем массив так, чтобы появилось две части так,

чтобы все элементы в левой части были не больше любого элемента в правой части (строка 4). После этого сортируем по отдельности левую и правую части (строки 5,6). Отметим особенность — левая часть начинается с начала массива (*from*) и заканчивается предыдущим элементом ($q - 1$), а правая начинается со следующего за q элементом, т.е. q больше не трогаем. Метод *Partition* разбивает данный массив на два подмассива и возвращает позицию опорного элемента.

```

1.  Partition (a, from, to){
2.      q=a[from];
3.      i =from;
4.      i1=from;
5.      i2=to-1;
6.      while (i1 < i2){
7.          while ( i2 > i and a[i2] >= q) i2--;
8.          a[i]=a[i2]; a[i2]=q;
9.          i=i2;
10.         while ( i1 < i and a[i1] <  q) i1++;
11.         a[i]=a[i1]; a[i1]=q;
12.         i=i1;
13.     }
14.     return i; // процедура возвращает номер элемента, который
                // обладает нужным свойством: все элементы справа
                // больше, все элементы слева не больше
17. }
```

Пояснения к алгоритму *Partition*. На вход алгоритм получает массив и границы фрагмента массива, который мы будем делить на части. Строки 2-5 производят инициализацию. В переменных $i1$, $i2$ хранятся границы необработанной части массива; в части массива с индексами $< i1$ элементы меньше, чем q , а в части массива $> i2$ все элементы не меньше, чем q . В процессе работы мы сужаем необработанную часть массива. Сначала в правой части ищем элемент массива, нарушающий правило, т.е. элемент, который меньше или равен q (строка 7). Когда его найдем, меняем местами его и опорный элемент (строка 8) и устанавливаем новую позицию опорного элемента (строка 9). Потом ищем “нарушителя” с начала массива и делаем перестановку. Потом продолжаем искать в конце и в начале необработанной части массива, и т.д., пока $i1$ и $i2$ не совпадут (строка 6), что означает, что необработанная часть сжалась до одного элемента. Этот элемент (опорный) занял свое окончательное положение.

На рис.2.4 показан пример работы процедуры *Partition*. Точный анализ алгоритма быстрой сортировки достаточно сложен, приведем здесь лишь основные положения. Среднее время работы алгоритма $T(N) = O(N \log N)$. В худшем случае время работы алгоритма квадратичное ($T(N) = O(N^2)$). Худший случай — когда массив уже отсортирован в прямом или в обратном порядке. Если мы заранее знаем об этом, лучше “испортить” массив,

	3	1	7	4	2	6	⇐
	3	1	7	4	2	6	⇐
⇒	2	1	7	4	3	6	
⇒	2	1	7	4	3	6	
	2	1	3	4	7	6	⇐
	2	1	3	4	7	6	⇐

Рис. 2.4: Пример работы процедуры Partition. Опорный элемент – 3. Сначала ищем справа элемент, который меньше 3. Как только находим – переставляем элементы и начинаем поиск слева элемента, который больше 3.

перетасовав его элементы в случайном порядке (что занимает N операций), и затем применить QSort ($O(N \log N)$). Избежать этого можно предварительно перетасовав массив за время $O(N)$. И даже в этом случае быстрая сортировка будет работать быстрее наивной. Как узнать, хорош ли массив для работы QSort? Упорядоченность массива оценивается по числу инверсий. Случайный, неупорядоченный массив длины N содержит приблизительно $N(N-1)/4$ инверсий (теория вероятностей!). Если число инверсий стремится к нулю – массив отсортирован в возрастающем порядке, если к N – в убывающем. К достоинствам быстрой сортировки относятся:

- Высокая скорость работы.
- Простота реализации.

В отличие от сортировки слиянием (Merge_Sort) алгоритм быстрой сортировки (QSort) одним недостатком. Это – неустойчивость (то есть ячейки с одинаковыми данными могут изменить положение друг относительно друга после сортировки). Этот недостаток не так важен при сортировке массивов чисел, но может оказаться важным при сортировке ключей, которые связаны с некоторыми данными. Упорядоченный массив тоже имеет свои недостатки. Добавление элементов в него требует усилий:

1. Найти место куда вставить ($T(N) = O(\log N)$)
2. Освободить место для вставки ($T(N) = O(N)$)
3. Вставить элемент

Кроме того, нет гарантий, что размер массива позволит вставку элемента. Мы уже поняли, что сортированный массив это хорошо. Но возникает вопрос, что лучше по мере добавления элементов вставлять их на правильное место или сначала собрать массив, а потом отсортировать? Оценим время, требуемое на сборку массива. На каждом шаге надо: (1) найти место для вставки ($T(N) = O(\log N)$); (2) сделать вставку ($T(N) = O(N)$). Итого для

сборки упорядоченного массива необходимо:

$$\sum_i (\alpha \cdot i + \beta \cdot \log i) = \alpha \cdot \frac{N(N+1)}{2} + \beta \cdot \sum_i \log i = O(N^2)$$

Свойства сортированного массива:

Упорядоченный массив	
Вставка	$O(L)$
Удаление	$O(L)$
Поиск	$O(\log L)$

Есть структуры данных лишённые этих недостатков, но обладающие другими, о них читайте ниже.

2.3 Списки

2.3.1 Односвязный список

Всем знаком принцип построения очереди в поликлинике: каждому следующему человеку достаточно знать, кто стоит перед ним, при этом люди не обязательно стоят строго один за другим, они могут, к примеру, отходить. На этом принципе построена такая структура данных, как связный список (односвязный список). Связный список состоит из элементов, каждый из которых несет какие-то данные (*data*). Кроме того, любой элемент имеет ссылку на следующий за ним элемент (*next*). Последний элемент списка ссылается на пустой элемент — *null*. В отличие от элементов массива, части связного списка могут лежать в разных местах памяти, не обязательно один за другим. Элемент связного списка можно описать так:

```

1.     Struct LinkedListItem {
2.         data d;
3.         LinkedListItem next;
4.     }
```

Чтобы определить связный список, достаточно указать первый элемент списка (*root*). На него не ссылается ни один другой элемент. Полезно знать последний элемент (который ссылается на *null*) и размер списка (рис.2.5). Итак, структура связного списка может быть описана:

```

1.     Struct LinkedList{
2.         LinkedListItem root;
3.         LinkedListItem last;
4.         int size;
5.     }
```

Рассмотрим основные операции для связного списка. Поиск. Поиск по связному списку определенных данных (*data d*) ведется от корня. В ходе поиска

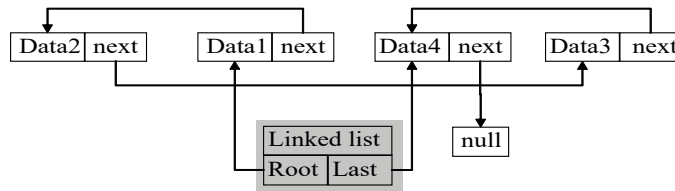


Рис. 2.5: Односвязный список

просматривается по очереди каждый элемент: есть в этом элементе нужные данные? — нет — идем дальше, есть — ура; так до последнего элемента (ссылающегося на null). Следовательно, алгоритм поиска не отличается от такового для массива. Понятно, что и работает он за то же время $T(N) = O(N)$, где N — величина списка

Упражнение 11. *Напишите алгоритм поиска элемента в связном списке.*

Вставка элемента. Плюс связного списка в том, что новый элемент в него можно добавлять всего за одну операцию, независимо от длины списка. Пример из жизни. Вы стоите в очереди в буфет (вы знаете за кем вы стоите). К вам подходит знакомый и он, естественно, встает вместе с вами. Теперь он стоит за тем, за кем стояли вы, а вы стоите за ним. Если мы знаем элемент, после которого нужно вставить новый, то необходимо просто переключить ссылки.

```

1. insertAfter (LinkedListItem elm, LinkedListItem after) {
2.     if(after==null){
3.         elm.next=null;
4.         root= elm;
5.         last=elm;
6.         size++;
7.         return;
8.     }
9.     elm.next = after.next;
10.    after.next = elm;
11.    size++;
12.    if(elm.next==null) last= elm;
13.    return
13. }
  
```

Комментарий к коду . Если мы указываем в качестве параметра *after* *null*, то это означает, что мы вставляем новый корневой элемент (строки 2 - 7), при этом список может быть не пустым, просто мы вставляем в самое начало. В частности, так можно создавать новый список. Иначе мы

просто перебрасываем ссылки (строки 9-12). Очевидно, что вставка элемента требует константного количества операций: $T(N) = O(1)$.

Удаление. Не получится удалить элемент из связанного списка так же просто и быстро, как вставить, поскольку каждый элемент “знает” только следующего соседа, но не предыдущего. Поэтому для удаления элемента придется пройти весь список от *root* до элемента, за которым следует удаляемый элемент:

```

1.  remove(LinkedListItem r){
2.      for(LinkedListItem l=root; l!=null; l=l.next){
3.          if (l.next == r){
4.              l.next=r.next;
5.              if(l==root) root=l.next;
6.              if(l.next==last) last=l;
7.              size --;
8.              return;
9.          }
10.     }
11. return ERROR!!!
12. }
```

Связный список, в отличие от массива, является структурой данных с не фиксированной длиной. Создавая новый список, мы задаем его *root*, к которому в дальнейшем можно добавлять сколько угодно новых элементов один за другим. Обсуждая массивы, мы упоминали, что фиксированная длина, в зависимости от ситуации, может быть как плюсом, так и минусом.

Когда фиксированная длина — плюс? В случае, если не хочется засорять память лишними пустыми элементами, предпочтительнее списки. Пример из Java: известная всем и часто используемая структура данных с не фиксированной длиной *Vector*, или массив в языке Perl, или список в Python. При создании нового экземпляра класса *Vector* занимается определенное количество памяти — около десяти элементов. Даже если в нашем векторе будет храниться только два элемента, остальные восемь останутся недоступными для других программ и объектов. В такой ситуации лучше было бы использовать массив длиной два элемента. Следовательно, “отъедание” памяти “про запас” — возможный минус структур данных с фиксированной длиной. Связный список таким свойством не обладает, он не занимает лишней памяти, поэтому его не фиксированная длина — только плюс. Но надо иметь в виду, что мы к каждому элементу добавили ссылку. Если поле *data* занимает один байт, то мы увеличили требование к памяти в 5 раз (вдобавок к одному байту данных мы храним четыре байта ссылки). Успокаивает одно — обычно поле данных достаточно большое. С другой стороны часто встречаются случаи, когда размер множества заранее известен, например, если вам надо создать какие-то данные, привязанные к

дню недели. Тогда очевидно, что надо иметь массив размером 7. Для многих языков программирования ограниченность размера массива скрыта от программиста. На самом деле все массивы имеют ограниченный размер. Если происходит переполнение размера массива, то система (интерпретатор кода) создает новый массив, например, в полтора раза больше, чем существующий, и переписывает данные из старого массива в новый. Это решение хоть и удобно, но не может считаться эффективным. Еще одним плюсом является локальность данных. Компьютерная архитектура современных компьютеров устроена так, что данные, расположенные рядом, загружаются в память с быстрым временем доступа, поэтому, например, элементы, расположенные вместе, обойти, как правило, быстрее.

Связный список	
Вставка	$O(1)$
Удаление	$O(L)$
Поиск	$O(L)$

2.3.2 Двусвязный список

Односвязный список хорош тем, что в него можно быстро добавлять новые элементы, но, к сожалению, по времени поиска и удаления элемента он не отличается от массива. Быстро удалять элементы не получается, поскольку каждый элемент “знает” только следующего соседа, но не предыдущего. Этому недостатка лишен двусвязный список — структура, аналогичная списку односвязному и отличающаяся только более сложной системой связей. Каждый элемент двусвязного списка хранит данные, а также две ссылки: на следующий элемент и на предыдущий элемент списка. Элемент двусвязного списка можно описать так:

```

1. BListItem{
2.     data d;
3.     BListItem next;
4.     BListItem prev;
5. }
```

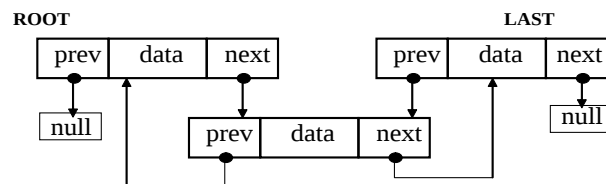


Рис. 2.6: Двусвязный список

Чтобы описать двусвязный список достаточно указать первый элемент `root`. Здесь также полезно хранить последний элемент и размер. Структура двусвязного списка

```

1. Struct BiLinkedList{
2.     BLinkedListItem root;
3.     BLinkedListItem last;
4.     int size;
5. }

```

Заметим, что и односвязный, и двусвязный список могут быть циклическими. Для односвязного цикла должно выполняться условие: $last.next = root$: Для двусвязного: $root.prev = last$; $last.next = root$. В подобных циклах теряется всякий смысл понятий $root$ и $last$, и обход списка можно начинать с любого элемента.

Поиск элемента в двусвязном списке не отличается ни по алгоритму, ни по затрачиваемому времени от поиска в односвязном списке либо в массиве. Двусвязный список можно быстро сортировать методом QSort, который в среднем работает за время порядка $T(N) = O(N \log N)$. Процедура Partition требует возможности быстро менять два элемента местами. В случае списков можно менять не элементы целиком (что дольше и не всегда возможно), только их данные, или просто переключить ссылки. Существуют алгоритмы сортировки двусвязного списка за время порядка $T(N) = O(N \log N)$. К сожалению, быстрый бинарный поиск по списку неосуществим, поскольку мы не можем найти середину списка. В массиве существовал индекс каждого элемента, и получение любого элемента, в том числе среднего, стоило $O(1)$ операций. В списке либо невозможно найти середину, либо, в лучшем случае, это будет стоить $O(N)$ операций, что сведет на нет всю пользу бинарного поиска.

Вставка элемента в двусвязный список, так же, как в односвязный, занимает $T(N) = O(1)$ операций:

```

1. insertElm(BLinkedListItem elm, BLinkedList after) {
2.     if(after==null){
3.         if( root == null )
4.             last=elm
5.         else
6.             root.prev = elm
7.             elm.next = root
8.             root = elm
9.             size++
10.            return
11.    }
12.    if(after.next != null)
13.        after.next.prev = elm
14.    else
15.        last = elm
16.    elm.prev = after

```

```

17.     after.next = elm
18.         size++;
19. }

```

Пояснение. Если мы вставляем перед корнем, то новый элемент должен стать корнем. Если корня нет, то список – пустой и новый элемент станет последним (строка 4), в противном случае у бывшего корня предыдущим элементом станет вновь пришедший (строка 6). Вне зависимости от наличия корня мы устанавливаем следующий элемент и размер. Если же вставка происходит не перед корнем (строки 12-18), то надо проследить за тем, не последним ли является тот элемент, после которого идет вставка. Если он не последний, то у него есть потомок, а у потомка есть предшественник, который был равен after. Этому предшественнику надо установить новую ссылку (строка 13), если же after был последним, то надо переопределить last (строка 15). Потом можно перебросить остальные ссылки.

Удаление элемента из двусвязного списка занимает $O(1)$ операций:

```

1.  remove(BLListItem r){
2.      BLListItem prevElm=r.prev
3.      BLListItem nextElm=r.next
4.      if(prevElm!=null)
5.          prevElm.next=r.next
6.      else:
7.          root=r.next
8.      if(nextElm!=null):
9.          nextElm.prev=r.next
10.     else:
11.         last=r.prev
12.     size--
13. }

```

Итак, основные характеристики двусвязного списка:

Двусвязный список	
Вставка	$O(1)$
Удаление	$O(1)$
Поиск	$O(L)$

2.3.3 Пример применения связного списка. Организация данных в памяти компьютера

Память компьютера линейна, то есть переменные, указатели и т.п. расположены один за другим. Выше мы уже обсуждали организацию памяти с помощью указателей на примере массива. Структура оперативной памяти в современных компьютерах очень сложна. Мы рассмотрим простейшую систему распределения памяти, на которой, к примеру, базировались ранние

версии операционной системы UNIX и DOS. Такая память организована в виде “кусков” разного размера. У каждого такого “куска” есть часть, несущая данные и так называемый контрольный блок (control block, CB). Контрольные блоки напоминают элементы связного списка — каждый из них несет информацию об объеме данных “подопечного” участка памяти и его состоянии (занят/свободен), а также ссылку на следующий контрольный блок (рис.2.7). Как “положить” в такую память, к примеру, массив объе-

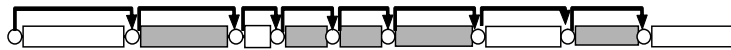


Рис. 2.7: Распределение памяти. Кружки обозначают контрольный блок, белые прямоугольники — свободную память, серые прямоугольники — занятую память.

мом 500 байт? Идем к первому CB спрашиваем: “участок памяти занят?”. Пусть он свободен, тогда спрашиваем: “достаточно ли места для нашего массива?”. Нет, этот участок памяти может хранить только 300 байт. Идем дальше, пока не найдем незанятый участок нужной длины. Если находим участок длины 1000 байт, делим его на два, вставляя новый CB. Если два соседних участка памяти свободны, они объединяются под контролем одного CB. Память, построенная таким образом может привести к феномену “дырявой памяти”. Бывают ситуации, когда объем свободной памяти большой, а положить в нее ничего нельзя — память разбита на очень маленькие кусочки (фрагментация памяти). Современные системы используют улучшенными схемами распределения памяти и управления ею.

Упражнение 12. *Напишите алгоритм распределения памяти. Необходимо две процедуры — `alloc(size)` — запрос на занятие памяти размера `size`. процедура возвращает указатель на предоставленную память. Освобождения памяти — `free(указатель на блок памяти)`. При этом имейте в виду, что `alloc` возвращает не CB, а именно свободное место, а `free` получает указатель на занятую память, а не на CB.*

2.3.4 Очередь и стек

Со списками связаны важные структуры данных — очереди. Они еще больше, чем списки, напоминают настоящую живую очередь. В поликлинике первый пришедший человек (первый добавленный в очередь элемент) первым заходит к врачу (выбывает из очереди). Следовательно, структура данных Очередь (Queue) построена по принципу FIFO — first input, first output. Раньше добавленный в очередь элемент раньше из нее и выйдет.

Очередь — это множество элементов, позволяющее добавлять один элемент и извлекать один элемент. При этом соблюдается правило: первым

извлекается элемент, пришедший первым. Очередь строится на базе односвязного списка, только элементы помнят не предыдущий элемент, а следующий (не так как в жизни — Вы помните ЗА кем Вы стоите в очереди. А не того, кто за Вами). Естественно необходимо знать начало и конец очереди — начало для того, чтобы извлекать из очереди элементы, а конец — для того, чтобы добавлять элементы. Касательно очереди существует две главные процедуры — `void put(elm e)` — добавляет элемент `e` в конец очереди и `elm get()` — берет элемент из начала очереди. В компьютере примером очереди служит очередь задач.

Стек (Stack) , как и очередь, — это множество элементов, позволяющее добавлять один элемент и извлекать один элемент. Однако стек построен по противоположному принципу: первым из него извлекается элемент, добавленный последним (LIFO — last input, first output, вспомните стопку тарелок). Для работы со стеком необходимы две основных процедуры: `void push (elm e)` — кладет элемент `e` в стек и `elm pop()` — извлекает элемент из стека. Стек очень широко используется в организации работы компьютера. Приведем простой пример из жизни. Представьте, что Вы читаете книгу. Вдруг зазвонил телефон. Вы отложили книгу и начали телефонный разговор. При этом Вы вернетесь к книге, когда закончите телефонный разговор, иными словами Вы положили чтение книги в стек. Теперь во время телефонного разговора позвонили в дверь. Вы извинились и поставили телефонный разговор в стек. Вы идете открывать дверь. В этот момент ваш любимый кот решил съесть вашу колбасу. Посетитель тоже отправился в стек. После того, как Вы разобрались с котом Вы извлекаете из стека посетителя и разбираетесь с ним. После ухода посетителя Вы извлекаете из стека то, что там лежит — телефонный разговор. Поговорив, вы снова извлекаете из стека книгу и продолжаете читать. Подобных событий в компьютере происходит множество. Можно привести два примера. Разбор формулы (см. табл. 2.3.4). Пусть дана формула:

$$2 + 3 + (4 * (2 + (2 * 3)) / 3) - 1$$

В этой формуле мы предполагаем, что нет приоритета операций, а порядок вычислений определяется только скобками. Как все это будет вычислять компьютер. Пусть у него есть регистр, в котором будет накапливаться результат. Итак, считываем первое число и помещаем в регистр. Затем считываем операцию и второе число. Производим сложение. Теперь у нас в регистре лежит 5. Далее считываем операцию и видим, что дальше идет скобка. Это сигнал к тому, чтобы содержание регистра и текущую операцию отложить до тех пор, пока эта скобка не закроется, т.е. положить содержимое регистра и операцию в стек. Далее считываем очередное число и операцию. Снова наталкиваемся на скобку — кладем регистр и операцию в стек, и т.д. На каком-то этапе мы встретим закрывающую скобку — это сигнал к тому, чтобы снять со стека регистр и операцию.

Таблица 2.1: пример разбора формулы

прочитано	рег.	опер.	стек	Комментарий
	0			
2	2			
2+	2	+		
2+3	5			2+3=5
2+3+	5	+		
2+3+(0		5+	Кладем регистр и операцию в стек
2+3+(4	4		5+	
2+3+(4*	4	*	5+	
2+3+(4*(0		4*,5+	Кладем регистр и операцию в стек
2+3+(4*(2	2		4*,5+	
2+3+(4*(2+	2	+	4*,5+	
2+3+(4*(2+(0		2+;4*,5+	Кладем регистр и операцию в стек
2+3+(4*(2+(2	2		2+;4*,5+	
2+3+(4*(2+(2*	2	*	2+;4*,5+	
2+3+(4*(2+(2*3	6		2+;4*,5+	2*3=6
2+3+(4*(2+(2*3)	8		4*,5+	снимаем со стека 2+; 2+6=8
2+3+(4*(2+(2*3))	32		5+	снимаем со стека 4*; 4*8=32
2+3+(4*(2+(2*3))/	32	/	5+	
2+3+(4*(2+(2*3))/3	10.67		5+	32/3=10.67
2+3+(4*(2+(2*3))/3)	15.67			снимаем со стека 5+; 5+10.67=15.67
2+3+(4*(2+(2*3))/3)-	15.67	-		
2+3+(4*(2+(2*3))/3)-1	14.67			15.67-1=14.67

Упражнение 13. Напишите структуры и алгоритмы для очереди и стека.

Упражнение 14. Разберите формулу $1+2*(3+4*(5+6*(7+8*9)))$

Упражнение 15. Напишите алгоритм разбора формулы, считая, что за одно считывание читается целиком число или символ операции или скобка и что стек уже реализован и имеет методы *push* и *pop*.

Аналогичная картина происходит, когда вызывается подпрограмма. В этом случае в стеке запоминается точка возврата и состояние текущей программы (все ее внутренние переменные). По окончании работы подпрограммы со стека снимается точка возврата и восстанавливается контекст. В частности при рекурсивном вызове, например программы `B_Search` переменная

i , объявленная в строке 2 алгоритма, появляется еще раз. Ее изменение никак не повлияет на состояние такой же переменной в вызывающей процедуре это просто две разные сущности, лежащие в разных местах памяти.

2.4 Бинарное дерево поиска

Бинарное дерево поиска — это дерево, все вершины которого упорядочены, каждая вершина имеет не более двух потомков (назовем их левым и правым), и все вершины, кроме корня, имеют родителя. Вершины, не имеющие потомков, называются листьями (рис.2.8). Для каждой вершины бинарного

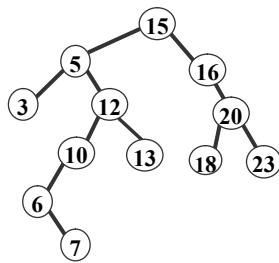


Рис. 2.8: Бинарное дерево поиска

дерева поиска выполняется правило: ключи всех элементов левого поддерева меньше ключа данной вершины, а все ключи в правом поддерева не меньше ключа данной вершины. Иначе это можно записать так:

left and all childs < this, если left существует this \leq right and all childs, если right существует

Какие типы данных можно хранить в виде бинарного дерева поиска? Сразу хочется сказать "числа". На самом деле не только. Можно хранить любые данные, на которых задана процедура сравнения. К примеру, если мы хотим хранить белки, нам нужно условиться, как эти белки сравнивать (по массе, заряду, коэффициенту седиментации и др.), что значит "этот белок строго больше того" (массивнее, более положительно заряжен, медленнее осаждается). Можно также хранить строки. Структура данных может быть представлена в виде:

```

1. BTreeElm{
2.     BTreeElm parent;
3.     BTreeElm left;
4.     BTreeElm right;
5.     Data d;
6. }
  
```

Корень дерева не имеет родителей, его `parent` представляет собой пустой элемент (`null`), соответственно, листья дерева не имеют ни левого, ни правого потомка. Очевидно, что любое поддерево бинарного дерева поиска является бинарным деревом поиска. Основное свойство любого дерева — его высота h . Для бинарного дерева поиска как для структуры данных это очень важная величина: от нее зависит время вставки, удаления и поиска элемента по дереву. В ”хорошем” дереве все эти времена логарифмические: $T_x(N) = O(\log_2(N))$, где x — вставка, поиск, удаление. Почему поиск занимает логарифмическое время, понятно — вспомним бинарный поиск по сортированному массиву. Вставка сводится к поиску нужной позиции в зависимости от вставляемого элемента. Удаление элемента сводится к поиску наименьшего элемента в правом поддереве (см. ниже.)

2.4.1 Поиск в бинарном дереве

Поиск в бинарном дереве осуществляется по следующему алгоритму:

1. Если ключ равен вершине, то Ура!!!
2. Если он меньше вершины, то переходим на левую дочернюю вершину.
3. Иначе переходим на правую дочернюю вершину.
4. Если дерево кончилось, то Увы...

Здесь показан алгоритм поиска в поддереве, начинающемся в `treeElm`:

```

1.  TreeSearch(treeElm, key){
2.      if(treeElm == null) return ‘not found’;
3.      if(treeElm.data == key ) return treeElm;
4.      if(treeElm.data < key)
5.          return TreeSearch(treeElm.right,key);
6.      else
7.          return TreeSearch(treeElm.left ,key);
8.  }
```

Для поиска в дереве надо вызвать эту процедуру для корня: `TreeSearch(root, key)`. Обратим внимание на то, что поиск происходит рекурсивно: если очередной элемент не найден, то ищем в соответствующем поддереве. Время работы $T(N) = O(h)$, где h — высота дерева.

2.4.2 Добавление элемента в бинарное дерево поиска

Бинарное дерево поиска — упорядоченная структура, поэтому ”просто так” элемент в него не добавишь. Для каждого нового элемента нужно искать место так, чтобы и после добавления элемента структура оставалась бинарным деревом поиска. Проще всего добавлять новые элементы в качестве листьев. В зависимости от ключа нового элемента, он вставляется в то или иное поддерево. Для вставки идем по дереву вниз так же, как и при поиске

до первого свободного места, и туда вставляем новый элемент в качестве листа дерева. Можно предложить две формы записи вставки элемента в дерево: рекурсивную и циклическую запись.

Рекурсивная форма:

```
1. Tree_Insert(tree_elm, new_elm){
2.     if(new_elm > tree_elm){
3.         if(tree_elm.right==null){
4.             tree_elm.right=new_elm;
5.             new_elm.parent=tree_elm;
6.             return;
7.         }
8.         Tree_Insert(tree_elm.right, new_elm);
9.     }
10.    else{ //Аналогично для left
11.        if(tree_elm.left==null){
12.            tree_elm.left=new_elm;
13.            new_elm.parent=tree_elm;
14.            return;
15.        }
16.        Tree_Insert(tree_elm.left, new_elm);
17.    }
18. }
```

Вставка с использованием цикла:

```
1. InsertTree(tree_elm, new_elm){
2.     while(true){
3.         if(tree_elm <= new_elm){
4.             if(tree_elm.right == null){
5.                 tree_elm.right=new_elm;
6.                 new_elm.parent=tree_elm;
7.                 return;
8.             }
9.             tree_elm = tree_elm.right;
10.        }
11.        else{ //Аналогично для left
12.            if(tree_elm.left == null){
13.                tree_elm.left=new_elm;
14.                new_elm.parent=tree_elm;
15.                return;
16.            }
17.            tree_elm = tree_elm.left;
18.        }
19.    }
20. }
```

Пояснение: Строчка 2 открывает бесконечный цикл. Мы из этого цикла выходим на строке 7 и на строке 15. Какая форма лучше?

- Рекурсивный тип алгоритма представляется более прозрачным.
- Часто можно записать алгоритм в виде цикла, но не всегда это легко сделать
- Использование цикла предпочтительно, поскольку несколько быстрее, и нет опасности переполнения стека

В любом случае, добавление элемента в дерево стоит $T(N) = O(h)$, где h — высота дерева.

Упражнение 16. *Напишите алгоритм поиска элемента в двоичном дереве в виде цикла.*

Упражнение 17. *Напишите алгоритм двоичного поиска в массиве в виде цикла.*

2.4.3 Поиск минимального (максимального) элемента в дереве

Дерево — это частный случай поддеревы. *Минимальный элемент в поддереве — это самый левый узел, соответственно, максимальный — самый правый узел.* Приведем псевдокод поиска минимального элемента в поддереве, растущем из элемента e :

```

1.  getMin(BTreeElm e){
2.      while(e.left!=null)
3.          e=e.left;
4.      return e;
5.  }
```

Аналогично для максимального элемента:

```

1.  getMax(BTreeElm e){
2.      while(e.right!=null)
3.          e=e.right;
4.      return e;
5.  }
```

Найти минимальный (максимальный) элемент в поддереве удастся за время порядка $T = O(h)$, где h — высота поддеревы.

2.4.4 Поиск следующего по возрастанию элемента

Часто встречается задача вывода элементов какой-либо структуры данных в порядке возрастания (убывания). Чтобы решить эту задачу для бинарного

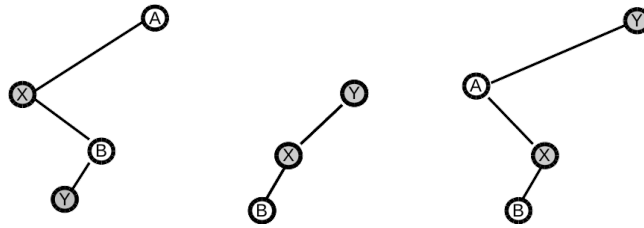


Рис. 2.9: Варианты поиска следующего элемента. X – текущий элемент; Y – следующий элемент. Левый граф – X имеет правое поддерево – следующий элемент – самый левый в правом поддереве; Средний граф – X не имеет правого поддерева, но есть правый родитель; правый граф – нет правого родителя, но есть правый прародитель

дерева поиска, нужно уметь искать следующий (предыдущий) по возрастанию элемент. Для данного элемента следующим будет самый левый элемент в правом поддереве, если такой существует. Такого элемент может и не быть в случае, если текущий элемент не имеет правого поддерева (например, является листом). В этом случае надо искать элемент среди предков. Если предок – левый, т.е. текущий элемент является правым потомком, то предок меньше, чем текущий элемент. Но у этого предка есть еще предки. Пройдя по цепочке предков мы можем обнаружить левого предка, а это значит, что он больше, чем текущий элемент.

```

1.  GetNext(BTreeElm e){
2.      if(e.right != null)
3.          return getMin(e.right);
4.      while(e.parent != null){
5.          if( e = e.parent.left)
6.              return e.parent;
7.          else
8.              e = parent
9.      return null;
10. }
```

Пояснение к псевдокоду. Строка 2 – проверка, есть ли правый ребенок. Если есть (строка 3), то ищем минимальный элемент в правом поддереве. Иначе проверяем родителя и прародителей. Если текущий элемент – левый ребенок, то для данного элемента следующим будет ближайший правый предок (строка 6). Если такового нет, то значит нет следующего элемента. Аналогично ищется предыдущий элемент.

```

1.  GetNext(BTreeElm e){
2.      if(e.right != null)
3.          return getMax(e.left);
```

```

4.  while(e.parent != null){
5.      if( e = e.parent.right)
6.          return e.parent;
7.      else
8.          e = parent
9.      return null;
10. }

```

Упражнение 18. *Напишите программу, которая печатает элементы дерева в возрастающем (убывающем) порядке.*

2.4.5 Удаление элемента из бинарного дерева поиска

В зависимости от расположения элемента в дереве алгоритм удаления будет разным. Есть три основных случая (рис.2.10):

- Элемент детей не имеет. Тогда просто берем и удаляем его (обнуляем ссылку на него у родителя).
- Элемент имеет одного ребенка. Элемент удаляется, а его ребенок усыновляется его родителем (переключаются ссылки)
- Элемент имеет двух детей. Находим следующий элемент (см. `GetNext`). Его удаляем со старого места и вставляем на место удаляемого элемента. (переключаем соответствующие ссылки).

Упражнение 19. *Почему в третьем случае следующий элемент всегда существует*

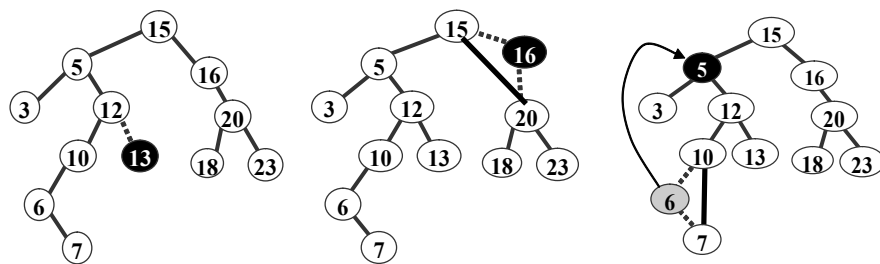


Рис. 2.10: Удаление элемента. Удаляемый элемент показан черным. Пунктиром показаны связи, которые исчезают, толстой линией — новые связи. Слева — удаления элемента без детей, в середине — удаление элемента, имеющего одного ребенка. Справа — удаление элемента, имеющего двух детей.

Псевдокод для удаление элемента:

```

1.  remove(BTreeElm e){
2.      if(e.left ==0 && e.right==0){ // нет детей

```

```

3.     if(e==root){           // удаляем корень, у которог нет детей
4.         root=null;
5.         return;
6.     }
7.     if(e==e.parent.left)   // удаляемый элемент -- левый лист
8.         e.parent.left=null;
9.     if(e==e.parent.right)  // удаляемый элемент -- правый лист
10.        e.parent.right=null;
11.    return;
12. }
13. if(e.parent.left==null){   // единственный правый ребенок
14.     e.parent.right=null;    // удаляем ссылку
15.     return;
16. }
17. if(e.parent.right==null){  // единственный левый ребенок
18.     e.parent.left=null;     // удаляем ссылку
19.     return;
20. }
21. else: //если есть и правый и левый ребенок
22.     y=getMin(e.right); //ищем миним. элемент в правом поддереве
23.     e.data=y.data
24.     remove(y); // удаляем объект
25. }

```

Поясним алгоритм. Идея заключается в том, чтобы перестроить дерево, а затем, если необходимо, переопределить данные на том месте, где был удаляемый элемент. Строки 2–12 — удаляемый элемент — лист. Строки 3–6 удаление корня, у которого нет детей. Дерево становится пустым. Строки 13–20 — удаляемый элемент имеет единственного ребенка. Строки 21–24 два ребенка. В этом случае ищем следующий по возрастанию элемент. Пересаживаем данные на новое место (строка 23) и удаляем этот элемент (строка 24). Можно (и лучше) вместо переноса данных переключить ссылки. С другой стороны, если в поле данных лежат не сами данные, а указатель на них, то предложенный вариант оптимален. Подводя итог алгоритму удаления элементов, заметим, что в худшем случае удаление элемента сводится к поиску наименьшего элемента в правом поддереве, что занимает $O(h)$, где h — высота дерева. В двух других случаях удаление элемента будет стоить еще меньше — $O(1)$.

Бинарное дерево поиска		
Задача	Время	
	средн.	худш.
Вставка	$O(\log(n))$	$O(n)$
Удаление	$O(\log(n))$	$O(n)$
Поиск	$O(\log(n))$	$O(n)$

2.4.6 Преобразование деревьев поиска

Нужно помнить, что логарифмические времена поиска, вставки, удаления возможны только при условии "хорошего" дерева. Время в худшем достаточно велико, поскольку в худшем случае дерево не сбалансировано (рис.2.11). Важно отметить, что оба дерева на этом рисунке соответствуют одному и тому же набору данных! Ясно, что для того, чтобы дерево поиска было эффективным, необходимо, чтобы оно имело небольшую высоту, и поэтому сильно ветвилось. Отметим, что деревья, представленные на рисунке содержат один и тот же набор данных. Это наводит на мысль, что, во-первых, существуют классы эквивалентности деревьев поиска, а, во-вторых, существуют преобразования, переводящие одно дерево поиска в другое, эквивалентное ему. Тогда, применив эти преобразования можно *сбалансировать* дерево. Таким преобразованием является вращение дерева.

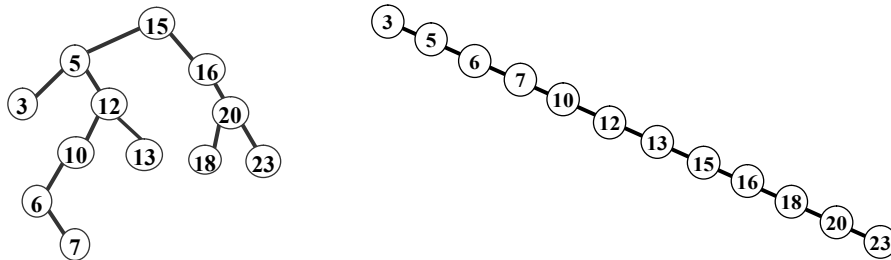


Рис. 2.11: Слева показано хорошо сбалансированное дерево (много узлов с двумя потомками), справа — несбалансированное дерево.

Дано: некое бинарное дерево поиска с поддеревьями (см. рис.2.12) a , b и c . Узел Y имеет в качестве дочерних поддеревья a слева и b справа. Узел X имеет дочерние узлы: Y слева и c справа.

Требуется: поменять X и Y местами.

Решение: Так как данное дерево является бинарным деревом поиска, то $Y < b < X$. Поэтому X может входить в правое поддерево вершины Y , а b — в левое поддерево вершины X , что не нарушит основное свойство бинарного дерева поиска. В связи с этим перестраиваем дерево таким образом (см. Рис.2.12).

Псевдокод вращений.

```

1. LeftRotate(X){
2.   Y=X.left;
3.   X.left= Y.right;
```

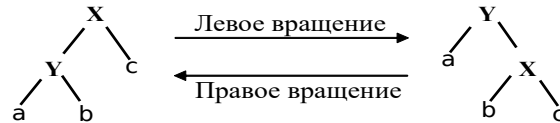



Рис. 2.12: Вращение деревьев поиска

```

4.   if(X.parent !=null){
5.       if(X.parent.left==X)
6.           X.parent.left=Y;
7.       else
8.           X.parent.right=Y;
9.   }
10.  Y.parent=X.parent;
11.  X.parent=Y;
12.  Y.right=X;
13. }

```

```

1.  RightRotate(Y){
2.      X=Y.right;
3.      Y.right=X.left;
4.      if(Y.parent != null){
5.          if(Y.parent.right==Y
6.              Y.parent.right=X;
7.          else
8.              Y.parent.left =X;
9.          }
10.     X.parent=Y.parent;
11.     Y.parent=X;
12.     X.left=Y;
13. }

```

2.4.7 Красно-черные деревья

Определение 1. *Дерево называется красно-черным (Red-Black tree, рис.2.13), если выполняются следующие условия:*

1. *дерево является бинарным деревом поиска*
2. *любая вершина имеет цвет — красный или черный*
3. *добавлены фиктивные листья, всегда черные. Поэтому любая настоящая вершина имеет двух детей*
4. *красная вершина имеет только черных детей*

5. по пути от корня до любого листа количество черных вершин одинаково. Это число называется черной высотой дерева (bh).

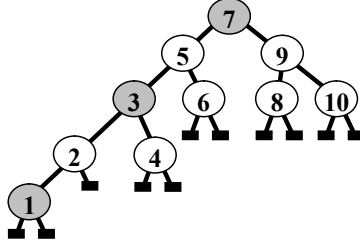


Рис. 2.13: Пример красно-черного дерева. "Красные" вершины отмечены серым. "черные" вершины отмечены белым. Чему равна черная высота?

Из свойства (4) следует, что на любом пути длины l не может быть больше $l/2$ красных вершин. У красно-черного дерева любое поддерево красно-черное.

Лемма 1 (о высоте). *Полная высота красно-черного дерева не превышает два логарифма числа вершин:*

$$h \leq 2 \cdot \log_2(n + 1)$$

где h — высота дерева, n — число вершин дерева.

Доказательство. Докажем лемму методом математической индукции.

Основание индукции: для $n = 1$ утверждение верно: $1 < 2$

Предположение индукции: Пусть для любого поддерева (левого и правого) утверждение леммы верно, т.е. $h \leq 2 \cdot \log_2(n + 1)$, или $n \geq 2^{h/2} - 1$.

Индукция: Пусть n_r — число вершин в правом поддереве, а n_l — в левом. $n = n_l + n_r + 1$ (число вершин дерева в точности равно сумме числа вершин в правом и левом поддеревьях плюс корень). Высота каждого поддерева не меньше, чем $h - 1$, иначе не выполнялось бы свойство одинаковости черной высоты по любому пути. По предположению индукции для n_r и n_l верны оценки:

$$n_r \geq 2^{(h-1)/2} - 1, n_l \geq 2^{(h-1)/2} - 1$$

Отсюда получаем:

$$\begin{aligned} n &= n_l + n_r + 1 \geq 2^{(h-1)/2} - 1 + 2^{(h-1)/2} - 1 \\ &= 2 \cdot 2^{(h-1)/2} - 1 = 2^{(h-1)/2+1} - 1 = 2^{(h+1)/2} - 1 \\ &\geq 2^{h/2} - 1 \end{aligned}$$

или $h \leq 2 \cdot \log_2(n + 1)$

□

Таким образом, время поиска, вставки и т.п. в красно-черное дерево всегда логарифмическое. Правда, после вставки или удаления может нарушиться одно из условий определения красно-черного дерева. Заметим, что "плохих" красно-черных деревьев не бывает. Нельзя сделать цепочку из красных вершин, а длинная цепочка из черных нарушит черную высоту. В этом преимущество красно-черных деревьев.

Поиск, вставка и удаление элементов из красно-черного дерева

Поиск элемента. Поиск элемента в красно-черном дереве ведется точно так же, как в обычном.

Добавление элемента. Новый элемент в черно-красное дерево добавляется так же, как и в любое бинарное дерево поиска и красится в красный цвет (чтобы не нарушать черную высоту). Если родителем нового элемента является черный элемент, ничего больше делать не нужно: черная высота не нарушена, дерево осталось красно-черным. Однако если родитель нового (красного) элемента тоже красный, возникает конфликт. Чтобы разрешить его и оставить дерево красно-черным, придется преобразовывать дерево — вращать и перекрашивать вершины. За добавлением новой красной вершины, пришедшейся "не к месту", следует поэтапная перекраска дерева, сопровождающаяся вращениями. У элемента могут быть "дальний" и "ближний" дядя. Проще всего объяснить, чем ближний дядя отличается от дальнего, на рисунке 2.14: В зависимости от цвета "дядей" добавленного элемента

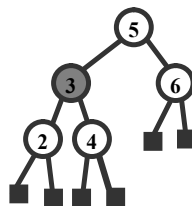


Рис. 2.14: Схема красно-черного дерева. Кружками показаны настоящие, квадратиками — фиктивные листья. Черная высота дерева — 3. Элемент 6 для элемента 4 является *ближним дядей*, а для элемента 2 — *дальним дядей*.

указанные элементы перекрашиваются определенным образом. К примеру, если дальний дядя пришедшего элемента черный, вращаем дерево на уровне родителя нового элемента влево (рис. 2.15). Разрешение конфликта в красно-черном дереве можно реализовать путем вращения (максимум два раза). При этом его черная высота сохраняется, а коллизия поднимается выше. В конце концов коллизия исчезнет, поскольку мы достигнем корня. Корень красим в черный цвет, и конфликт исчерпан. На предыдущих этапах

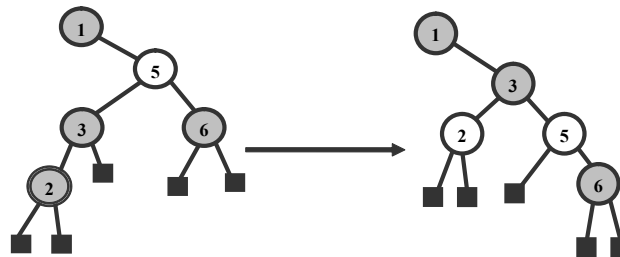


Рис. 2.15: Один из примеров перемещения конфликта. Добавлен элемент 2. После вращения элемент 2 перекрашивается. Основное свойство сохранено.

мы этого сделать не могли, поскольку нарушился бы баланс между левым и правым поддеревьями у корня. Существует много различных вариантов вращения и перекрашивания в зависимости от цвета родителя и “дядей”. Детали этого процесса мы опустим и отправим читателя к специальной литературе. В любом случае смысл изменений дерева состоит в перемещении конфликта (красная вершина имеет красного ребенка) к корню дерева, где его можно спокойно разрешить.

Удаление элемента. Удаление красной вершины ничего не меняет, удаление черной вершины влияет на черную высоту, и кроме того, может привести к появлению красного ребенка у красного родителя. Здесь также возникает несколько случаев. Удаляемый элемент – настоящий лист. Если у него нет “брата”, то он неизбежно красный (почему?) и его удаление не имеет проблем – делаем как обычно. Если он черный и у него есть “брат”, то брат – черный (почему?). Тогда удаляем элемент, брата перекрашиваем в красный. Если родитель был красным, то делаем его черным и все в порядке. Но, если родитель был черным, то делаем его дважды-черным и при вычислении черной высоты будем ее считать дважды. Таким образом основное свойство красно-черного дерева сохраняется. Что с делать с дважды-черной вершиной – разберемся позднее. Удаляемый элемент имеет одного настоящего ребенка. Тогда этот ребенок – красный, а удаляемый элемент неизбежно черный. Пересаживаем на его место ребенка и красим его в черный цвет – все в порядке. Если удаляемый элемент имеет двух настоящих детей, то поступаем как обычно – находим минимальный элемент в правом поддереве, удаляем его из дерева (может быть получив дважды-черную вершину) и ставим на место действительно удаляемой вершины. Цвет определяем по цвету удаляемой вершины. Теперь – что делать с дважды-черной вершиной. С помощью вращений и перекрашивания мы можем поднять ее вверх до корня. тогда корень можно безбоязненно перекрасить из дважды-черной в просто черную – и все правила красно-черного дерева восстановятся.

Красно-черное дерево	
Вставка	$O(\log(n))$
Удаление	$O(\log(n))$
Поиск	$O(\log(n))$

2.5 Хеш-таблицы

Обращение к элементу массива происходит за время $O(1)$. Представим себе, что мы как-то пронумеровали ВСЕ возможные ключи. Тогда процедура добавления, поиска и удаления элемента была бы простой — достаточно по ключу определить номер и посмотреть есть ли в соответствующем элементе массива данные, или положить новые данные, или удалить данные, которые там уже лежат. Беда только в том, что, как правило, множество возможных ключей очень велико и нам никогда не удастся построить массив такого размера с массой пустых ячеек. Тем не менее, допустим, у нас есть отображение g множества ключей на множество целых неотрицательных чисел:

$$g : \{key\} \rightarrow N$$

Это отображение называется хеш-функцией. Во многих случаях можно построить такое отображение. Например, есть простой способ отображения строк на пространство целых чисел. Каждый символ кодируется, скажем, байтом. Тогда можно написать:

$$g(s) = s[0] + 256 \cdot s[1] + 256^2 \cdot s[2] + \dots$$

Разумеется, ни для какой строки длиннее 4 мы не сможем уместить такое число в переменную типа `int` и нельзя даже представить себе массив, способный хранить целочисленные образы всех строк. Если же вместо отображения на пространство всех чисел использовать отображение (не обязательно взаимно-однозначное) на небольшое подмножество множества целых неотрицательных чисел:

$$h : \{key\} \rightarrow 0, 1, \dots, m - 1$$

то для хранения данных можно было бы использовать массив размера m . Данные можно хранить в простой таблице, которая каждому значению хеш-функции ставит в соответствие какие-то данные. Простейшим вариантом такого ограниченного отображения будет

$$h(key) = g(key) \pmod m$$

Однако, поскольку несколько ключей могут соответствовать одному хеш-значению, возможно возникновение *коллизий*, т.е. ситуаций, когда несколько ключей захотят записаться в одну клетку таблицы. Для выхода из этой ситуации существует 2 способа. Первый заключается в том, что в ячейках таблицы лежат не просто ключи, а связный список ключей. Такой способ разрешения коллизий называется закрытой адресацией.

Добавление элемента заключается в вычислении хеш-функции, затем происходит обращение к ячейке таблицы, и, если она пустая, то создается связный список из одного элемента, если там список уже есть, то добавляем элемент к списку. Время $O(1)$ - время добавления в связный список.

Поиск. Вычисляем хеш-функцию, и обращаемся к соответствующей ячейке таблицы. Далее ищем в связном списке, лежащем в этой ячейке. Время $O(1+l)$, где l - длина списка.

Удаление. Сначала находим элемент в связном списке, затем его удаляем из связного списка. Время $O(1+l)$, где l - длина списка. Эффективность работы хеш-таблицы определяется хеш-функцией. Хорошая хеш-функция должна удовлетворять двум основным требованиям: она должна быстро вычисляться, и она должна порождать хорошие ключи для равномерного распределения элементов по таблице. Поговорим о втором условии. С введением хеш-функций решилась проблема переполнения памяти: хеш-значений ограниченное число (m), а проблема коллизий решена с помощью связного списка. Обычно длина такого списка невелика и не превышает числа строк в хеш-таблице. Если хеш-функция достаточно хорошая, списки в ячейках короткие. Иначе элементы с большой вероятностью попадут в одно и то же место таблицы. Приведем примеры хеш-функций:

- Вырожденная функция с одним хеш-значением :

$$h(key) \equiv 1$$

Понятно, что такая функция неприменима. В одной ячейке с хеш-значением 1 будет находиться длинный список из всех элементов. Это сводит хеш-таблицу к связному списку.

- Если у нас определена функция отображения множества ключей на множество неотрицательных целых чисел, то можно предложить такую хеш-функцию:

$$h(key) = g(key) \% m$$

где $\%$ обозначает остаток от деления. В таком случае, чем больше число m , тем короче будут списки в каждой ячейке. Понятно, что принимает значения от 0 до m , поэтому чем меньше число, тем меньше возможных хеш-значений, тем больше элементов соответствует каждому из них. Однако делать m слишком большим нельзя — это элиминирует все преимущества хеш-таблицы над простой таблицей ключей.

Вопрос построения эффективных хеш-функций является предметом достаточно сложных математических исследований. Важнейшим свойством правильной хеш-функции является равновероятность появления хеш-значений на множестве ключей. Поэтому для построения правильной хеш-функции необходимо знать распределение вероятностей появления тех, или иных ключей. Если хеш-функция равномерная, то эффективность операций: добавление,

удаление и поиск в среднем определяется коэффициентом заполненности таблицы $\alpha = n/m$, где n — количество объектов, а m — количество значений хэш-функции. $h: T = O(1 + \alpha)$. Единица в оценке — определяет время вычисления хеш-функции. Вообще говоря, это время зависит от размера ключа. Например в случае если в качестве ключей используются строки, то время вычисления хеш-функции может зависеть от размера строки (ключа). Добавление в хэш-таблицу занимает $O(1)$, так как это время добавления в конец связанного списка.

Второй способ заключается в том, что при возникновении коллизии (ячейка хэш-таблицы занята) просматриваются последующие ячейки до тех пор, пока не найдется пустая. При этом если таблица закончилась переходим в первую ячейку и продолжаем поиск пустой. Решением проблемы полностью заполненной таблицы может быть динамическое увеличение размера хеш-таблицы, с одновременной ее перестройкой. Данный вид хэш-таблиц получил название открытой адресации. Для нее коэффициент заполненности $\alpha = n/m$, обозначается несколько иначе: n — число объектов вставленных в таблицу, а m — общее число ячеек в ней.

Упражнение 20. *Напишите основные алгоритмы работы с хеш-таблицами.*

Связанная с биоинформатикой хеш-функция: отображает сочетание двух букв, обозначающих нуклеотиды, в число от 0 до 15

хеш-функция			
aa	0	ga	8
ac	1	gc	9
ag	2	gg	10
at	3	gt	11
ca	4	ta	12
cc	5	tc	13
cg	6	tg	14
ct	7	tt	15

Пусть данные `data` — это адрес начала определенного динуклеотида в нуклеотидной последовательности, а ключ `key` — этот динуклеотид. Тогда при данной последовательности

```

g a t c g c t c c g t c t g c t a c t g c
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

```

Хеш-таблица будет иметь такой вид:

hash	data	hash	data
0		8	0
1	16	9	4 ↔ 13 ↔ 19
2		10	
3	1	11	9
4		12	15
5	7	13	2 ↔ 6 ↔ 10
6	3 ↔ 8	14	12 ↔ 18
7	11 ↔ 14 ↔ 17	15	

Рассмотрим один пример. Заполнение таблицы. Первому динуклеотиду **ga** соответствует хеш-значение 8. Кладем в таблицу адрес 0. Второму динуклеотиду **at** соответствует ключ 3. Кладем в таблицу адрес 1. И так далее, пока не дойдем до конца последовательности. Если возникает коллизия (например, для динуклеотида **tc**), то организуем связный список. Третья ячейка в ней оказалась пустой, в некоторых ячейках возникла коллизия. Длина последовательности небольшая, поэтому и в каждую ячейку попало всего по несколько элементов. Число хеш-значений данной хеш-функции ограничено (16), поэтому при работе с реальными последовательностями длина связных списков в каждой ячейке будет велика. Теперь эту конструкцию можно использовать для поиска. Если мы хотим найти динуклеотид **ct** (его позиция). Мы берем ключ (строку "ct", обозначающую динуклеотид) и находим ее хеш-значение (7). В ячейку хеш-таблицы с хеш-значением 7 находим связный список со всеми координатами. В связи с этим в реальной работе предпочтительнее использовать хеш-функции с большим числом значений. При работе BLAST задействована хеш-функция, отображающая не ди-, а 12-нуклеотиды на пространство целых чисел. У такой функции 4^{12} значений, что оптимально для работы с большими нуклеотидными последовательностями. В случае аминокислотных последовательностей используется хеш-таблица, основанная на триплетках аминокислотных остатков. Современный BLAST использует достаточно умную схему хранения данных и схему хеширования, которые позволяют достичь высокой эффективности.

2.6 Заключительные замечания

В таблице 2.2 показаны сравнительные характеристики рассмотренных структур данных, в том числе "накладные расходы" на хранение дополнительных указателей, максимальные минимальные и средние времена выполнения операций. Замечание: $(n/2) \cong O(n)$

Таблица 2.2: Времена работы основных действий для различных структур данных

	Доп. указат.	Поиск			Добавление			Удаление		
		сред	Max	Min	сред	Max	Min	сред	Max	Min
Простой массив	0	n	n	1	1	1	1	n	n	1
Сортированный массив	0	$\log n$	$\log n$	1	n	n	1	n	n	1
Односвязный список	1	n	n	n	1	1	1	n	n	1
Двусвязный список	2	n	n	n	1	1	1	1	1	1
Двоичное дерево поиска	3	$\log n$	n	1	$\log n$	n	1	$\log n$	n	1
Красно-черное дерево	3 + 1 бит	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$
Хеш-таблица	1	$1 + \alpha$	n	1	1	1	1	$1 + \alpha$	n	1

Глава 3

Алгоритмы для строк

3.1 Определения, постановка задачи. Наивный алгоритм

В биоинформатике (и не только) очень часто встречается задача поиска подстроки. Сначала дадим несколько определений.

Строка S — массив символов s_0, s_1, \dots, s_{n-1} , принадлежащих некоторому алфавиту A . *Подстрока* $S[i, j]$ — фрагмент массива $s_i, s_{i+1}, \dots, s_{j-1}$, размер подстроки меньше или равен размеру текста. *Префикс* $S[i]$ — подстрока $S[0, i]$. *Суффикс* $S[i]$ — подстрока $S[i, n]$. Очевидно, что строка является своим префиксом и суффиксом. Это тривиальные префикс и суффикс (см. рис. 3.1). Любой нетривиальный префикс или суффикс называется еще *собственным*.

Префиксы	{	a a b a	}	Суффиксы
		a a b a s a		
		a a b a s a s b b a		
		a a b a s a s b b a s a b a s		
		a s a s b b a s a b a s		
		a s		

Рис. 3.1: Префиксы суффиксы. Верхние 3 строки показывают префиксы, нижние — суффиксы

Задача поиска образца. Даны две строки T (текст) и P (образец, паттерн). Нужно узнать, встречается ли подстрока, такая же, как P , внутри T , иными словами, найти хотя бы одно точное вхождение P в T . Существует также задача поиска числа точных вхождений P в T . Образец P встречается в тексте T со сдвигом s , если $T[s + 1, \dots, s + m] = P[1, \dots, m]$. Если

P встречается в T со сдвигом s , то s называется допустимым, или корректным, сдвигом. Итак: *Дано*: длинная строка $T = t_0, t_2, \dots, t_{n-1}$ (текст) и короткая строка (образец, паттерн) $P = p_0, p_2, \dots, p_{m-1}$, $m \leq n$. *Требуется найти*: непрерывный фрагмент строки T , полностью совпадающий с P , т.е. подстроку $T[i, i + m]$, такую что $t_i = p_0, t_{i+1} = p_1, \dots, t_{i+m-1} = p_{m-1}$. Наивный алгоритм: идем по строке T , сдвигаясь каждый раз на одну позицию, проверяем, не начало ли это нужной нам подстроки либо до ошибки, либо до конца строки. Ниже приведена процедура `SubStrCmp` — процедура сравнения подстроки $T[i, i + m]$ с образцом P . Если эта подстрока равна P , возвращается `true`, если хоть один символ отличается — возвращается `false`

```

1. SubStrCmp (T, i, P){
2.     for (j=0; j<= Length(p); j++){
3.         if (t[i+j] != p[j]) return false;
4.     }
5.     return true;
6. }
7. NaiveSearchPattern (T, P){
8.     for(int i=0; i<Length(t)-Length(p)+1; i++){
9.         if (SubStrCmp(T, i, P) = True) return i;
10.    }
11.    return ‘‘Not found’’;
12. }
```

Описанный алгоритм работает просто, но долго. В наихудшем случае наивный алгоритм поиска образца по строке работает за $O(m \times n)$ операций. Таких случаев можно придумать несколько, например для строки и образца

```

T=“aaaaaaaaaaaaa”
P=“aaaaax”
```

T не содержит ни одного вхождения P , но включает много похожих подстрок, отличающихся от P только последним символом. При поиске процедура `SubStrCmp` постоянно проходит почти всю длину, обрываясь постоянно на последнем символе, что занимает много времени. С учетом того, что роль строки T в биоинформатических задачах часто играют нуклеотидные последовательности, не отличающиеся особой краткостью, время работы алгоритма может быть велико. Однако, если призвать здравый смысл, то картина может оказаться еще более удивительной. Вопрос: насколько будет отличаться время поиска слова, например, "cat" в геноме плазмиды rUC19 (2686 пар оснований) и в геноме человека (ок. 3 миллиардов пар оснований). Очевидно, что и там и там слово найдется очень быстро и время поиска не зависит от длины текста! Отметим, что здесь речь идет о поиске первого вхождения слова. Если же нам надо найти все вхождения, то время будет заведомо пропорционально длине текста — просто чем длиннее текст, тем больше слов мы найдем. Оценим *среднее время* работы алгоритма на строке

3.1. ОПРЕДЕЛЕНИЯ, ПОСТАНОВКА ЗАДАЧИ. НАИВНЫЙ АЛГОРИТМ69

T длиной n с образцом P длиной m . Среднее время поиска равно

$$\begin{aligned} & (\text{количество вызовов } \text{SubStrCmp} = n - m) \times \\ & (\text{среднее время исполнения } \text{SubStrCmp}). \end{aligned}$$

Для оценки среднего времени работы процедуры `SubStrCmp` можно применить схему Бернулли. Элементарным событием будет служить сравнение каждого следующего символа строки T с каждым следующим символом образца P . При этом условие независимости элементарных событий выполнено: каждое следующее сравнение не зависит от предыдущего. У любого события возможны только два исхода: успех (символы совпали) и неудача (символы не совпали). Среднее время работы `SubStrCmp` зависит от параметра p — вероятности успеха, равной для каждого элементарного события. Пусть:

$$\begin{aligned} p & \text{ — вероятность совпадения символов} \\ q & = 1 - p \text{ — вероятность несовпадения символов} \end{aligned}$$

Время работы процедуры зависит от числа шагов. Работа процедуры обрывается при первом несовпадении, и чем больше совпадений идет подряд до первого несовпадения, тем дольше работает процедура. Вероятность того, что работа `SubStrCmp` обрывается на первом шаге, равна q (вероятности несовпадения символов в первой позиции). Вероятность обрыва на втором шаге равна pq (чтобы оборваться на втором шаге, нужно пройти первый шаг, то есть на первом шаге должно быть совпадение, вероятность чего равна p). Таким образом, среднее время работы процедуры равно $q + 2pq + 3p^2q + \dots + (m-1)qp^{m-2} + mp^{m-1}$. Подобную сумму мы уже вычисляли при анализе поиска в массиве:

$$E = \frac{1 - p^m}{1 - p}$$

Предположим, что появление букв в тексте и в образце равновероятны. Тогда вероятность совпадения символов из строки T и образца P равна $1/|A|$ где A — алфавит, к которому относятся символы, входящие в T и P , $|A|$ — мощность алфавита A . Если алфавит большой, значение p маленькое. Тогда величиной p^m можно пренебречь, поэтому средняя скорость работы `SubStrCmp` имеет порядок

$$T(\text{SubStrCmp}) \approx \frac{1}{q} = O(1)$$

При выводе этой формулы мы делали следующие предположения:

- Строки T и P независимы
- Алфавит A имеет состоит из достаточно большого числа символов (имеет большую мощность)

Поскольку нам надо применить процедуру `SubStrCmp` n раз, среднее время наивного поиска, равно $O(n)$. Из выведенной формулы понятно, что чем больше мощность алфавита и чем меньше длина строки, в которой мы ищем, тем быстрее (в среднем) идет поиск. Однако в реальности далеко не всегда частоты встречи символов одинаковые, кроме того строки могут быть и плохими (как в приведенном примере), т.е. условия поставленной задачи не обязательно удовлетворяют данным предположениям. Поэтому нужно искать другие алгоритмы.

3.2 Алгоритм Рабина-Карпа

Пусть алфавит состоит из цифр. Тогда любую строку можно представить в виде целого числа. Для того, чтобы сравнить две строки надо просто сравнить эти целые числа. Если строки не слишком длинные, то сравнение двух чисел можно сделать за одну операцию. Надо только преобразовать строку цифр в число:

$$N(\alpha_0\alpha_1\alpha_2 \dots \alpha_{n-1}) = \alpha_0 \cdot A^{n-1} + \alpha_1 \cdot A^{n-2} + \dots + \alpha_{n-1} \cdot A^0$$

где A — размер алфавита. Например, преобразование строки "3243" в число будет выглядеть так:

$$N(3243) = 3 \cdot 10^3 + 2 \cdot 10^2 + 4 \cdot 10^1 + 3 \cdot 10^0$$

Алгоритм поиска подстроки можно представить теперь так: преобразовываем образец P в число. Потом каждую подстроку текста T преобразовываем в число и сравниваем числа. Видно, что преобразование последовательности цифр в число требует времени $O(L)$, где L — длина последовательности. Поэтому пока идея не работает — общее время поиска образец будет требовать времени порядка $O(m \times n)$ (как обычно, m — длина образца, n — длина текста). Однако есть одно соображение, которое позволяет сделать алгоритм более эффективным. Пусть нам известно преобразование подстроки, начинающейся с позиции i : $N(i) = N(t_i t_{i+1} \dots t_{i+m})$. Тогда можно легко вычислить $N(i+1)$. Для этого надо отбросить первую цифру и дописать цифру после последней:

$$t_0 t_1 \dots t_i \overbrace{t_{i+1} \dots t_{i+m-1} t_{i+m}}^{N(i)} \dots t_{n-1}$$

$N(i+1)$

Чтобы отбросить первый символ числа-строки надо взять число по модулю A^{m-1} , чтобы добавить символ справа надо умножить число на A и прибавить число, соответствующее добавляемому символу:

$$N(i+1) = (N(i) \bmod |A|^{m-1}) \cdot |A| + N(t_{i+m}) \quad (3.1)$$

Таким образом алгоритм становится более эффективным. Сначала вычисляем числа, соответствующие образцу P и префиксу текста $T[1, m]$. Затем

в цикле сравниваем числа, потом сдвигаем в соответствии с 3.1 и повторяем сравнение.

```

1. int StrToNum(String p, int l){
2.     n=0;
3.     for(i=0; i<l; i++){
4.         n=n*A+toInt(p[i]);
5.         return n;
6.     }
7. }
8. int RabinMatch(String t, String p){
9.     M=StrToNum(p,m);
10.    N=StrToNum(t,m);
11.    nn=power(A,m-1);
12.    for (i=0; i < n - m ; i++){
13.        if(M==N) return i;
14.        N=(N mod nn)*A + toInt(t [i]);
15.    }
16.    return ‘‘Not found’’;
17. }

```

Пояснения к алгоритму. Строки 1-7 вычисляют число по строке длины l . A — размер алфавита. Строки 9 и 10 вычисляют числа для образца P и префикса $T[0, m]$. Дальше в цикле происходит сравнение чисел (строка 13) и пересчет числа (строка 14). Здесь мы использовали две функции `toInt(a)` — преобразование символа в целое число и `power(a, b)` — возведение числа a в степень b . Описанный алгоритм, носит имена Рабина и Карпа. Нетрудно заметить, что время поиска не зависит от длины образца, каждое следующее число получается за время порядка $O(1)$ и поэтому время работы имеет порядок $T = O(m + n)$.

Ограничения работы алгоритма: Алгоритм нельзя использовать на больших образцах, поскольку соответствующее число может ”не влезть” ни в один тип данных, и на больших алфавитах, так как невозможно будет отобразить буквенную строку в цифровую той же длины. Однако эту проблему можно обойти. Можно вычислять значения p и $N(i+1)$ по модулю q . Значение p по модулю можно вычислять за $O(m)$, а все значения $N(i+1)$ по модулю q за $O(n-m+1)$. Значение q выбирается таким образом, чтобы величина $|A|^m q$ помещалась в компьютерное слово, а рекуррентное соотношение (строка 14) исправляется так, чтобы оно работало по модулю q . Алгоритм Рабина-Карпа очень хорош для поиска в нуклеотидных последовательностях — алфавит содержит 4 буквы, поэтому использование `int` позволяет искать образцы длиной до 16, а применение `long` — до 32. Кроме того операции умножения, деления и `mod` заменяются на соответствующие логические операции и сдвиг, которые исполняются гораздо быстрее, чем умножение, деление и вычисление остатка.

Упражнение 21. Реализуйте на любом доступном языке алгоритм Раббина-Карпа для поиска образца в нуклеотидной последовательности. Примите во внимание, что алфавит – четырехбуквенный, поэтому каждый символ можно кодировать двумя битами. В этом случае операции умножения и деления заменяются на двоичные сдвиги.

Упражнение 22. Реализуйте на любом доступном языке вариацию алгоритма Раббина-Карпа для поиска образца в нуклеотидной последовательности с заданным количеством ошибок (замен).

3.3 Алгоритм Кнута-Морриса-Пратта

Вспомним сложный случай, на котором наивный поиск работает максимальное время:

Строка P aaabaaaaaaaaabcaaaaaaaaaaaaaaaaaaaaaa
 Паттерн T aaaa

Представим, что мы в процессе поиска мы получили наложение:

```

a a a b a a a a a a a a a a b
      a a a a x
  
```

Шли, шли и "потерпели неудачу" на букве x ; нужно сдвигаться и искать дальше. Но ведь понятно, что проверять сходство первых букв, при следующем сдвиге не нужно — они заведомо совпадут, а мы проделаем лишнюю работу. Или другой пример:

```

a a a b a b a b a a a a a a b
      a b a b a x a
  
```

Видно, что делать следующую проверку, сдвинувшись на один символ, бессмысленно — совпадения не будет; для следующей попытки правильно сдвинуть на 2. Получается, что, глядя заранее на *структуру* паттерна, можно сказать на какое количество символов нужно сдвигаться после "неудачи". Если q первых символов совпало, а в позиции $q + 1$ есть несовпадение, то из анализа только образца можно сказать насколько имеет смысл сдвигать образец.

3.3.1 Префикс-функция $sp_i(P)$

Определение 1. Префикс-функцией $sp_i(P)$ называется длина наибольшего собственного суффикса подстроки $P[0, i - 1]$ совпадающего с префиксом P : $P[0, sp_i] = P[i - sp_i, i]$. Собственным называется суффикс, не совпадающий со всем словом.

Пример: возьмем слово `abcabcabba` и найдем все sp_i

$sp_0 = 0$		abcabcabba
$sp_1 = 0$		a bcabcabba
$sp_2 = 0$		ab cabcabba
$sp_3 = 0$		abc abcabba
$sp_4 = 1$	подслово префикс суффикс	abca bcabba abca abca
$sp_5 = 2$	подслово префикс суффикс	abcab cabba abcab abcab
$sp_6 = 3$	подслово префикс суффикс	abcabc abba abcabc abcabc
$sp_7 = 4$	подслово префикс суффикс	abcabca bba abcabca abcabca
$sp_8 = 5$	подслово префикс суффикс	abcabcab ba abcabcab abcabcab
$sp_9 = 0$		abcabcabb a
$sp_{10} = 1$	подслово префикс суффикс	abcabcabba abcabcabba abcabcabba

Отметим некоторые особенности префикс-функции для этого примера. При размере префикса 0 значение sp_0 равно 0, поскольку если нет префикса, то и длина максимального суффикс-префикса равна 0. При размере префикса 1 sp_1 также равен 0, поскольку нас интересуют только собственные суффикс-префиксы. Это выполняется для *любого* образца. Для префиксов размером от 4 до 8 значение префикс-функции увеличивается от 1 до 5. Это связано с тем, что образец содержит точный повтор. Для префикса длиной 9 нет суффикс-префикса, так что sp_9 равен 0. При длине префикса 10 sp_{10} равен 1. Отметим, что для случая sp_7 и sp_8 суффиксы пересекаются с префиксами. Это допустимо.

3.3.2 Алгоритм поиска образца

Лемма 1. Если при наложении образца на текст первое несовпадение обнаружилось в позиции i образца, то мы можем безопасно сдвинуть образец на величину

$$shift = i + 1 - sp_i$$

Доказательство. Приведенное утверждение доказывается от противного. Предположим, что мы передвинули образец на число символов $k < i + 1 - sp_i$, и при этом произошло совпадение (см. Рис3.2): Это значит, в частности,

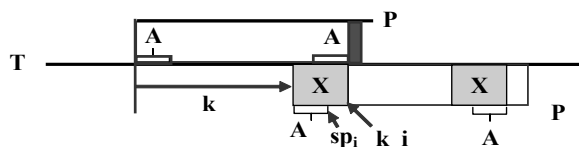


Рис. 3.2: К доказательству безопасности сдвига по Кнуту-Моррису-Пратту. A — суффикс-префикс

что фрагмент $X = P[0, i - k]$ совпадает с текстом. Но при предыдущем наложении у нас совпал без ошибок фрагмент образца $P[0, i]$. Поэтому слово X совпадает с фрагментом $P[k, i]$. Длина этого слова равна $i - k + 1 > i - 1 + sp_i - i + 1 = sp_i$. Это противоречит определению sp_i как *максимального* суффикса-префикса. Противоречие доказывает утверждение. \square

Очевидно, что после сдвига нет необходимости снова проверять символы, принадлежащие соответствующему суффикс-префиксу, поскольку они совпадают согласно определению. Поэтому мы не смотрим на некоторые символы много раз. С другой стороны, можно привести пример, когда на один символ из текста мы будем смотреть много раз (рис.3.3). Видно, что в самом худшем случае, когда надо сдвигаться на 1 позицию мы на символ x посмотрели 5 раз. Однако надо иметь в виду, что в этом случае на остальные символы образца мы посмотрели по одному разу. Поэтому полное число сравнений не превышает длину текста более, чем вдвое. Эти простые соображения можно формализовать и строго доказать, что общее число сравнений в этом алгоритме $N_{KMP} \leq 2n$. Осталось только найти префикс-

a	a	a	a	x	a	a	a	a	x
a	a	a	a	a					
	a	a	a	a	a				
		a	a	a	a	a			
			a	a	a	a	a		
				a	a	a	a	a	
					a	a	a	a	a
1	1	1	1	5	1	1	1	1	5

Рис. 3.3: Анализ сравнений алгоритма КМП. В нижней строке показано, сколько раз мы смотрели на символы из текста

функцию. Этим занимается так называемый препроцессинг.

Упражнение 23. Напишите псевдокод для алгоритма Кнута-Морриса-Пратта поиска образца, при условии, что массив sp дан.

3.3.3 Препроцессинг Кнута-Морриса-Пратта

Поиск $sp_i(P)$ каждый раз, когда мы сдвигаем паттерн относительно строки занимает много лишнего времени. Функция $sp_i(P)$ отражает свойства образца P и не зависит от текста. Поэтому правильнее перед поиском подстроки, равной образцу, сделать препроцессинг — найти все $sp_i(P)$ для этого образца и записать в таблицу. Наивный алгоритм препроцессинга заключается в том, чтобы для каждого i искать заново значение $sp_i(P)$. Наивный алгоритм препроцессинга:

```

1. int[] Naive_Spi_Search(String P){
2.     for(i=1;i<P.length,i++) {
3.         len = 0;
4.         l = i-1;
5.         while(l > 0 && len==0){
6.             s = substring(P,i-1,i);
7.             f = substring(P,0,l);
8.             if(s==f){
9.                 len=1;
10.            }
11.            l--;
12.        }
13.        sp[i]=len;
14.    }
15.    return sp
16. }
```

Для каждого значения i (размера префикса) мы определяем максимальный суффикс-префикс. При этом начинаем искать с самых длинных возможных sp . Строка 4: если бы мы начинали цикл с i , все найденные максимальные префиксы, совпадающие с суффиксами, были бы тривиальными. Строка 5: мы сокращаем размер sp . Если на какой-то итерации величина len стала отличной от нуля, то это значит, что мы нашли sp и искать более короткие суффикс-префиксы не имеет смысла. В строке 8 мы сравниваем суффикс и префикс. Алгоритм работает кубическое время, поскольку цикл строки 2 крутится m раз, цикл строки 5 крутится также порядка m раз (в среднем $m/2$). Внутри этого цикла есть сравнение строк (строка 8), которое также требует времени порядка m . Таким образом, время работы наивного алгоритма имеет порядок $O(m^3)$, то есть время работы кубическое. Столь долго работающий алгоритм мы с негодованием отвергаем. Есть способ найти все sp_i для строки длиной m за линейное время $O(m)$. "Умный" алгоритм препроцессинга не "смотрит" многократно на каждую букву строки.

Лемма 2. *Отметим два свойства суффикс-префикса. Во-первых, $sp_i < i$. Это следует из того, что sp — собственный суффикс-префикс. Во-вторых, если l — длина какого-либо (не обязательно максимального) суффикс-префикса подстроки $P[0, i]$, то sp_l определяет длину другого суффикс-префикса этого*

же под слова. Более того, не существует суффикс-префикса с длиной меньше l и больше sp_l .

Доказательство. Действительно, поскольку есть суффикс-префикс длиной l , то верно по определению $P[0, l] = P[i - l, i]$, но sp_l — означает, что существует суффикс-префикс у под слова $P[0, l]$, т.е. $P[0; sp(l)] = P[l - sp(l); l] = P[i - l; i - l + sp(l)] = P[i - sp(l); i]$ (см.рис 3.4). Вторая часть утверждения о том, что не существует суффикс-префикса длиной меньше l и больше sp_l следует из условия максимальности суффикс-префикса. \square



Рис. 3.4: К доказательству леммы 2. Поскольку есть суффикс-префикс длины l , то по определению есть совпадающие фрагменты (светло-серые). Поскольку sp_l есть длина суффикс-префикса, то темно-серые прямоугольники совпадают. А это значит, что существует суффикс-префикс длины sp_l .

Следствие 1. $sp(sp(\dots sp(sp(i))))$ являются суффиксами-префиксами и других суффиксов-префиксов нет. Здесь и далее для удобства мы используем обозначение $sp(i) = sp_i$

Представим, что мы нашли sp_j для всех префиксов длиной меньше $j < i$. Тогда не представляет труда найти sp_i . Действительно, если $P[sp_i] = P[i]$, то $sp_i = sp_{i-1} + 1$. Если же совпадения нет, то воспользуемся леммой, и посмотрим на более короткие суффиксы-префиксы. Поэтому надо проверить совпадают ли $P[sp(sp(i-1) + 1)]$ с $P[i]$. Если совпадает, то $sp(i) = sp(sp(i-1)) + 1$; Если не совпадает, то надо смотреть предыдущий суффикс-префикс и т.д. (см. рис.3.5). Алгоритм препроцессинга.

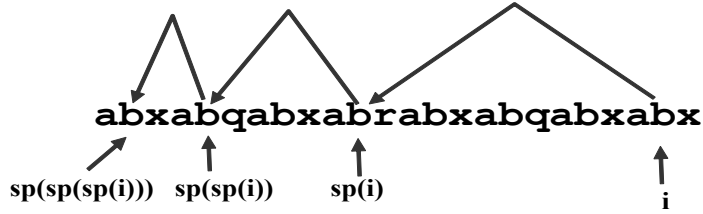


Рис. 3.5: Иллюстрация препроцессинга.

```

1. KMP_Preprocessing(P){
2.     sp[0] = 0;
3.     for (int i=1; i < m-1; i++){
4.         k = sp[i-1];

```

```

5.     while(k>0 && P[i] != P[k]){
7.         k = sp[k-1];
8.     }
9.     if (P[i] == P[k]){
9.         k++;
10.    }
11.    sp[i] = k;
12. }
13. return sp
14. }
```

Можно строго доказать, что время работы алгоритма препроцессинга линейно по длине образца $T = O(|P|)$. Заметим, что итоговое количество итераций цикла `while` определяет асимптотику алгоритма. Теперь отметим, что `k` на каждом шаге увеличивается не более чем на 1. Значит, максимально возможное значение `k` = `n-1`. Поскольку внутри цикла `while` значение `k` лишь уменьшается, получается, что `k` не может суммарно уменьшится больше чем `n-1` раз. В итоге цикл `while` выполнится не более `n` раз, что дает оценку времени алгоритма $O(n)$. Таким образом, время работы алгоритма Кнута-Морриса-Пратта:

$$T_{KPM} = O(|T| + |P|)$$

3.4 Конечные автоматы

Определение 2. *Детерминированным конечным автоматом называют следующее множество:*

$$\{Q, q_0 \in Q, A \subseteq Q, \Sigma, \delta\}$$

где Q — конечное множество, называемое множеством всех *состояний* конечного автомата; q_0 — некое выделенное состояние, называемое *начальным состоянием* конечного автомата; A — подмножество множества состояний (финальное или *допускающее множество состояний*); Σ — конечный алфавит; $\delta : Q \otimes \Sigma \rightarrow Q$ — *функция перехода*, которая отображает множество пар вида "символ-состояние" на множество состояний. Аргументами этой функции служат некий символ алфавита и состояние конечного автомата, значением — новое состояние конечного автомата.

Автомат начинает работать из начального состояния q_0 . Далее автомат считывает по одному символу из входной строки. Полученный на вход символ изменяет состояние автомата на новое из конечного набора возможных состояний $Q = q_0, q_1, \dots, q_m$ в соответствии с имеющейся функцией переходов δ . Нетрудно заметить, что конечный автомат является частным случаем машины Тьюринга, которая во-первых пишет на ленту то, что прочитала

(не изменяет символа в тексте), а во-вторых движется только слева направо. Как видно, здесь ключевым элементом конечного автомата является функция переходов. Множество состояний, алфавит и функцию переходов можно представить в виде таблицы переходов:

состояния	алфавит		
	a	b	c
0*	0	2	1
1	2	3	1
2	3	0	1
3+	0	2	1

* — начальное состояние
+ — допускающее состояние

Часто работу конечного автомата изображают в виде графа переходов (см. рис.3.6). Обратите внимание на то, что в конечном автомате из каждого состояния выходит столько стрелок, сколько букв в алфавите. Пусть на

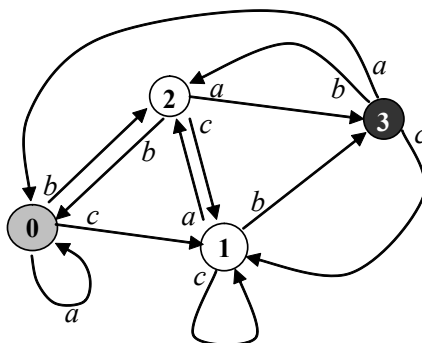


Рис. 3.6: Изображение конечного автомата в виде графа переходов. Начальное состояние 0 (выделено серым). Допускающее состояние 3 (темно-серое).

вход автомата подается строка:

a	a	a	b	c	b	c	c	b	a	a	a	b	a	b	b	c	b	c	a	a	a	b	c	c	a
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Автомат пройдет через следующую последовательность состояний. Начальное состояние 0 (по определению) первые три символа *a* не изменят состояния, поскольку по букве *a* переход происходит из состояния 0 в состояние 0. Затем автомат перейдет в состояние 2 - символ *b* переводит автомат из состояния 0 (текущее состояние) в состояние 2, и т.д. Последовательность состояний следующая:

a	a	a	b	c	b	c	c	b	a	a	a	b	a	b	b	c	b	c	a	a	a	b	c	c	a
0	0	0	2	1	3	1	1	3	0	0	0	2	3	2	0	1	3	1	2	3	0	2	1	1	2

Автомат 5 раз прошел через допускающее состояние. Вообще-то приведенный автомат не понятно зачем нужен. И то, что он несколько раз прошел через допускающее состояние нам ничего не говорит. На самом деле конечные автоматы — достаточно мощный инструмент для распознавания текстов. Как мы увидим позже, конечные автоматы можно применять для поиска сразу нескольких слов и для поиска регулярных выражений, которые составляют основу для разбора текстов. Обозначим Σ^* множество всех слов, которые можно породить из алфавита. Функция $\varphi : \Sigma^* \rightarrow Q$ — *функция конечного состояния*: $\varphi(w)$ — конечное состояние автомата при предъявлении ему слова w (иными словами это состояние, в которое придет автомат после обработки w). Автомат допускает строку w тогда и только тогда, когда $\varphi(w) \in A$. Функцию φ можно определить *рекуррентно*:

- $\varphi(\varepsilon) = q_0$ — начальное состояние ε — пустое слово);
- $\varphi(wa) = \delta(\varphi(w), a)$ для любых $w \in \Sigma^*$, $a \in \Sigma$ (т.е. мы находились в состоянии $\varphi(w)$ и совершили переход в a .)

Ряд алгоритмов поиска подстрок начинаются с построения конечного автомата, который в некотором тексте T находит все вхождения образца P .

3.4.1 Поиск образца с помощью конечного автомата

Конечные автоматы используются для поиска подстроки в строке. Они позволяют находить не только единичное вхождение образца в строку, но и фиксировать все подстроки, равные образцу, сколько бы их ни было, как, впрочем, и алгоритм Кнута-Морриса-Пратта. Пример: Хотим найти в какой-либо строке паттерн вида $abab$. Пусть эта строка состоит только из символов "a" и "b", для построения конечного автомата будем пользоваться алфавитом из тех же символов. Конечный автомат, изображенный

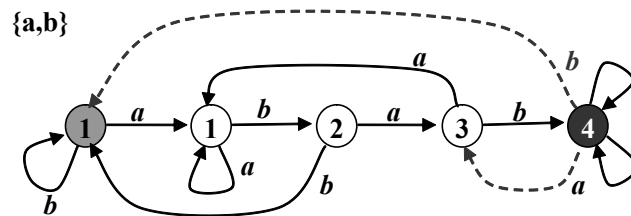


Рис. 3.7: Автомат для поиска образца $abab$. Автомат, включающий только черные стрелки, ищет первое вхождение образца $abab$ в строку. Если в автомате добавить стрелки, исходящие из допускающего состояния (пунктирные), то автомат будет искать все вхождения образца.

на рис.3.7 был придуман напряжением мысли. А существует ли регулярный способ построения конечных автоматов для поиска вхождений образца?

Первым шагом для построения автомата, соответствующего строке-образцу P является построение так называемой суффикс-функции.

Определение 3. Суффикс-функцией называется вспомогательная функция $\sigma : \Sigma^* \rightarrow \{0, 1, \dots, m\}$ — отображение слов на множество чисел. По определению, $\sigma(w)$ — длина максимального (не обязательно собственного) суффикса слова w , который совпадает с префиксом образца P . Σ^* — множество всех слов, построенных на заданном алфавите.

Важно, что для каждого образца своя суффикс-функция. На первый взгляд суффикс-функция $\sigma(w)$ чем-то напоминает суффикс-префикс функцию sp , которую мы рассматривали в алгоритме Кнута-Морриса-Пратта. Однако есть существенная разница. Суффикс-функция определена на множестве всех слов, а суффикс-префикс функция определена на множестве чисел $\{0, 1, \dots, m\}$, поэтому их сравнивать бессмысленно. $\sigma(w)$ является отображением множества Σ^* на $\{0, 1, \dots, m\}$, так, что величина $\sigma(w)$ равна длине максимального префикса P , который является суффиксом строки w :

$$\sigma(w) = \max\{k : w \subset P_k\}$$

Ясно, что если $\sigma(w) = m$, то конец слова w полностью совпадает с образцом. Простой пример суффикс-функции для образца $P = \mathbf{ab}$:

$$\sigma(\varepsilon) = 0$$

$$\sigma(\mathbf{accabca}) = 1$$

$$\sigma(\mathbf{bcabcab}) = 2$$

Очевидные свойства суффикс-функции:

- если длина P равна m , то $\sigma(x) = m$ тогда и только тогда, когда P является суффиксом x
- если x является суффиксом y , то $\sigma(x) \leq \sigma(y)$.

Здесь мы будем использовать обозначения: P_q — префикс длины q образца P , а $P_q\mathbf{a}$ означает, что к слову P_q справа приписана буква \mathbf{a} . Теперь определим конечный автомат, соответствующий образцу P , таким образом:

Множество состояний	$Q = \{0, 1, \dots, m\}$
Начальное состояние	$q_0 = 0$
Единственное допускающее состояние	$A = m$
Функция переходов δ	$\delta(q, \mathbf{a}) = \sigma(P_q\mathbf{a})$

Вот пример построения функции переходов для образца $\mathbf{abbabba}$

P_i	$P_i a$	$\sigma(P_i a)$	$P_i b$	$\sigma(P_i b)$
abbabba	aa	1	ab	0
abbabba	aba	1	abb	3
abbabba	abba	4	abbb	0
abbabba	abbaa	1	abbab	5
abbabba	abbaba	1	abbabb	6
abbabba	abbabba	7	abbabbb	0
abbabba	abbabbaa	1	abbabbab	5

Докажем, что построенный конечный автомат действительно способен находить образец P и ничего больше. Для этого докажем пару лемм и теорему.

Лемма 1. $\sigma(xa) \leq \sigma(x) + 1$

Доказательство. Допустим $\sigma(xa) > \sigma(x) + 1$. Тогда отбросим последний символ от наибольшего суффикса xa , совпадающего с префиксом P . Тогда получим суффикс строки x , который длиннее $\sigma(x)$ и является префиксом P , что противоречит определению (рис.3.8). \square

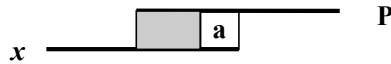


Рис. 3.8: К доказательству леммы 1

Лемма 2. Пусть $q = \sigma(x)$. Тогда $\sigma(xa) = \sigma(P_q a)$.

Доказательство. Действительно, поскольку $\sigma(xa) \leq \sigma(x) + 1$, $\sigma(xa)$ не изменится, если отбросить начало от строки xa , оставив только последние $q + 1$ символов (поскольку $q = \sigma(x)$), то $\sigma(xa) = \sigma(P_q a)$ \square

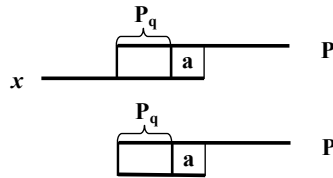


Рис. 3.9: К доказательству леммы 2

Теорема 1. $\varphi(T_i) = \sigma(T_i)$, где T_i - суффикс слова T длиной i : $T_i = T[0 \dots i]$.

Доказательство. Доказательство по индукции. При $i = 0$ -

$$\varphi(\varepsilon) = \sigma(\varepsilon) = 0 \quad \text{если строка пустая}$$

Покажем, что если утверждение теоремы верно при некотором i , то будет верно и при $i + 1$. Пусть \mathbf{a} — это символ строки T в позиции i , то есть следующий за префиксом T_i символ. Тогда

$$\begin{aligned} \varphi(T_{i+1}) &= \delta(\varphi(T_i), \mathbf{a}) && \text{по свойству функции конечного состо-} \\ & && \text{яния } \varphi(w\mathbf{a}) = \delta(\varphi(w), \mathbf{a}) \\ \delta(\varphi(T_i), \mathbf{a}) &= \delta(\sigma(T_i), \mathbf{a}) && \text{по предположению индукции равенство} \\ & && \varphi(T_i) = \sigma(T_i) \text{ верно для } i \\ \delta(\sigma(T_i), \mathbf{a}) &= \delta(q, \mathbf{a}) = \sigma(P_q\mathbf{a}); && q = \sigma(T_i) \text{ — обозначение для } q. \text{ Функ-} \\ & && \text{ция переходов автомата по определе-} \\ & && \text{нию задается так: } \delta(q, \mathbf{a}) = \sigma(P_q\mathbf{a}) \\ \sigma(P_q\mathbf{a}) &= \sigma(T_{i+1}) && \text{согласно лемме 2 о суффикс-функции} \end{aligned}$$

Таким образом доказано, что $\varphi(T_{i+1}) = \sigma(T_{i+1})$ □

Следствие 1. *Следствие. Автомат приходит в допускающее состояние, только если суффикс строки T_i совпадает с образцом P .*

Поиск подстроки с применением конечного автомата работает за время порядка $O(n + T(\sigma))$, где n — длина строки, в которой ищем вхождение, $T(\sigma)$ — время построение суффикс-функции. Это немного больше, чем время работы алгоритма Кнута-Морриса-Пратта. Однако у конечных автоматов есть большое преимущество: они позволяют искать вырожденные или неточные вхождения.

3.4.2 Поиск регулярных выражений. Недетерминированные конечные автоматы.

Все изученные нами способы поиска подстроки (наивный алгоритм, алгоритмы Рабина-Карпа, Кнута-Морриса-Пратта, поиск с помощью конечных автоматов) служат для поиска точного вхождения в строку. Если же мы хотим найти не точный паттерн, а регулярное выражение, описывающее сразу несколько образцов? Регулярные выражения используются для сжатого описания некоторого множества строк с помощью шаблонов, без необходимости перечисления всех элементов этого множества. Пример регулярного выражения: $\mathbf{R} = (\mathbf{abb}) * (\mathbf{c}) + (\mathbf{abc})$.

*	Символ или группа символов может появиться несколько раз или ни разу
+	Символ или группа символов должна появиться хоть один раз.
выражение, после которого не стоит никаких знаков	Символ или группа символов должна появиться ровно один раз.

Нашему регулярному выражению соответствуют строки $ssssabc$ и $abbabc$ и много других строк. Как найти любую из них пользуясь регулярным выражением? *Недетерминированный конечный автомат* отличается от детерминированного тем, что результатом функции перехода является некоторое множество состояний (возможно пустое). Поэтому недетерминированный конечный автомат в каждый момент времени находится не в одном состоянии a в множестве состояний. Кроме того, переходы не обязательно определены для всех символов из входного потока. Вместо этого мы введем дополнительный переход ε , который отвечает "считыванию пустого символа", т.е. не снимает символ из строки. Распознавание образца состоит из

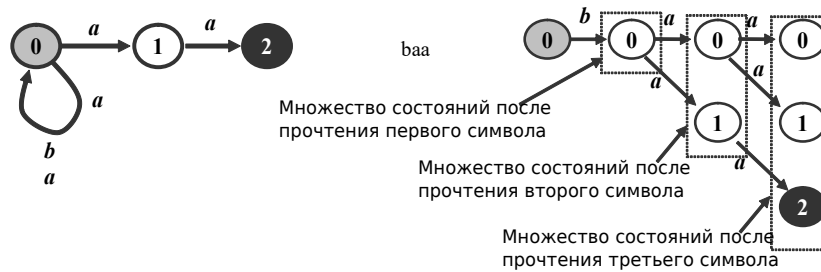


Рис. 3.10: Недетерминированный конечный автомат и множества состояний в процессе обработки строки baa .

двух этапов — построение автомата и по образцу и работа автомата. При построении автомата используется рекурсивное определение регулярного выражения R :

$$R ::= \langle \text{символ} \rangle \mid RR \mid (R) \mid R+ \mid R^*$$

Эта формула означает, что регулярное выражение — это либо просто символ, либо регулярное выражение, за которым следует другое регулярное выражение, либо другое регулярное выражение заключенное в скобки, либо регулярное выражение, после которого идет знак $+$, либо регулярное выражение, за которым следует знак $*$. Такая рекурсия позволяет разложить регулярное выражения на элементарные составляющие. Для каждого составляющего можно построить конечный автомат, а затем эти автоматы можно объединять. При этом если один из элементов заканчивается символом $*$, то от начала (начального состояния) соответствующего автомата в его конец (во все допускающие состояния) идет ε -стрелка, что означает, что соответствующий фрагмент можно пропустить. Если после элементарного выражения стоит $+$ или $*$, то из всех допускающих состояний идет ε -стрелка в начальное состояние, что означает, что группу можно повторить (рис. 3.11). Теперь опишем, как этот автомат работает. При прочтении слова автомат может перейти в некоторое множество состояний. Введем некоторые обозначения: N_0 — множество, куда можно попасть из начального состояния по стрелкам с буквой ε ; N_{i+1} — множество состояний, в которое автомат

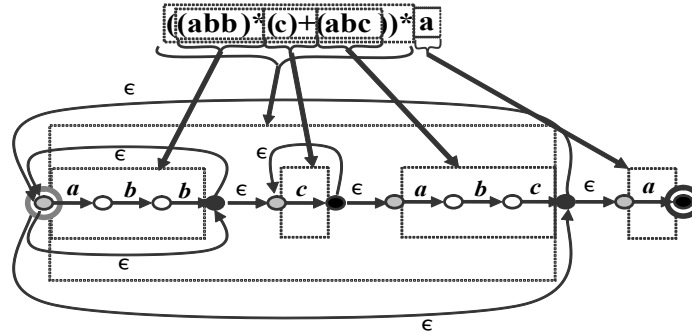


Рис. 3.11: Пример построения конечного автомата для регулярного выражения. Детали переходов в самых внутренних автоматах ОПУЩЕНЫ! Начальное и допускающие состояния общего автомата выделены двойными кругами. Начальные состояния — серые, допускающие — черные.

может попасть из предыдущего множества состояний, прочитав $i + 1$ -ю букву строки T . Если на каком-то шаге автомат нашел соответствие (match), то множество состояний автомата имеет непустое пересечение с множеством допускающих состояний. Если очередное множество состояний пусто, то автомат отвергает строку. Разбор строки заключается в перестраивании множества N_i . Ясно, что размер множества N_i не превышает длину образца, поэтому время разбора строки оценивается как $T = O(n \times m)$. Рассмотрим пример работы автомата для регулярного выражения $((ab)^*aa)+b$ на тексте $aaabaab$ (рис.3.12). Начальное состояние 0 имеет ϵ -переход в состояние 2. Поэтому вначале автомат находится сразу в двух состояниях — 0,2. При считывании буквы a состояние 0 переходит в состояние 1, а состояние 2 — в состояние 3. Считывание второй буквы a убивает состояние 1, а состояние 3 переходит в состояние 4. Но состояние 4 имеет ϵ -переход в состояние 0, поэтому множество состояний становится 4,0. Далее идет третья буква a , которая убивает состояние 4, а состояние 0 переводит в состояние 1. Буква b переводит состояние 1 в состояние 2, а ϵ -переход из состояния 2 порождает состояние 0. Теперь множество состояний автомата 2,0. Следующая буква переводит автомат в множество состояний 3,1. Очередная буква a переводит состояние 3 в состояние 4 и в состояние 0 (по ϵ -переходу), а состояние 1 погибает. Наконец, последняя буква убивает состояние 0, а состояние 4 переводит к допускающему состоянию 5. Минус описанного автомата — он ищет регулярное выражение только в префиксе строки. Чтобы искать регулярное выражение в любом месте строки нужно модифицировать регулярное выражение R : $R' = [\xi]R$, где ξ обозначает все буквы алфавита Σ и переход в первое состояние(самого себя). Конечный автомат примет такой "кудрявый" префикс (рис. 3.13).

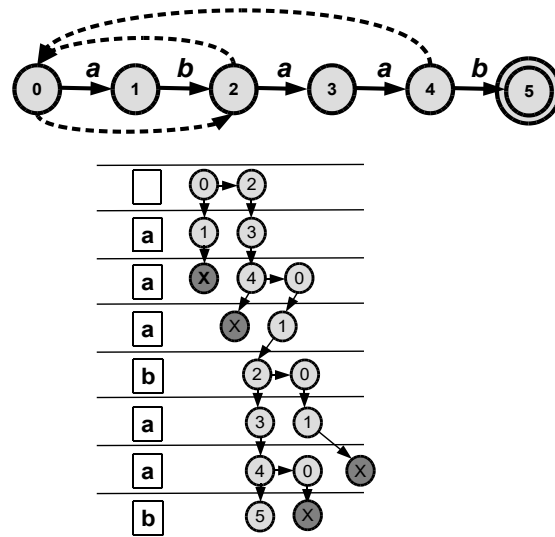


Рис. 3.12: Вверху показан автомат для регулярного выражения $((ab)^*aa)+b$. Пунктирными линиями обозначены ϵ -переходы. Внизу – последовательность множеств состояний при разборе строки $aaabaab$

3.5 Поиск многих образцов. Алгоритм Ахо-Корасик

А если перед нами стоит такая задача (ее еще называют задачей точного поиска или задачей поиска по групповому запросу): задано множество образцов $P = \{P_1, P_2, \dots, P_z\}$. Требуется обнаружить все вхождения в текст T хотя бы одного образца из P . Здесь неприменимы недетерминированные конечные автоматы. Можно, конечно, предложить проверить каждый из образцов по очереди, но их может быть очень много. Алгоритм классического решения задачи точного сопоставления множеств был предложен в 1975 г. Альфредом Ахо (Alfred V. Aho) и Маргарет Корасик (Margaret J. Corasick). Маргарет Корасик — женщина, ее фамилия не склоняется!), поэтому неверно название "алгоритм Ахо-Корасика". Это эффективный алгоритм, работающий за время $T = O(n + \sum_i m_i)$, где m_i — длина образца P_i . Языком L называется некий набор слов над алфавитом Σ . Пусть $\Sigma = \{a, b, \dots, ю, я\}$. Возьмем такой язык над этим алфавитом:

{абрау, абрек, брусника, брусья, руслан, русь}

Для языка L можно построить язык всех слов, содержащих хотя бы одно слово исходного языка в качестве подслово. Будем обозначать такой язык L' . Все слова можно расположить друг под другом в лексикографическом порядке. При это станет очевидным, что некоторые из них содержат общий префикс (отмечен серым):

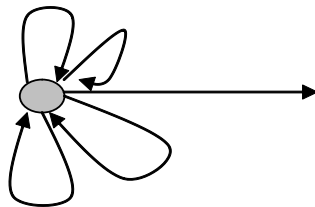


Рис. 3.13: Префикс текста для поиска образца с помощью недетерминированного автомата

абрау
абрек
брусника
брусья
руслан
русь

Список слов можно представить в виде дерева, где общие префиксы не повторяются.

абр –ау
–ек
брус–ника
–ья
рус –лан
–ь

Представим автомат для поиска слов нашего языка в виде дерева рис.3.14. Каждый узел этого дерева отвечает состоянию автомата. На ребрах записаны соответствующие символы (*метки ребер*). Листья этого дерева – допускающие состояния. Корень – начальное состояние. Показаны не все стрелки. Любая стрелка из любого состояний с "не той" буквой будет вести в начальное состояние. Ясно, что такой автомат не всегда будет работать правильно. Например, в тексте "абрусья" не будет найдено слово брусья. Поэтому следует поступить по аналогии с конечным автоматом для распознавания одного слова, т.е. нам необходимо провести обратные стрелки в какие-то разумные вершины. Это такие вершины, для которых префикс (путь от корня до вершины) совпадает с суффиксом той части слова, которая была принята до момента неудачи. Для того, чтобы нарисовать стрелки обратных переходов. Для этого введем функцию неудач. Функция неудач — куда надо идти, если идти вперед некуда, т.е. для каждой вершины определено: $f(v) : V \rightarrow V$. Метку ребра, ведущего в вершину v будем обозначать $g(v)$. Для каждой вершины (состояния) можно определить ее уровень — расстояние до корня. Ясно, что функция неудач должна идти в вершину меньшего уровня (догадайтесь, почему?), но далеко не всегда должна идти в корень: $level(f(v)) < level(v)$ если $level(v) > 0$. Алгоритм построения

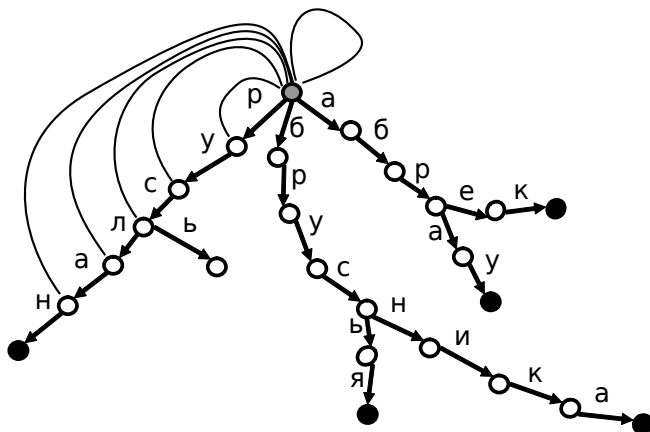


Рис. 3.14: Автомат для распознавания группы слов (языка). Показаны не все стрелки. Любая стрелка из любого состояний с "не той" буквой будет вести в начальное состояние

обратных переходов (функции неудач) основан на просмотре вершин по уровням. Ясно, что если, находясь в корне, мы не примем буквы (это значит, что ни одно из слов языка на эту букву не начинается) то мы должны в корне остаться. Отсюда первое правило: $f(0) = 0$. Если, приняв первую букву какого-то слова из языка у нас нет продолжения, то ясно, что надо вернуться в корень (собственный суффикс слова длины 1 равен 0). Отсюда второе правило: $\forall v : level(v) = 1 \ f(v) = 0$. Если на более высоком уровне нас постигла неудача, то надо посмотреть, куда мы шли в случае неудачи из предыдущей вершины. Если мы шли в корень, то из текущей вершины мы можем принять букву в качестве первой буквы слова из языка (если такое слово есть). Поэтому следующее правило: $if(f(v_{prev}) = 0)$ ищем такую вершину u , что $u : level(u) = 1 \ \& \ g(u) = g(v)$. Если находим, то определяем $f(v) = u$, иначе у нас нет слова, начинающегося с буквы $g(v)$, потому $f(v) = 0$. Если же предыдущая вершина указывала в случае неудачи на какую-то другую вершину $w^0 = f(v_{prev})$, то надо проверить, а не совпадают ли метки ребра в текущую вершину $g(v)$ и метка какого-либо ребра, исходящего из w (не совпадает ли суффикс подслова с префиксом другого слова). Иными словами $if(f(v_{prev}) \neq 0)$ ищем вершину u такую, что $w^0 = f(v_{prev})$; и $\exists u : (u_{prev} = w^0) \ \& \ (g(w^0) = g(v))$. Если такая вершина нашлась, то $f(v) = u$. Если такая вершина не нашлась, то смотрим функцию неудач $w^1 = f(w^0) = f(f(v))$ и проверяем метки ребер, исходящих из w^1 : $\exists u : (u_{prev} = w^1) \ \& \ (g(w^1) = g(v)) \Rightarrow f(v) = u$. Так продолжаем до тех пор, пока не найдем соответствующий префикс, совпадающий с суффиксом, либо не придем в корень. Построение функции неудач аналогично препроцессингу Кнута-Морриса-Пратта.

Итак алгоритм:

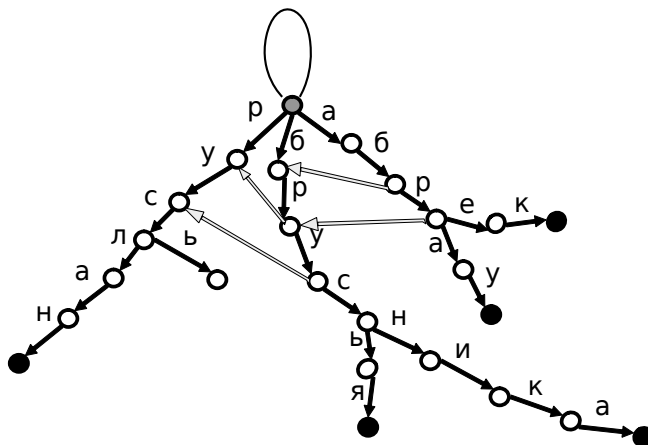


Рис. 3.15: Автомат Ахо-Корасик. Обратные ребра (неудачи) обозначены серым. Ребра, ведущие в корень не показаны.

```

1. f(0)=0;
2. for(level from 1){
3.   for(v: level(v) = level){
4.     f(v)=root;
5.     if(level > 0){
6.       w = f(v.prev);
7.       while(w != root) {
8.         for(u: u.prev = w){
9.           if(g(u) = g(v)) f(v)=u;
10.        }
11.        if(f(v) = root) w = f(w)
12.        else w = root
13.      }
14.    }
15.    if(f(v)=root){
16.      for(u: level(u) = 1){
17.        if(g(v)=g(u)): f(v)=u;
18.      }
19.    }
20.  }
21. }
22.}

```

Пояснение. Строка 1 определяет функцию неудач для корня. Строка 2 определяет цикл по уровням. Строка 3: определяет цикл по вершинам уровня. Строка 4: назначаем функцию неудач равной корню (root), впрочем по-

том можем ее изменить (строка 10). Далее пытаемся найти такую вершину, в которую мы можем пойти по функциям неудач из предыдущей вершины (строка 6), и из которой есть ребро с меткой $g(v)$ (строки 9–11). Цикл поиска (строка 6) крутится пока не определится функция неудач (условие $f(v) = root$), либо, пока не дойдем до корня (условие $w! = root$). Если до сих пор функция неудач не определилась (строка 15), то пытаемся начать новое слово (строки 16–18).

Важное замечание. функция неудач показывает стрелки, на которых не написан символ (метка). Это значит, что при таком переходе символ на самом деле не снимается и после перехода в новое состояние мы должны еще раз сделать проверку этого символа. Исключения представляют только переходы неудач, исходящие из корня. При таких переходах символ снимается из входного потока и более не проверяется. Поэтому в автомате Ахо-Корасик есть три вида переходов (стрелок) – стрелки с метками, стрелки без метки и без снятия символа и переходы без метки, но со снятием символа. На самом деле последний тип перехода есть только из корня, и на нем можно поставить метку – любой символ, кроме тех, которые стоят на 'правильных' стрелках.

Замечание. Если множество образцов состоит из одного слова, то автомат Ахо-Корасик идентичен алгоритму Кнута-Морриса-Пратта. В частности, построение функции неудач устроено так же, как препроцессинг Кнута-Морриса-Пратта. Автомат, построенный по алгоритму Ахо-Корасик (рис.3.15), имеет более "приятный" вид, чем автомат (3.14). В нем нет "избыточных" состояний. Такой автомат называют еще деревом Ахо-Корасик. Проверим с помощью нашего автомата, принадлежит ли слово "абруслабрек" к языку L' . Проиллюстрируем шаги работы автомата с помощью таблицы. Состояниям будем индексировать префиксами слов.

Автомат в состоянии	На входе строка	Автомат переходит в состояние
""	а	"а"
"а"	б	"аб"
"аб"	р	"абр"
"абр"	у	"бру"
"бру"	с	"брус"
"брус"	л	"русл"
"русл"	а	"русла"
"русла"	б	"аб"
"аб"	р	"абр"
"абр"	е	"абре"
"абре"	к	"абрек"

Оценка времени работы автомата Ахо-Корасик . Автомат Ахо-Корасик является частным случаем недерминированного автомата. При его работе,

если нет совпадения по прямому пути, мы начинаем довольно много прыгать. Однако, основываясь на тех же соображениях, что при оценке работы алгоритма Кнутта-Морриса-Пратта, можно показать, что на каждый символ текста мы посмотрим $N_{AK} \leq (k+1)$, где k — число образцов. Можно это время существенно сократить, если мы преобразуем недетерминированный автомат в детерминированный, когда на каждой стрелке есть определенный символ. В этом случае на каждый символ текста мы посмотрим ровно один раз.

Упражнение 24. *Преобразуйте автомат Ахо-Корасик в детерминированный конечный автомат*

Автомат Ахо-Корасик имеет широкий спектр применений. Например, поиск компьютерных вирусов часто использует этот алгоритм. Вся база данных вирусов собирается в автомат Ахо-Корасик, после чего этот автомат применяется для сканирования компьютера. В биоинформатике также часто применяется этот автомат для поиска множества возможных образцов. Идеи, лежащие в основе автомата Ахо-Корасик можно использовать также для статистического анализа частот встречаемости мотивов.

3.6 Поиск многих образцов в тексте. Препроцессинг текста

В предыдущих параграфах мы рассмотрели задачу поиска подстроки в длинном тексте. При этом подстрока (в общем случае образцы) заранее готовились — строились массивы суффикс-префиксов, создавались разного рода конечные автоматы. При этом подготовленный образец можно многократно применять для поиска в разных текстах. Однако иногда полезно провести подготовку текста, чтобы потом можно было эту подготовленную структуру использовать для поиска разных образцов. Например, хорошо бы подготовить последовательность генома, чтобы потом в нем эффективно искать разного рода образцы. Рассмотрим все ту же основную задачу поиска точного вхождения строки-образца P в длинном тексте T .

3.6.1 Суффиксные массивы

Образец присутствует в тексте, если и только если с него начинается один из суффиксов текста. Наивный алгоритм поиска на самом деле заключался в последовательном сравнении образца со всеми суффиксами — прикладывание образца к тексту в позиции i эквивалентно сравнению образца и начала суффикса. Рассмотрим, например текст `mississippi`. Чтобы в нем найти слово, скажем `sip`, мы прикладываем образец к началам суффиксов:

0	mississippi
1	ississippi
2	ssissippi
3	sissippi
4	issippi
5	ssippi
6	sippi
7	ippi
8	ppi
9	pi
10	i

Поиск образца в этой интерпретации напоминает поиск в массиве (не сортированном), что требует времени порядка $T = O(n \times t)$, где t — среднее время одного сравнения. Это время не превышает длину образца. Но в начале книги мы узнали, что искать в сортированном массиве гораздо эффективнее объект можно найти за логарифмическое время. Отсортируем массив суффиксов в лексикографическом порядке:

10	i
7	ippi
4	issippi
1	ississippi
0	mississippi
9	pi
8	ppi
6	sippi
3	sissippi
5	ssippi
2	ssissippi

В таком массиве можно найти образец за время порядка $T = O(\log n \times t) \leq O(\log n \times m)$, что намного быстрее, чем с помощью алгоритма Кнута-Морриса-Пратта. Здесь надо заметить, что совсем нет необходимости реально помнить суффиксы в отдельном массиве — ведь у нас есть исходный текст. Поэтому в массиве достаточно иметь позиции суффиксов. Хотя поиск происходит быстро надо потратить время на построение суффиксного массива. Время сортировки массива составляет величину порядка $n \times \log n \times tc$, где tc — среднее время сравнения суффиксов (чтобы понять кто больше). Таким образом, чтобы провести поиск одного образца в тексте требуется время порядка $T = O(n \times \log_2 n \times t \times tc)$, что больше, чем время, необходимое для поиска образца алгоритмом Кнута-Морриса-Пратта. Однако при массовом поиске этот метод имеет значительное преимущество.

Замечание. Время построения суффиксного массива. Очевидно, что строить его прямым способом, сортируя суффиксы, дело долгое. Существует ряд

достаточно понятных алгоритмов, работающих за $O(\log n \times n)$. Есть те, которые работают за $O(n)$ однако они весьма запутаны. За $O(n)$ можно также построить суффиксный массив по суффиксному дереву при обходе в глубину (см. далее).

3.6.2 Суффиксное дерево

При поиске в суффиксном массиве мы тем не менее несколько раз прикладывали образец к суффиксам. Однако, при поиске слова, скажем, "sip" разумно было бы начинать сравнение прямо с той части массива, который начинается с буквы 's', а при поиске слова "miss" разумно было бы начинать поиск в том месте массива, где лежат слова, начинающиеся с буквы 'm'. В алгоритме Ахо-Корасик все слова были организованы в виде дерева. Суффиксы также можно организовать в виде дерева. Однако, чтобы в дальнейшем не было путаницы введем еще символ конца строки, например символ \$. В результате получим такое дерево (рис. 3.16). Основные свойства

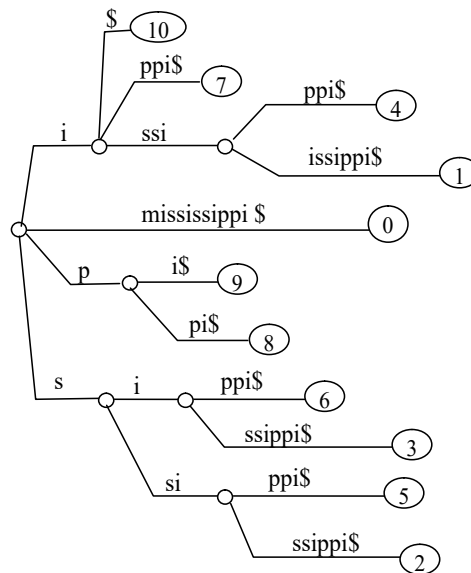


Рис. 3.16: Суффиксное дерево для слова mississippi

суффиксного дерева:

- Листьев столько же, сколько и позиций в слове
- На ребрах определены строки (метка ребра)

- Каждый путь от корня до листа — суффикс
- На листьях указаны позиции начал соответствующих суффиксов
- Можно построить такое дерево за время $T = O(n)$, где n — длина слова. В это надо поверить, поскольку, хотя соответствующий алгоритм существует, он достаточно сложен и громоздак, чтобы его здесь описывать. Желающие могут обратиться к книге Гасфилда.

Разберемся теперь, как хранить суффиксное дерево, используя линейную память. Для этого оставим в суффиксном дереве только вершины, имеющие более одного ребенка. Вместо строки $[t_i \cdots t_j]$ для ребра будем хранить ссылку на соответствующий сегмент слова. В описанном виде суффиксное дерево называется сжатым. Заметим, что, так как теперь каждая внутренняя вершина является вершиной разветвления, она добавляет к дереву как минимум один лист. Листьев в любом суффиксном дереве столько же, сколько букв в слове (пусть $n > 1$), поэтому число внутренних вершин в сжатом суффиксном дереве варьирует от 1 до $n-1$. Таким образом, всего вершин и ребер в сжатом суффиксном дереве будет линейное зависеть от длины и дерево будет занимать линейную память.

3.6.3 Алгоритмы построения суффиксного дерева

Наивный алгоритм построения суффиксного дерева предполагает, что одновременно есть доступ ко всему тексту. Однако для многих приложений нам не хочется анализировать сразу весь текст, а строить дерево по мере чтения текста. Такие алгоритмы называются *on-line* алгоритмами. On-line подход построения суффиксного дерева строит дерево последовательно для всех префиксов слова, получая в итоге полное дерево. Прежде чем обратиться к описанию алгоритма, введем понятие *неявного* суффиксного дерева. Неявное суффиксное дерево — это суффиксное дерево для текста без знака \$ на конце. Некоторые суффиксы в нем заканчиваются на ребрах или во внутренних вершинах, и их номер нигде не хранится. В таком дереве листьев, как правило, меньше, чем позиций в слове. Обычное суффиксное дерево называется еще явным. Есть слово $T = t_0 t_1 \cdots t_n$. Построим его суффиксное дерево с помощью *on-line* алгоритма. Будем строить деревья по очереди для всех префиксов слова. На первом шаге строим дерево для слова, состоящего из одной буквы t_0 . На следующем шаге получаем букву t_1 , теперь нам нужно строить дерево для слова $t_0 t_1$. Не будем строить заново, а достроим уже имеющееся дерево. Вот план работы алгоритма над словом $t_0 t_1 \cdots t_n$

0 Строим суффиксное дерево для t_0

1 Расширяем его до дерева для t_1

...

(рис.3.17): $n-1$ Расширяем дерево для $t_0 t_1 \cdots t_{n-2}$ до дерева для $t_0 t_1 \cdots t_{n-1}$.

n Расширяем дерево для $t_0 t_1 \cdots t_{n-1}$ до дерева для $t_0 t_1 \cdots t_n$ (делаем из неявного дерева явное)

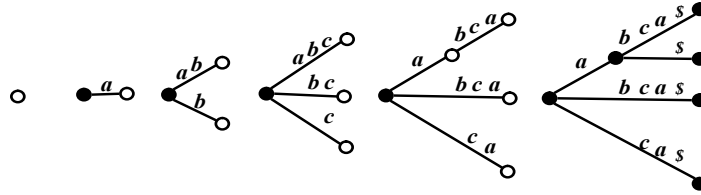


Рис. 3.17: Построение суффиксного дерева слова "abca" с помощью on-line алгоритма. Первые пять деревьев являются неявными, последнее — явное суффиксное дерево

смотрим i -й этап работы алгоритма. На этом этапе мы перестраиваем неявное дерево для префикса строки $T_i = t_0 t_1 \dots t_{i-1}$ в неявное дерево для префикса $T_i = t_0 t_1 \dots t_{i-1}$. С этой целью для каждого j от 0 до $i-1$ находим в суффиксном дереве конец суффикса $t_j \dots t_{i-1}$. Далее продлеваем его буквой t_i , если необходимо. При этом действуем по одному из следующих трех правил:

1. Правило продления. Если, прочитав суффикс $t_j \dots t_{i-1}$, мы пришли в лист суффиксного дерева, удлиним ребро, ведущее в этот лист, добавляя к строке, записанной на ребре, новую букву t_i .
2. Правило ответвления. Прочитав суффикс $t_j \dots t_{i-1}$, мы можем остановиться не на листе, а в какой-нибудь внутренней вершине или даже прямо на ребре. Тогда определяем новую "точку роста" и определяем на ней начало соответствующего суффикса. Если остановились на ребре, а следующая буква — не t_i , придется разбить ребро на две части. Новая вершина добавляется после буквы t_{i-1} . От этой вершины будут отходить два ребра: остаток старого и новое ребро, несущее букву t_i .
3. Пустое правило. Если, прочитав суффикс $t_j \dots t_{i-1}$, мы видим, что дальше уже есть нужная нам буква t_i , не создаем ничего нового, кроме точки роста.

Оценка времени работы алгоритма. Для построения суффиксного дерева слова длиной n on-line алгоритм должен пройти n шагов. На i -ом шаге мы продлеваем i суффиксов (к примеру, на втором шаге мы продлеваем два суффикса: $t_0 t_1$ и t_1). Продление каждого j -го суффикса на i -ом шаге занимает $O(i-j)$ операций. Действительно, счетчик j на каждом шаге меняется от нуля до i . Рассмотрим опять второй шаг работы алгоритма: нам нужно продлить слово длиной 2 ($i = 2$) буквой t_1 . Счетчик j сначала равен нулю, и мы продлеваем суффикс $t_0 t_1$ (длиной $i-j$), это также займет $O(i-j)$ операций, поскольку читаем суффикс побуквенно. Затем $j = 1$, продлеваем слово t_1 . Когда $i = j = 2$, продлеваем пустое слово. В итоге продление каждого суффикса на i -ом шаге стоит $O(i-j)$ операций, где j для каждого i

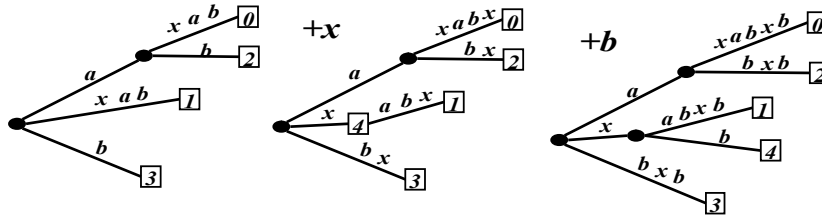


Рис. 3.18: К алгоритму построения суффиксного дерева. При приходе нового символа мы прежде всего продолжаем "точки роста". Потом проверяем, не привело ли появление нового символа к появлению новых точек роста. Точки роста отмечены квадратиками

меняется от 0 до i . С учетом сказанного составим выражение для скорости работы алгоритма над словом длины n :

$$T(n) = O\left(\sum_{i=0}^n \sum_{j=0}^i (i-j)\right) = O(n^2)$$

Существуют алгоритмы, формирующие суффиксное дерево слова длиной n за линейное время $O(n)$, но они значительно сложнее кубического алгоритма. Один из них строится модификацией описанного on-line подхода, другой (алгоритм Мак-Крейта) начинает работу с пустого дерева и добавляет суффиксы, начиная с самого длинного. Алгоритм Мак-Крейта не является on-line алгоритмом, т.е. для его работы необходима вся строка целиком. Мы не будем обсуждать здесь эти алгоритмы.

3.6.4 Сравнение суффиксных деревьев и суффиксных массивов

Суффиксные массивы считаются более простым способом хранения текста, чем суффиксные деревья. Действительно, алгоритм построения суффиксного массива гораздо проще, чем любой алгоритм построения суффиксного дерева. Огромный плюс суффиксных массивов — их размер в памяти определяется только размерами текста T и никак не зависит от его алфавита, в отличие от суффиксных деревьев. Несмотря на то, что время построения массива хуже, чем дерева, суффиксное дерево требует большого количества памяти — $T = O(|T| \cdot |\Sigma|)$, где $|\Sigma|$ — мощность алфавита. В некоторых задачах о поиске подстрок алфавит очень велик (например, естественные языки с большими алфавитами, сравнение изображений, где алфавит — все возможные значения цвета пикселя) и иногда память, которую занимает суффиксное дерево, делает его неприемлемым. Суффиксные массивы гораздо более рационально используют память: для их хранения достаточно одной строки длины $|T|$ и $4 \cdot |T|$ байтов для хранения целых чисел суффикс-

ного массива. При этом поиск по суффиксному массиву работает почти так же быстро как поиск по суффиксному дереву.

3.6.5 Примеры применения суффиксных деревьев

Поиск подстроки. Знакомая нам уже задача: поиск подстрок, равных образцу P длиной m , в тексте T длиной n . Некоторые уточнения: текст фиксирован, ищем много разных образцов, текст длиннее образца, в образце не встречается знак конца строки '\$'. В таком виде задача поиска выгоднее всего решается с помощью суффиксных деревьев, при условии достаточного объема памяти (иначе лучше задействовать суффиксные массивы, см. ниже). Примеры задачи: поиск разных французских слов в тексте "Войны и мира" Л.Н. Толстого; биоинформатический пример — поиск регуляторных последовательностей в эукариотическом геноме. Для решения задачи построим суффиксное дерево для текста T . Будем читать паттерн вдоль дерева от корня. Если в какой-то момент не сможем прочитать следующую букву паттерна, значит, в тексте T ни разу не встречается строка, равная P . Допустим, что такая строка в тексте все-таки есть, тогда, прочитав эту строку, мы приходим либо во внутреннюю вершину v , либо останавливаемся на ребре (прийти в лист мы не можем, поскольку в паттерне не встречается знака '\$'). Если остановились на ребре, проходим до ближайшей вершины v вниз по дереву. Далее читаем числа на листьях потомков вершины v . Эти числа — номера суффиксов, начинающихся с подстроки P , а значит, индексы вхождений P в текст T . Такой поиск занимает приблизительно $T = O(m)$ операций. Точно оценивая, заметим, что, прочитав все буквы паттерна по дереву за $O(m)$ операций, нужно еще дойти до ближайшей вершины, а потом узнать у ее потомков номера соответствующих суффиксов. Поэтому более точная оценка времени работы поиска говорит о $O(m + l)$, где l — число листьев у вершины v . С помощью описанного поиска мы находим позиции всех вхождений паттерна P в текст T . Быстрее работает алгоритм, который ищет только первое вхождение слова в текст. Если каждая вершина суффиксного дерева "знает" число, хранящееся в ее ближайшем ребенке, дойдя до вершины v по прочтении слова P , мы узнаем индекс первого вхождения P в текст T . Этого часто достаточно для решения задачи (особенно если мы хотим просто знать, есть или нет в некоем тексте некое слово). Алгоритм поиска первого вхождения работает в точности за $O(m)$.

Поиск повторов. Любая внутренняя вершина суффиксного дерева соответствует повтору. Позиции на всех листьях-потомках этой вершины есть точные повторы.

Поиск наибольшего общего подслоа нескольких слов. Кроме задачи о поиске образца в фиксированной строке, суффиксные деревья помогают решать задачу о поиске максимального общего подслоа. Пусть у нас есть два слова, к примеру, нуклеотидные последовательности $T_1 = \text{"atgcat"}$

является решением задачи, в нашем случае это строка "atgc". Алгоритм работает за $O(|T_1| + T_2|)$ — за линейное время. Согласитесь, это лучше, чем тетричное время. Понятно, что общее подслово можно искать не только в двух, но и в большем количестве слов. Для трех слов придется вводить, кроме длинных и коротких, еще и "средние" суффиксы, начинающиеся во втором слове.

3.7 Сложность текста

Вы всегда легко запоминаете телефонные номера своих друзей? Поспорим, что есть такие номера, что, увидев, запомните сразу, а есть такие, что всегда приходится записывать. Сложность здесь заключается не обязательно в вашей памяти. Телефонные номера — это тексты. Все тексты разные: есть "Война и мир", а есть геном *E.coli*; телефонные номера — тоже разные. Как оценить сложность текста? Вот, к примеру, два примера нуклеотидных текстов:

```
aaaaaaaaaaattttttttt
atatatatatatat
```

Оба текста простые, потому что их просто описать. Вместо первого текста можно сказать: "12 раз **a**, 8 раз **t**". Второй текст еще проще — "7 раз **at**". А попробуйте описать такой текст:

```
caactgaaactgttgactta
```

По длине он такой же как первый, но гораздо сложнее. Такой текст проще прочитать (назвать все буквы), чем как-то описывать. В 1960-х годах русский математик А.Н.Колмогоров поставил вопрос: "Какова внутренняя сложность описания строки двоичных символов?". Из ответа на этот вопрос сегодня мы с вами знаем, что сложность текста это — длина самой короткой программы, которая генерирует данный текст. Такая длина называется сложностью текста по Колмогорову и обозначается $K(s)$, где s — текст (строка). Введем некоторые формальные понятия. Во-первых, будем для простоты рассматривать двоичный алфавит 0,1. Множество всех слов над этим алфавитом обозначим $\Sigma = \{0, 1\}^*$. Рассмотрим множество $E = \{x, y\}$ пар слов такое, что:

1. если $(x, y_1) \in E$ и $(x, y_2) \in E$, то $y_1 = y_2$
2. Существует программа Π такая, что $\forall (x, y) \in E : \Pi(x) = y$

Такое множество и программа называются системой описания текста. На самом деле достаточно одной программы и ее области определения. Во-первых, это похоже на архивирование-разархивирование. Здесь x — архив, y — разархивированный файл, а программа Π — разархиватор. Очевидно, что далеко не любой файл является архивом чего-либо (здесь и появляется ограничение на область определения программы). Во-вторых, здесь нет

взаимной однозначности — разным x может соответствовать один y . А раз так, тот можно поставить вопрос о наиболее коротком описании строки.

Определение 4. Сложностью $K_{\Pi}(y)$ текста y относительно программы Π называется минимум:

$$K_{\Pi}(y) = \min_x (\Pi(x) = y)$$

Системы описания текста можно сравнивать. Говорят, что система Π_1 не хуже системы Π_2 , если существует константа C , не зависящая от слова x такая, что

$$\forall x : K_{\Pi_1}(x) \leq K_{\Pi_2}(x) + C$$

Если есть две системы описания Π_1 и Π_2 , то может быть на одних словах первая система более эффективна, а на других — вторая порождает более короткие описания. Есть теорема, что если есть две системы описания, то можно построить третью, которая будет не хуже двух исходных. Здесь нельзя забывать о константе! За каждый переход на новую систему описания надо платить добавлением константы в размер описателя. Тем не менее, можно поставить вопрос об описателе (программе), который был бы не хуже всех остальных и о размере описания текста с точки зрения этого описателя. Любая программа имеет конечный размер. Поэтому количество программ счетно, а значит их можно пронумеровать. Если задан текст, то можно попробовать описать оптимальным способом всеми программами, и в начало описания поставить номер соответствующей программы. Таким образом можно построить универсальный описатель.

Определение 5. Сложностью по Колмогорову называется длина самого короткого описания текста относительно некоторого универсального описателя.

Свойства сложности по Колмогорову.

Свойство 1. Сложность по Колмогорову не вычислима, т. е. не существует алгоритма, который глядя на последовательность за конечное время вычислит его сложность.

Свойство 2. Количество слов W сложности K может быть оценена как:

$$2^{K-C} \leq W \leq 2^{K+1}$$

Случайный текст по Колмогорову — текст, любое описание которого не короче самого текста. Конечно, эта характеристика относится к любым текстам, а не только к двоичным. Удобно относить ее к нуклеотидным текстам. Заметим, что большинство последовательностей, относящихся к структурным генам, имеют высокую сложность, тогда как некодирующие последовательности (сателлитные, например, центромерные и теломерные участки) могут на многие килобазы состоять из повторов короткого олигонуклеотида. Различные архиваторы данных работают как раз описывая текст внутри

архивируемого файла подобной программой. Fasta-файлы быстро пакуются любым архиватором, поскольку в длинных текстах над маленьким алфавитом всегда много повторов. Кстати, знаете, чем отличается написанный человеком "случайный" нуклеотидный текст от созданного природой? В рукотворном тексте почти никогда нет длинных повторов вида `tttttttttttt`, поскольку человеку подсознательно кажется, что наличие таких повторов говорит о неслучайности текста. Пример программы-описателя. В нее могут входить два вида элементарных команд, а именно:

Добавить символ в конец текста

Скопировать блок текста из уже сгенерированной последовательности

Сгенерируем с помощью этих команд первый текст (пусть нумерация в строке начинается с единицы):

```
add a
copy 1 символа с позиции 1
copy 2 символа, начиная с позиции 1
copy 4 символа, начиная с позиции 1
copy 4 символа, начиная с позиции 1
add t
copy 1 символ с позиции 13
copy 2 символа начиная с позиции 13
copy 3 символа начиная с позиции 13
```

На выходе `aaaaaaaaaaatttttttt`. Фрагменты низкой сложности часто встречаются в нуклеотидных и аминокислотных последовательностях. Поэтому при поиске сходства, скажем программой BLAST, находки достаточно протяженных фрагментов низкой сложности зачастую занимает верхние строки выдачи, затеняя собой биологически-значимые результаты. Для того, чтобы избежать нерелевантных результатов, программа BLAST имеет встроенный модуль подавления фрагментов низкой сложности (они заменяются на букву `x`). Поэтому Вы никогда не найдете гомологов для последовательности

```
aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa
```

Впрочем, эту опцию можно отключить.

3.8 Посимвольное сжатие (код Хаффмана)

Алгоритмы, предназначенные для решения задач оптимизации, обычно представляют собой последовательность шагов, на каждом из которых предлагается некоторое множество выборов. В *каждом алгоритме* всегда делается выбор, который кажется самым лучшим в данный момент - то есть проводится локально оптимальный выбор. Важным приложением методов

жадного программирования является разработка кодов Хаффмана для сжатия данных. Коды Хаффмана - очень эффективный метод посимвольного сжатия данных, который, в зависимости от характеристик этих данных, обычно позволяет сэкономить от 20% до 90% объема. В жадном алгоритме Хаффмана используется таблица, содержащая частоты появления тех или иных символов. С помощью этой таблицы определяется оптимальное представление каждого символа в виде бинарной строки. Сначала все символы сортируются по убыванию частоты. Затем строится дерево следующим образом: в корень кладется суммарное количество всех встреченных в тексте символов. Полное количество символов хочется разбить на две части так, чтобы они были бы наиболее близки по количеству. В нашем случае (3.1) мы разбиваем на две части так: 45 (символ a) + 55 (остальные символы). Левое ребро помечаем меткой 0, правое - меткой 1. Левое ребро приводит к окончательному символу a, а правое с меткой 1 - к множеству остальных символов. Остальные символы также делим на две наиболее близкие по количеству группы: 25(c+b)(метка 0) и 30(d+e+f) (метка 1). Для каждой группы повторяем процедуру $25 = 13(c) + 12(b)$; $30 = 14(f+e) + 16(d)$. В конце концов получаем дерево, на листьях которого записаны символы. Код каждого символа это последовательность меток ребер, ведущих к символу. Для хра-

Таблица 3.1: Файл данных содержит только символы a-f с указанными частотами, а также указанными для них кодовыми словами фиксированной и переменной длины.

символ	a	b	c	d	e	f	длина записи
количество (тыс.)	45	13	12	16	9	5	100
Фиксированный код	000	001	010	011	100	101	300
код Хаффмана	0	101	100	111	1101	1100	$45 * 1 + 13 * 3 + 12 * 3 + 16 * 3 + 9 * 4 + 5 * 4 = 224$

нения файла, приведенного в качестве примера, понадобится 224 000 битов. К примеру, при использовании кода фиксированной длины для хранения того же файла нам бы понадобилось 300 000 битов (3 бита на символ * 1000 символов). Таким образом, мы видим, код переменной длины позволяет сэкономить 25 % объема. Оптимальный код файла всегда представлен полным бинарным деревом, каждый узел которого имеет по два дочерних узла. Теперь, имея код можно закодировать сообщение, например aacaeafd. Это сообщение будет закодировано следующим образом 00100011011100111. Если мы получили такой набор 0 и 1, сможем ли мы раскодировать этот текст? читаем 0, идем по дереву налево, сразу получаем символ a. Возвращаемся в корень дерева. Снова читаем 0 и опять получаем символ a. Читаем 1 - идем направо. Там узел. Читаем 0 - идем налево читаем 0 опять налево - получаем c. Дальнейший ход раскодирования показан в таблице. Стрелки указывают на выбор движения по дереву: 1 - переход по правой ветке, 0 - переход

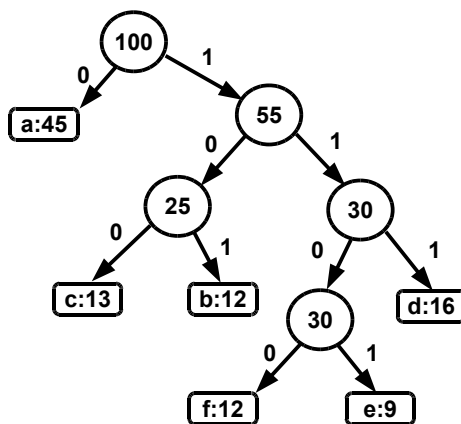


Рис. 3.20: Дерево, соответствующее оптимальному коду, для схемы кодирования на рис 3.19 .

	0 ←	0 ←	1 →	0 ←	1 →	1 →	1 →
	a	a	0 ←	a	1 →	1 →	1 →
по левой ветке			0 ←		0 ←	0 ←	1 →
			c		1 →	0 ←	d
					e	0 f	

Отметим удивительную особенность - несмотря на переменную длину кода (в отличие от генетического кода) мы никогда не сбиваемся. Однако любая ошибка на 0/1 приводит к ошибкам в тексте (может съехать даже весь текст).

Упражнение 25. Найдите частоты встречаемости аминокислот в белках и создайте код Хаффмана для белков. Оцените степень сжатия.

3.9 Сжатие RLE

Распространенный метод сжатия RLE (Run Length Encoding), который заключается в записи вместо последовательности одинаковых символов одного символа и их количества, является подклассом данного алгоритма. Слово "AAAAAA" будет закодировано RLE в виде (A,7), а алгоритмом Лемпеля-Зива в виде "A"(-1,6). этот

3.10 Сжатие информации по Лемпелю-Зиву

Алгоритм Лемпеля-Зива — это универсальный алгоритм сжатия данных без потерь. Алгоритм основан на поиске и кодировании повторов в тексте. Общая идея алгоритма: если в уже прочитанном слове уже встречалась подобная последовательность символов, причем запись о ее длине и смещении

от текущей позиции короче чем сама эта последовательность, то в выходной файл записывается ссылка (смещение, длина), а не сама последовательность. К примеру, строку КОЛОКОЛ_ОКОЛО_КОЛОКОЛЬНИ короче записывается в памяти компьютера так: "КОЛО"(-4,3)"_(-5,4)"0_(-14,7)"ЪНИ". Для сжатия данных текста T длиной n строится суффиксное дерево (время построения $O(n)$). Что такое повтор в суффиксном дереве? Если ребро l ведет в вершину, разветвляющуюся на два ребра, строка, записанная на ребре l , встречается в слове два раза, если в три ребра — три раза и т.п. В изобра-

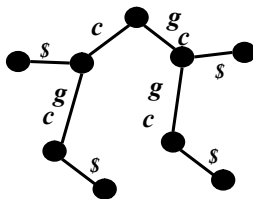


Рис. 3.21: Суффиксное дерево

женном на рис.3.21 суффиксном дереве для слова "gscgc" ребро "gc" ведет в вершину разветвления на два потомка. Подслово gc встречается в слове "gscgc" два раза. То же можно сказать и о подслове "c". Алгоритм Лемпеля-Зива использует on-line подход, он читает и кодирует текст частями. Для начала строится суффиксное дерево части слова. Затем следующая порция текста как бы "протаскивается" через существующее суффиксное дерево — процесс, сходный с продлением ветвей при построении суффиксного дерева. Пусть уже упаковано $i - 1$ символов. Нам надо вычислить следующую команду (позиция, сдвиг) или сгенерировать следующую букву. Берем суффикс $S(i, n)$ и пропускаем его через существующее дерево (как при поиске слов). Либо находим на ребрах дерева максимальное слово, которое является префиксом $S(i, n)$, и тогда ставим команду скопировать это слово. В противном случае генерируем новую букву. Суффиксное дерево можно генерировать одновременно с поиском.

3.11 Сжатие информации по Лемпелю-Зиву-Вельчу

Описанный в предыдущем параграфе алгоритм имеет большую теоретическую значимость, однако достаточно труден для реализации. В 1984 году был предложен еще один алгоритм сжатия, основанный на использовании словарей. В алгоритме, во-первых, используется побитовое кодирование с переменной длиной кода. Во-вторых, по ходу алгоритма строится словарь фраз, ранее встретившихся.

3.11.1 Сжатие

Обозначим W – самое длинное слово в словаре, совпадающее с суффиксом входного потока на очередном этапе работы алгоритма. Алгоритм:

- Иницилируем словарь алфавитом
- Считываем очередной символ a
- Если $a=\$,$ то выдать код k и завершить
- Если Wa есть в словаре присвоить фразе соответствующий код k
- Если нет, выдать код k и записать слово Wa в словарь

Заметим, что код пишется только, если новое слово не найдено, либо в конце текста. Разберем на примере. Пусть входной поток (входная строка) имеет вид $abacabadabacaba\$$. Алфавит состоит из 5 символов $a, b, c, d, \$$ (признак конца строки). Вначале иницилируем словарь алфавитом:

начальный словарь

слово	код
a	0
b	1
c	2
d	3
\$	4

- Пришла буква 'a'. Слово **есть** в словаре. Ничего не пишем, текущий код приравниваем $k=0, W='a'$.
- Пришла буква 'b'. текущее слово равно $W='ab'$. Его **нет** в словаре. Поэтому пишем текущий код (0) в выходной поток, записываем слово 'ab' в словарь. Слово забываем $W='b', k=1$.
- Пришла буква 'a'. Слово **есть** в словаре – ничего не пишем, текущий код $k=0, W='a'$.
- Пришла буква 'c', $W='ac'$. Слова W **нет**, пишем текущий код 0, записываем слово 'ac' в словарь. Сбрасываем $W='c', k=2$
- Пришла буква 'a', $W='ca'$. Слова W **нет**, пишем текущий код 2, записываем слово 'ac' в словарь. Сбрасываем $W='a', k=0$
- Пришла буква 'b', $W='ab'$. Слово W **есть**, запоминаем его код $k=5$

Работа алгоритма показана на рис.3.22 В результате работы алгоритма получаем $LZW(abacabadabacaba\$) = 01025039864$. Сжали 15 символов до 11. Здесь есть доля лукавства, поскольку мы кроме сжатого текста должны были передать и словарь. Однако, как мы увидим ниже, это не обязательно. Кроме того, на самом деле при записи кодов используется техника кодов переменной длины, что позволяет еще существенно усилить сжатие.

кодирование						словарь	
input	Wa		словарь	код	out	слово	код
a	a	есть		0	-	a	0
b	a b	нет	ab → 5	1	0	b	1
a	b a	нет	ba → 6	0	1	c	2
c	a c	нет	ac → 7	2	0	d	3
a	c a	нет	ca → 8	0	2	\$	4
b	a b	есть		5	-	ab	5
a	ab a	нет	aba → 9	0	5	ba	6
d	aba d	нет	ad → 10	3	0	ac	7
a	d a	нет	da → 11	11	3	ca	8
b	a b	есть		5	-	aba	9
a	ab a	есть		9	-	ad	10
c	aba c	нет	abac → 12	7	9	da	11
a	c a	есть		8	-	abac	12
b	ca b	нет	cab → 13	5	8	cab	13
a	b a	есть		6	-		
\$	\$	end			6		

Рис. 3.22: Пример работы алгоритма Лемпеля-Зива-Вельча

3.11.2 Декодирование

Итак, мы получили сжатое сообщение без словаря. Алфавит мы знаем, так что по крайней мере начальный словарь нам уже известен. Будем декодировать и параллельно восстанавливать словарь. Общий алгоритм такой: строим начальный словарь из алфавита. Далее действуем по алгоритму, который одновременно пишет декодированное сообщение и восстанавливает словарь: извлекаем из входного потока первый код k_0 и декодируем его $k_0 \rightarrow X$. Это сделать всегда можно, поскольку первый код соответствует одиночному символу, который есть в словаре.

- Считать очередной код k из сообщения и декодировать его и записать в выходной поток. Из Y извлекаем первый символ y
- Если y – это конец сообщения, то выдать X
- Если фразы Xy нет в словаре, вывести X , а Xy занести в словарь.

Покажем на том же примере как происходит декодирование (рис. 3.23). Итак, у нас есть сообщение 01025039864 и словарь.

- код 0. Он есть в словаре и ему соответствует символ 'a'. Пишем его в выходной поток. $y='a'$
- код 1. Он есть в словаре и ему соответствует символ 'b'. Пишем его в выходной поток и записываем слово $Y='ab'$ в словарь. Присваиваем $y='b'$

Декодирование						
input	X	y	Y	Xy	словарь	out
0		a		a		a
1	a	b	b	ab	ab → 5	b
0	b	a	a	ba	ba → 6	a
2	a	c	c	ac	ac → 7	c
5	c	a	ab	ca	ca → 8	ab
0	ab	a	a	aba	aba → 9	a
3	a	d	d	ad	ad → 10	d
9	d	a	aba	da	da → 11	aba
8	aba	c	ca	abac	abac → 12	ca
6	c	b	ba	cb	cb → 13	ba
4	\$	b				\$

Рис. 3.23: Декодирование

- код 0. Есть в словаре 'a'. Пишем в выходной поток 'a', записываем в словарь 'ba', $y='a'$
- код 2. Есть в словаре 'c'. Пишем 'c', записываем в словарь 'ac', $y='c'$
- Код 5. Есть в словаре 'ab'. Пишем 'ab'. $y='a'$, $X='ab'$
- Код 0. Есть в словаре 'a'. Пишем 'a'. в словарь пишем $X='aba'$
- Код 3. Есть в словаре 'd'. Пишем 'd'. $y='d'$
- ...

Глава 4

Преобразование Барроуза-Уиллера и картирование на геном

В последнее время с появлением мощных методов секвенирования стала более актуальной задача картирования фрагментов на геном. Для того, чтобы эффективно картировать последовательности на геном удобно использовать подготовленный геном, например, геном, представленный в виде суффиксного дерева или суффиксного массива. Но размеры геномов эукариот весьма велики и не позволяют разместить эти конструкции в памяти компьютера. Поэтому необходимы новые структуры. Весьма популярными в настоящее время являются алгоритмы, связанные с преобразованием Барроуза-Уиллера. Эта глава является продолжением главы про алгоритмы на строках.

4.1 Преобразование Барроуза Уиллера

4.1.1 Концепция

Рассмотрим все циклические перестановки текста T (рис.4.1А). Упорядочим их лексикографически и выделим последнюю колонку в этой матрице. Слово, содержащееся в этой колонке и есть результат преобразования. На самом деле вся процедура очень похожа на концепцию суффиксного массива, только мы здесь работаем не с суффиксами, а с циклическими перестановками. Отметим, что, поскольку работаем с циклическими перестановками, то каждая строка устанавливает связь последнего символа с первым. Например, глядя на первый и последний элемент в строке 5 рис.4.1А, мы понимаем, что где-то в тексте есть пара букв si . После сортировки эта связь не теряется. Обозначим:

A	B
1 mississippi\$	12 \$mississippi i
2 ississippi\$m	11 i\$mississipp p
3 sissippi\$mi	8 ippi\$mississ s
4 sissippi\$mis	5 ississippi\$miss s
5 ississippi\$miss	2 ississippi\$m m
6 sissippi\$missi	1 mississippi \$
7 sippi\$missis	10 pi\$mississip p
8 ippi\$mississ	9 ppi\$mississi i
9 ppi\$mississi	7 sippi\$missis s
10 pi\$mississip	4 sissippi\$mis s
11 i\$mississipp	6 sissippi\$missi i
12 \$mississippi	3 ssissippi\$mi i

Преобразование Барроуза-Уиллера: mississippi\$ \rightarrow ipssm\$piissii

Рис. 4.1: Слева – циклические перестановки текста, справа – лексикографически упорядоченные циклические перестановки текста. Последняя колонка матрицы – преобразование Барроуза-Уиллера.

T – исходный текст

S – строка из первых символов в отсортированном массиве циклических суффиксов

BWT – строка из последних символов в отсортированном массиве циклических суффиксов

Будем отслеживать номер появления символа в строке $O(S_i)$. Например, символ s на четвертой позиции BWT является вторым появлением символа s в этой строке: $O(BWT_4) = 2$. Преобразование Барроуза-Уиллера обратимо, т.е. зная только результат преобразования (строку), мы можем восстановить исходную строку. Для этого есть несколько подходов – разной степени вычислительной сложности и сложности алгоритма. Простейший алгоритм (рис. 4.2) состоит в следующем. В начале мы строим столбец из пустых строк. Затем последовательно повторяем две операции:

1. Добавить слева столбец BW
2. Отсортировать полученные слова

После первого этапа мы в столбце имеем упорядоченный набор всех букв в тексте. После второго этапа мы получаем набор всех пар букв, которые встречаются в тексте и во всех его циклических перестановках. На k шаге мы имеем полный набор k -грамм, которые встречаются в тексте и в его циклических перестановках. На последнем этапе мы получаем набор всех циклически переставленных вариантов текста, причем в первой строке будет исходный текст, который будет идти сразу после спец-символа $\$$.

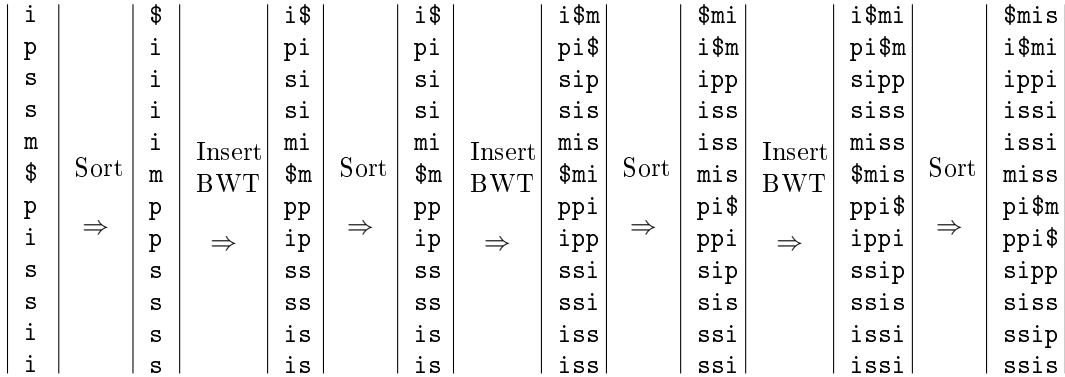


Рис. 4.2: Восстановление исходного текста

4.1.2 Быстрое восстановление исходного текста

Обозначим:

T – исходный текст

S – строка из первых символов в отсортированном массиве циклических суффиксов

BWT – строка из последних символов в отсортированном массиве циклических суффиксов

Будем отслеживать номер появления символа в строке $Oc(S_i)$. Например, символ s на четвертой позиции BWT является вторым появлением символа s в этой строке: $Oc(BWT_4) = 2$.

Лемма i - появление символа в последнем столбце указывает ровно на тот же текстовый символ, что i - появление в первом столбце, например первая буква p в правой колонке (2-я строка) и первая буква p в левой колонке (7-я строка) соответствует одной и той же второй букве p в основном тексте (эта буква подчеркнута на рис.4.1B). Отобразим номера появления символа в последней колонке на первую колонку. Для этого введем массив LF (last-first), в котором

$$LF[k] = i; \quad i : Oc(S_i) = Oc(BWT_k)$$

т.е будем указывать номер строки, в которой номер появления символа в первой колонке равен номеру появления символа в текущей строке и в последней колонке (рис.4.1.2). Кроме того введем массив o номеров появления символа в строке BWT:

$$o[k] = Oc[BWT_k]$$

Введем еще один дополнительный массив размеров в длину алфавита. $C[\alpha] = k$, где k – количество строк, предшествующих первому появлению α в первой позиции суффиксов. Иными словами, это номер строки, в которой впер-

вые появился символ в префиксе минус 1. Длина массивов LF и O равна длине исходной строки. Нетрудно догадаться, что если $o[i] = k$, то

C		LF		O
\$:0	\$BANANA	2	\$BANANA	1
A:1	A\$BANAN	6	A\$BANAN	1
	ANAN\$BAN	7	ANAN\$BAN	2
	ANANAN\$B	5	ANANAN\$B	1
B:4	BANANAN\$	1	BANANAN\$	1
N:5	NANAN\$BA	3	NANAN\$BA	2
	NANAN\$BA	4	NANAN\$BA	3

Рис. 4.3: Дополнительные массивы

$LF[i] = C(BWT_i) + k$. Поэтому построить массив LF можно просто и эффективно. Отметим, что массив LF определяет предыдущий символ – просто потому, что все построение основано на циклических перестановках исходного текста. Отсюда – алгоритм восстановления:

```

1  T=""; r=1; c=BWT[r];
2  while(c!='$'){
3      T=cT;
4      r=LF[r];
5      c=BWT[r];
6  }
```

Пояснения. Строка 1 – инициация. Вначале встаем на первую строку и определяем последний символ строки, поскольку первая строка в суффиксном массиве отвечает единственному символу конца строки и в первой строке всегда лежит исходная строка с приписанным в начале \$. Потом циклически повторяем следующее: приписываем очередной символ в начало текста (строка 4), находим номер строки с предыдущим символом; определяем предыдущий символ. Если мы прочитали символ конца строки (\$), то цикл заканчивается. Покажем, как это работает (рис.4.4)

4.1.3 Поиск образца с помощью преобразования BWT

Найти все вхождения образца P в текст T эквивалентно тому, что найти в суффиксном массиве фрагмент, который отвечает суффиксам, начинающимся с P. На рис.4.5. Обозначим:

$L(W)$ – номер первого суффикса в суффиксном массиве, начало которого совпадает с образцом W

$U(W)$ – номер последнего суффикса, начинающегося с W.

Введем еще один двумерный массив: $O[k, \alpha]$ – количество букв α выше позиции k в BWT. При этом если в позиции k есть буква α , то она также учитывается.

A		
	BWT	LF
\$BANANA	A	2
A\$BANAN	N	6
ANA\$BAN	N	7
ANANA\$B	B	5
BANANA\$	\$	1
NA\$BANA	A	3
NANA\$BA	A	4

B			
	T	r	c
1		1	A
2	A	2	N
3	NA	6	A
4	ANA	3	N
5	NANA	7	A
6	ANANA	4	B
7	BANANA	5	\$

Рис. 4.4: Пошаговый процесс восстановления строки. Слева – исходные данные матрица Барроуза-Уиллера (для сведения, на самом деле не используется), Строка BWT – записанная сверху вниз и массив LF. Справа изменение параметров T – текст, r – номер позиции в BWT, c – текущий символ

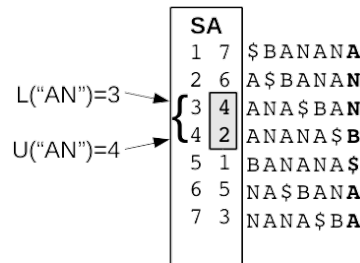


Рис. 4.5: Поиск всех вхождений образца эквивалентен поиску фрагмента суффиксного массива

	BWT	\$	A	B	N
\$BANANA	A	0	1	0	0
A\$BANAN	N	0	1	0	1
ANA\$BAN	N	0	1	0	2
ANANA\$B	B	0	1	1	2
BANANA\$	\$	1	1	1	2
NA\$BANA	A	1	2	1	2
NANA\$BA	A	1	3	1	2

Рис. 4.6: Массив O – количество букв соответствующего типа выше заданной позиции

Лемма Если мы знаем границы суффиксного массива для слова W мы можем легко вычислить границы для слова aW :

$$L(aW) = C(a) + O(L(W) - 1, a) + 1 \tag{4.1}$$

$$U(aW) = C(a) + O(U(W), a) \tag{4.2}$$

Эта лемма позволяет индуктивно вычислять границы суффиксного массива, в пределах которых находятся адреса вхождений.

```

1  a=P[m]; L=C[a]; U=C[a+1]    // a+1 -- следующий символ в алфавите
2  for(i=m-1; i>1; i--){
3      a=P[i];
4      L=C[a]+0[L-1,a]+1;
5      U=C[a]+0[U,a];
6  }
7  return (L,U);

```

Поиск происходит за время, линейное по длине паттерна, причем, в отличие от суффиксного дерева нет необходимости проходить дерево до конца, чтобы перечислить все вхождения – мы сразу получаем список позиций. Если паттерн отсутствует, то верхний индекс будет меньше нижнего.

Поиск образца с заменами При поиске образца с заданным числом замен на самом деле строится дерево замен (рис.4.8). В каждой позиции, кроме первой и последней, допускается замена, если только счетчик числа замен не исчерпан. При этом размер дерева может быть весьма большим, но, поскольку поиск происходит одновременно с построением дерева, есть большой шанс, что на некотором этапе очередная ветка засохнет, т.е. станет ясно, что значения границ в суффиксном массиве противоречива: $L > U$.

4.1.4 Оптимизация памяти

Оценим память, необходимую для всех структур.

- L – для запоминания BWT
- $L \cdot 4$ – для запоминания суффиксного массива SA
- $L \cdot |\Sigma| \cdot 4$ – для запоминания массива O
- $|\Sigma| \cdot 4$ – для запоминания массива C

Таким образом, по размеру памяти использование BWT для поиска образцов ничем не лучше суффиксного дерева. Первое, что можно сделать – это не запоминать весь массив O , а запомнить только его часть. Например, давайте запомним значения этого массива в каждой n -ой позиции, а необходимое значение элемента массива вычислять. Для того, чтобы их вычислить надо просмотреть часть строки BWT, в худшем случае – $n/2$ позиций. Пусть $n = 32$ и нам необходимо определить значение $O[i = 78, A]$. Тогда вычисляем $j_0 = i/32 = 2$; $j_1 = i \bmod 32 = 14$. Теперь можно взять $o = O[j_0 = 2, A]$ и пройти $j_1 = 14$ шагов вдоль строки BWT, начиная с позиции $i_0 = j_0 * n$, подсчитывая встраивающиеся буквы A. Если значение j_1 окажется больше половины n , то просмотр можно вести сверху, вычитая по 1 каждый раз, когда встретится буква A. Вот псевдокод этого вычисления.

Исходные данные

#	SA		BWT	O		
				\$ (0)	A (1)	B (5)
1	11	\$АВАВВАВАВВ	В	0	0	1
2	6	АВАВВ\$АВАВВ	В	0	0	2
3	1	АВАВВАВАВВ\$	\$	1	0	2
4	8	АВВ\$АВАВВАВ	В	1	0	3
5	3	АВВАВАВВ\$АВ	В	1	0	4
6	10	В\$АВАВВАВАВ	В	1	0	5
7	5	ВАВАВВ\$АВАВ	В	1	0	6
8	7	ВАВВ\$АВАВВА	А	1	1	6
9	2	ВАВВАВАВВ\$А	А	1	2	6
10	9	ВВ\$АВАВВАВА	А	1	3	6
11	4	ВВАВАВВ\$АВА	А	1	4	6

Изменения границ L и U по шагам

$$\begin{array}{l}
 1 \quad L(A) = 2; U(A) = 5 \\
 \hline
 L(BA) = C(B) + O(1, B) + 1 = \\
 5 + 1 + 1 = 7 \\
 2 \quad U(BA) = C(B) + O(5, B) = 5 + \\
 4 = 9 \\
 \hline
 L(ABA) = C(A) + O(6, A) + 1 = \\
 1 + 0 + 1 = 2 \\
 3 \quad U(ABA) = C(A) + O(9, B) = \\
 1 + 2 = 3
 \end{array}$$

Изменения границ суффиксного массива

1 \$АВАВВАВАВВ	1 \$АВАВВАВАВВ	1 \$АВАВВАВАВВ
2 АВАВВ\$АВАВВ	2 АВАВВ\$АВАВВ	2 АВАВВ\$АВАВВ
3 АВАВВАВАВВ\$	3 АВАВВАВАВВ\$	3 АВАВВАВАВВ\$
4 АВВ\$АВАВВАВ	4 АВВ\$АВАВВАВ	4 АВВ\$АВАВВАВ
5 АВВАВАВВ\$АВ	5 АВВАВАВВ\$АВ	5 АВВАВАВВ\$АВ
6 В\$АВАВВАВАВ	6 В\$АВАВВАВАВ	6 В\$АВАВВАВАВ
7 ВАВАВВ\$АВАВ	7 ВАВАВВ\$АВАВ	7 ВАВАВВ\$АВАВ
8 ВАВВ\$АВАВВА	8 ВАВВ\$АВАВВА	8 ВАВВ\$АВАВВА
9 ВАВВАВАВВ\$А	9 ВАВВАВАВВ\$А	9 ВАВВАВАВВ\$А
10 ВВ\$АВАВВАВА	10 ВВ\$АВАВВАВА	10 ВВ\$АВАВВАВА
11 ВВАВАВВ\$АВА	11 ВВАВАВВ\$АВА	11 ВВАВАВВ\$АВА

Рис. 4.7: Процесс поиска паттерна АВА в тексте АВАВВАВАВВ. Вверху показаны Суффиксный массив, собственно суффиксы (нужны для справки), массивы BWT и O

1 CalcO(i,a){

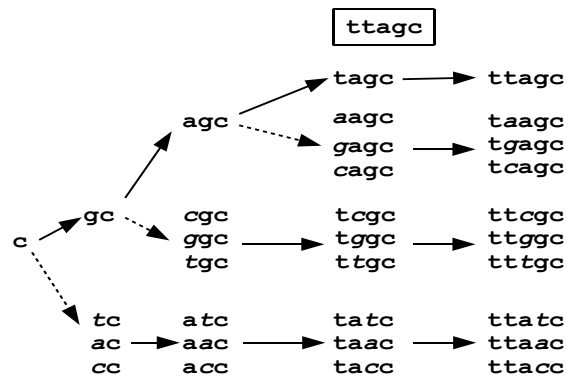


Рис. 4.8: дерево замен

```

2     j0==i/n; j1=i%n;      // здесь целочисленное деление
3     if(j1 <= n/2){
4         0c=0[j0,a];
5         for(k=j0*n+1; k<=i; k++)
6             if(BWT[k] == a) 0c++;
7         return 0c;
8     }
9     else{
10        0c=0[j0,a]
11        for(k=(j0+1)-1; k>=i; k--)
12            if(BWT[k] == a) 0c--;
13        return 0c;
14    }
15 }

```

Обычно используют шаг прореживания, равный $n = 128$. Это дает экономию памяти в 128 раз, но за это приходится платить временем – надо сделать до 64 шагов вдоль BWT. Для случая генома человека для запоминания массива O необходима память $L \cdot 4 \cdot 4/128 = 0.4Gb$. Здесь первая четверка – это размер алфавита, а вторая четверка – размер целого числа в байтах. Использование некоторых дополнительных приемов позволяет сократить размер необходимой памяти еще примерно в 2 раза. Если мы не используем прореживание, то нам нет необходимости помнить саму строку BWT, однако при использовании прореживания мы используем эту строку. Поэтому ее тоже надо помнить. В случае нуклеотидных последовательностей у нас 4-буквенный алфавит. Это значит, что каждый символ может быть закодирован двумя битами, тем самым сжав размер BWT в 4 раза. Геном человека в этом случае займет $0.75Gb$. BWT можно эффективно сжать.

Программистские приемы

Здесь мои представления, как можно эффективно сделать подсчет числа букв. Для того, чтобы вычислить значение элементов массива O в позициях не кратных 128 нет необходимости делать 64 шага, а достаточно сделать определенное небольшое количество побитовых операций, намного меньшее, чем 64. Сначала запишем массив BWT в виде двух битовых массивов – один будет отвечать за то, нуклеотид в соответствующей позиции является ли пурином, а второй – за то, что нуклеотид в соответствующей позиции является "сильным" (рис.4.9).

seq	accgttga	gcatgatc	tggtgcag
SW	01110010	11001001	01101101
RY	10010011	10101100	01101011

Рис. 4.9: битовое представление последовательности. Чтобы узнать, какой нуклеотид находится в позиции, скажем, 5 смотрим на значение битов в этих массивах. Имеем: SW=1, RY=0, т.е. нуклеотид является пиримидином и "сильным", значит, это 'с'

Теперь, используя такое представление, посчитаем количество букв 'а' в первой группе (в первом байте). Для этого построим битовый массив (байт), в котором 1 стоят в позициях, где буква 'а', т.е. позиции, где стоят пурины и при том "слабые" нуклеотиды. Это сделать просто с помощью побитовых операций, которые исполняются за одну команду. Сначала выделим позиции, где стоят слабые, т.е. позиции, где в первом массиве записаны 0 (w) и отметим позиции, где во втором массиве стоят 1, т.е. пурины (r). Побитовое логическое умножение этих двух битовых массивов (логическое "и") даст искомый результат.

$$\begin{aligned}
 w &= \neg SW[0] &= (10001101) & \text{ нам нужны позиции, где 0, поэтому отрицание} \\
 r &= RY[0] &= (10010011) & \text{ нам нужны позиции, где 1} \\
 a &= w \&r &= (10000001)
 \end{aligned}$$

Теперь, чтобы посчитать количество букв 'а' до, скажем, 5-й позиции надо убить единичные биты после этой позиции, т.е. логически умножить на (11111000). Получим единички только там, где нас интересует ((10000000)). Слова-маски – типа (11111000) для отбора нужных позиций можно затабулировать – таких масок будет 64. Осталось подсчитать количество ненулевых битов в слове (10000000). Для слов небольшой разрядности, например, байтов, можно их просто затабулировать. Слова размерностью, например, 64 битовых разряда можно представить, как 8 байтов и по очереди обратиться к соответствующему массиву и просто просуммировать количество ненулевых разрядов. На самом деле, можно хранить количество битов для всех двух-байтовых слов – это займет всего 64к. Тогда подсчет количества ненулевых битов сведется к 4 обращениям к массиву. Итого, для интерполяции значений массива O нам нужно выполнить до двух отрицаний, два

логических умножения и до 4 обращений к массиву – т.е. практически мгновенно.

4.1.5 Сжатие суффиксного массива

Осталось найти способ экономно запомнить сам суффиксный массив. Для этого используется еще одна конструкция – функция последования Ψ . Если у нас есть суффикс $s_k, s_{k+1}, \dots, \$$ и у него есть позиция в суффиксном массиве i , то по определению $\Psi[i]$ – есть позиция следующего суффикса $s_{k+1}, \dots, \$$ (рис.4.10). Поскольку $SA[i]$ – позиция некоторого суффикса, то $SA[\Psi[i]] + 1$ – позиция следующего суффикса. Поэтому по определению:

$$SA[\Psi[i]] = SA[i] + 1$$

Функция Ψ эквивалентна суффиксному массиву – зная функцию Ψ можно

#	SA	Ψ	Ψ^{-1}	Suffix
1	11	3	6	\$ABABVABABV
2	6	8	7	ABABV\$ABABV
3	1	9	1	ABABVABABV\$
4	8	10	8	ABV\$ABABVAB
5	3	11	9	ABVABABV\$AB
6	10	1	10	V\$ABABVABAB
7	5	2	11	VABABV\$ABAB
8	7	4	2	VABV\$ABABVA
9	2	5	3	VABVABABV\$A
10	9	6	4	VV\$ABABVABA
11	4	7	5	VVABABV\$ABA

Рис. 4.10: Функция последования Ψ

легко восстановить суффиксный массив. Она обладает рядом замечательных свойств. Например, она частично отсортирована. Поскольку все значения массива Ψ – разные, то можно определить обратную функцию Ψ^{-1} :

$$\Psi^{-1}[i] = x \Leftrightarrow i = \Psi[x]$$

Тогда, поскольку $SA[\Psi[i]] = SA[i] + 1$, то, обозначив $i = \Psi^{-1}(q)$ можно записать:

$$SA[q] = SA[\Psi^{-1}(q)] + 1$$

Применив это равенство несколько раз можно по части суффиксного массива восстановить его недостающую часть:

$$SA[k] = SA[(\Psi^{-1})^j(k)] + j \quad (4.3)$$

т.е. если мы знаем SA для некоторого значения, то мы можем вычислить значение SA для другого значения. Проблема в том, что вместо массива SA

мы должны помнить массив Ψ^{-1} . Здесь опять приходит на помощь BWT. Мы можем вычислить значение Ψ^{-1} с помощью того, что у нас уже есть:

$$\Psi^{-1}(i) = C(BWT[i]) + O(BWT[i], i) \quad (4.4)$$

Теперь мы запомним в суффиксном массиве только каждый 32 элемент (это значит, что для суффиксного массива нам надо $L/8$ байт). Если нам надо получить элемент суффиксного массива для k не кратном 32 мы применяем функцию Ψ^{-1} до тех пор, пока не встретим доступное значение суффиксного массива:

```

1 Psi_1(k){
2   a=BWT[k]
3   return C[a] + CalcO(k,a);
4 }
5 GetSA(k){
6   j=0;
7   while(k%32 !=0){
8     k=Psi_1(k);
9     j++;
10  }
11  return SA(k)+j;
12 }
```

Заметим, что к суффиксному массиву мы обращаемся только в самом конце поиска и то, если мы найдем искомое вхождение. Если, к тому же, мы программе скажем, что неоднозначное картирование нас не устраивает, то мы только однажды для каждого фрагмента будем искать адрес по указанному алгоритму.

Применив эти приемы, мы, заплатив дополнительным временем, смогли сократить объем необходимой памяти. Для генома человека необходимая память составляет примерно $L/4 + L/8 + L/8 = L/2 = 1.5Gb$.

4.1.6 Итоги

Итак, общая схема алгоритма картирования прочтений на геном состоит из двух этапов.

Препроцессинг На первом этапе мы делаем подготовку генома. Она делается один раз и до смены релиза генома полученные массивы не меняем. Поэтому особенно упираться в скорость работы нет особой необходимости. В результате получаем:

- компрессированный массив BWT
- массив $C[\alpha]$
- прореженный массив $O[i, \alpha]$

- прореженный суффиксный массив *SA*

Эти подготовленные данные можно потом использовать многократно.

Поиск Поиск образца *P* начинаем с конца и с помощью уравнений (4.1) вычисляем границы в суффиксном (настоящем) массиве, где лежат соответствующие суффиксы. Продвигаясь по образцу к началу, границы перестраиваются и в конце концов находим интересующий нас участок суффиксного массива. Если количество образцов больше заданного порога, то репортируем множественное совпадение. Иначе с помощью уравнений (4.3, 4.4) вычисляем значения суффиксного массива. Ниже приведен псевдокод процедуры поиска.

```

1  a=P[m]; L=C[a]; U=C[a+1]    // a+1 -- следующий символ в алфавите
2  for(i=m-1; i>1; i--){
3      a=P[i];
4      L=C[a]+Calc_0(L-1,a)+1;
5      U=C[a]+Calc_0[U,a];
6      if(L > U) return 'not found';
7  }
8  if(U-L > N_Max) return 'too many matchs'
9
10 for(i=L; i<=U; i++)
11     print GetSA(i);

```

Здесь мы обращаемся к описанным ранее функциям интерполяции массивов O ($Calc_0(i, a)$) и SA ($GetSA(i)$)

Глава 5

Алгоритмы для графов

Определение 1. Графом называется пара множеств $G = \{V, E\}$ и отображение $g : V \otimes V \rightarrow \{E, null\}$, то есть каждому элементу из множества E сопоставляется пара элементов из множества V . V называется множеством вершин графа, а E — множеством ребер графа. Если отображение g коммутативно ($g(v_1, v_2) = g(v_2, v_1)$), то граф называется неориентированным, иначе граф ориентированный. Заметим, что обычно функция g обладает свойством $g(v, v) = null$, т.е. нет ребер из вершину в себя же. Иначе мы имеем дело с псевдографом. Граф $G' = \{V', E'\}$ называется подграфом графа G , если $V' \subset V$ и $E' \subset E$.

Здесь мы будем рассматривать только *конечные графы*, т.е. такие графы, в которых множество вершин (и ребер) *конечно*. Ниже приведены еще некоторые определения, которые могут понадобиться.

путь (цепь) последовательность вершин, в которой каждая вершина, кроме последней, соединена с предыдущей вершиной ребром.

цикл путь, в котором первая и последняя вершины совпадают.

простой путь путь, ребра в котором не повторяются.

элементарный путь простой путь, вершины в котором не повторяются.

ациклический граф граф, не содержащий циклов.

взвешенный граф граф, на ребрах которого определены веса.

сильно связная компонента подграф G' данного графа, такой, что для любых вершин $v_1, v_2 \in G'$ есть путь из $v_1 \rightarrow v_2$, а между любой вершиной $v_3 \in G'$ и любой вершиной $v_4 \notin G'$ пути нет.

связная компонента Для неориентированных графов связная компонента. Для ориентированных графов — связная компонента симметризованного графа. Симметризованный граф — это граф, полученный из

ориентированного в котором к которому для каждого ребра (u, v) добавлено ребро (v, u)

двусвязная компонента подграф G' данного графа, такой, что для любых вершин $v_1, v_2 \in G'$ существует цикл $v_1 \in cycle$ и $v_2 \in cycle$.

связный граф граф, в котором для любых вершин v_1, v_2 есть путь из v_1 в v_2 . Связный граф имеет одну связную компоненту.

сильно связный граф ориентированный граф, у которого из любой вершины в любую имеется путь.

дерево связный ациклический граф.

лес граф, у которого все связные компоненты — деревья.

полный граф граф, любые две вершины которого соединены ребром (такой граф — сам себе максимальная клика).

клика подмножество вершин графа, такое, что между любой парой вершин есть ребро (полный подграф).

смешанный граф граф, содержащий ориентированные и неориентированные ребра.

двудольный граф (k-дольный граф) граф, вершины которого можно разбить на 2 (k) не пересекающихся подмножества, так, что не будет ребер, соединяющих элементы одного и того же подмножества.

индекс вершины неориентированного графа число ребер, выходящих из вершины.

индекс вершины ориентированного графа число, зависящее от количества ребер, входящих и выходящих из вершины. Одно входящее в вершину ребро прибавляет к этому числу единицу, одно выходящее — отнимает единицу.

ε -граф граф, в котором число ребер равно $|E| = (1 + \varepsilon)|V|, \varepsilon \ll 1$

5.1 Способы представления графов

Большинство привыкло рассматривать и изображать графы в таком виде, как изображено на рис.5.1. Безусловно, это очень наглядный способ, тем не менее, непонятный компьютеру. Удобно хранить графы в виде таблицы, где в столбцах и строках заданы вершины графа, если на пересечении столбца и строки есть отметка (+ или, удобней, true), между соответствующими ребрами есть ребро. Изображенный нами справа граф в виде таблицы представляется так:

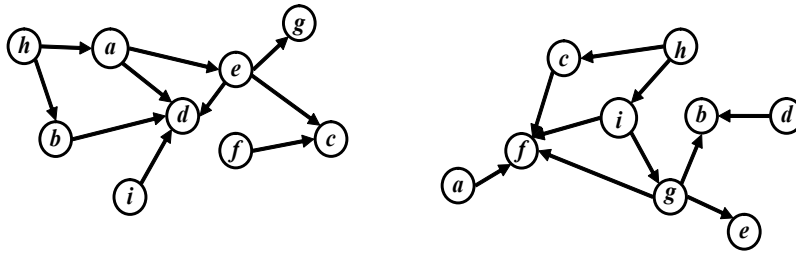


Рис. 5.1: Пример графического представления графов. Одинаковые (изоморфные) ли графы представлены на рисунке?

	a	b	c	d	e	f	g	h	i
a	0	0	0	1	1	0	0	0	0
b	0	0	0	1	0	0	0	0	0
c	0	0	0	0	0	0	0	0	0
d	0	0	0	0	0	0	0	0	0
e	0	0	1	1	0	0	1	0	0
f	0	0	1	0	0	0	0	0	0
g	0	0	0	0	0	0	0	0	0
h	1	1	0	0	0	0	0	0	0
i	0	0	0	1	0	0	0	0	0

Это ориентированный граф. Если бы граф был неориентированным (рис.5.2), таблица приняла бы несколько иной вид:

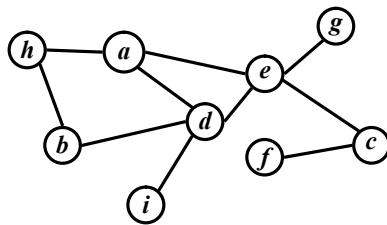


Рис. 5.2: неориентированный граф

	a	b	c	d	e	f	g	h	i
a	0	0	0	1	1	0	0	1	0
b	0	0	0	1	0	0	0	1	0
c	0	0	0	0	1	1	0	0	0
d	1	1	0	0	1	0	0	0	1
e	1	0	1	1	0	0	1	0	0
f	0	0	1	0	0	0	0	0	0
g	0	0	0	0	1	0	0	0	0
h	1	1	0	0	0	0	0	0	0
i	0	0	0	1	0	0	0	0	0

Если на пересечении строки и столбца подобной таблицы указано true (1), значит, из вершины в строке идет ребро в вершину в столбце. Подобная таблица называется матрицей смежности графа. Матрица смежности полного графа будет включать только единицы. Как хранить такую матрицу в памяти? Либо в виде двумерного массива, либо закодировав с применением алгоритма Лемпеля-Зива. Ни первое, ни второе в данном случае удобным не кажется. Для записи и хранения графа можно составить список ребер этого графа. Смежными назовем вершины, соединенные ребром (в неориентированном графе). Для приведенного выше примера неориентированного графа список смежных вершин будет выглядеть так: ad; ae; bd; ed; eg; es; fc; ha; hb; id;. Если граф ориентированный, порядок в парах важен. Еще граф можно хранить в виде нескольких связанных списков смежных вершин. Из каждой вершины графа пишем связный список, состоящий из вершин, смежных данной вершине (вершин, с которыми данная вершина связана ребром, для неориентированного графа, и вершин, в которые идет “стрелка” из данной вершины для ориентированного графа, *список соседей*). Тогда представление ориентированного графа будет следующим:

вершина	Список смежных вершин	вершина	Список смежных вершин
<i>a</i>	<i>d</i> → <i>e</i>	<i>f</i>	<i>c</i>
<i>b</i>	<i>d</i>	<i>g</i>	
<i>c</i>		<i>h</i>	<i>a</i> → <i>b</i>
<i>d</i>		<i>i</i>	<i>d</i>
<i>e</i>	<i>d</i> → <i>g</i> → <i>c</i>		

Для неориентированного графа:

вершина	Список смежных вершин	вершина	Список смежных вершин
<i>a</i>	<i>d</i> → <i>e</i> → <i>h</i>	<i>f</i>	<i>c</i>
<i>b</i>	<i>d</i> → <i>h</i>	<i>g</i>	<i>e</i>
<i>c</i>	<i>e</i> → <i>f</i>	<i>h</i>	<i>a</i> → <i>b</i>
<i>d</i>	<i>a</i> → <i>b</i> → <i>e</i> → <i>i</i>	<i>i</i>	<i>d</i>
<i>e</i>	<i>a</i> → <i>c</i> → <i>d</i> → <i>g</i>		

Графы являются удобным языком дискретной математики, и многие задачи формулируются в форме задач на графах. В частности, многие задачи биоинформатики имеют графовую постановку. Примеров можно привести множество. Задача выравнивания последовательностей формулируется как задача поиска кратчайшего пути в графе. Задача поиска ортологов может быть сформулирована как задача поиска клики в графе. Задача построения белковых семейств формулируется как задача поиска связной компоненты в графе или как задача поиска двусвязной компоненты, или как задача поиска клик в графах. Задачи кластеризации, восстановления последовательности по результатам секвенирования, задачи, связанные со скрытыми Марковскими моделями и т.д. и т.п. Часто графы имеют гигантские размеры. Например, при выравнивании двух последовательностей длиной по тысяче символов возникает граф размером порядка миллиона вершин. Поэтому построение эффективных алгоритмов на графах является актуальной задачей.

5.2 Обход графа в ширину

Многие задачи на графах требуют просмотра всех вершин в некотором порядке. Есть два основных способа обхода графа — обход в ширину и в глубину. Где у графа ширина, а где глубина, сказать сложно, но под обходом графа в ширину подразумевают просмотр всех соседей выбранной вершины, потом соседей их соседей и т.п. Для обхода графа в ширину используется такая структура данных, как очередь (first input — first output). Для работы алгоритма добавим к вершинам атрибут — цвет вершины:

- белый — вершина, которую еще не видели
- серый — вершина, которая находится в стадии обработки
- черный — вершина, для которой обработка завершена

Вначале все вершины белые, поскольку на них еще ни разу не смотрели. Затем работаем по алгоритму:

1. Берем произвольную вершину, помещаем в очередь и красим в серый цвет.
2. Изымаем очередную вершину из очереди (красим в черный цвет) и все смежные не пройденные (белые) вершины помещаем в конец очереди и красим в серый цвет
3. Если очередь не пуста, переходим к 2.
4. Если очередь пуста, но остались не просмотренные вершины, берем из них произвольную, помещаем в очередь и переходим к 2
5. Если очередь пуста и все вершины графа черные, заканчиваем работу — граф обойден

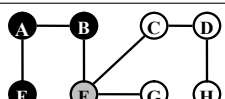
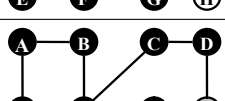
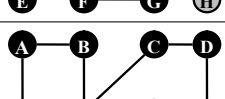
В таблице 5.2 показан порядок работы алгоритма обхода графа в ширину. Псевдокод алгоритма поиска в ширину:

```

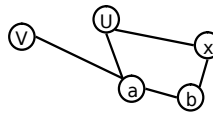
1.  WideSearch(){
2.    for(each v in V){
3.      if(v.color==white) VertexWideSearch(v);
4.    }
5.  }
6.  VertexWideSearch(Vertex v){
7.    Queue q;
8.    q.put(v);
9.    v.color=gray;
10.   while(!q.empty){
11.     v=q.get();
12.     v.color=black;
13.     for(w in Neighbours(v)){
14.       if(w.color==white){
15.         q.put(w);
16.         w.color=gray;
17.       }
18.     }
19.   }
20. }
```

Пояснение. Здесь две подпрограммы. Первая ищет очередную белую вершину, чтобы с нее стартовать. Вторая подпрограмма осуществляет собственно поиск в ширину. Строка 2 означает, что мы просматриваем все вершины. Поскольку множество вершин не обязательно организовано в массив, то цикл не объявлен явным перечислением значений индекса. Строка 7 объявляет очередь, в которой будут стоять вершины. Стартовую вершину помещаем в очередь. Пока очередь не пуста (цикл в строках 10-20) выполняем следующее. Снимаем вершину из очереди, красим в черный, просматриваем всех соседей (цикл в строках 13-17) и помещаем белых соседей в очередь. Время работы алгоритма $T = O(|V| + |E|)$. Действительно, мы смотрим по одному разу на каждую вершину в цикле строки 2 и смотрим на каждое ребро в цикле строки 13. Видно, что наиболее подходящим способом хранения графа для этой задачи является хранение списков соседних вершин. Впрочем, часто бывает, что ребра графа не заданы явно, а определяются “на лету”, как функция вершин. В этом случае время работы алгоритма может быть больше — $T = O(|V| \cdot |V|)$, поскольку цикл строки 13 заменится на просмотр всех вершин графа и появится проверка дополнительного условия являются ли вершины v и w соседями. Описанный алгоритм может помочь в решении задач анализа графов с целью выявления их структуры и вычисления ряда характеристик. С помощью поиска в ширину можно находить связанные компоненты графа. Когда в какой-то момент обхода графа очередь становится пустой, мы просмотрели одну связную компоненту. Если в графе

Таблица 5.1: Пример обхода графа в ширину.

Граф	Очередь	Действие
	A	Берем первую попавшуюся вершину (A), кладем в очередь и красим серым.
	BE	Снимаем вершину A из очереди, помечаем черным, кладем в очередь соседей — E и B
	EF	Снимаем первую вершину из очереди (B), красим в черный цвет, а ее соседей (F) кладем в очередь.
	F	Снимаем первую вершину из очереди (E), красим в черный цвет, а ее соседей (а их нет!) кладем в очередь.
	CG	Снимаем первую вершину из очереди (F), красим в черный цвет, а ее соседей (CG) кладем в очередь.
	GD	Снимаем первую вершину из очереди (C), красим в черный цвет, а ее соседей (D) кладем в очередь.
	D	Снимаем первую вершину из очереди (G), красим в черный цвет, а ее соседей (а их нет!) кладем в очередь.
	H	Снимаем первую вершину из очереди (D), красим в черный цвет, а ее соседей (H) кладем в очередь.
		Очередь пуста и нет не просмотренных вершин! Мы обошли весь граф

остались не просмотренные вершины, значит, имеется, как минимум, еще одна связная компонента. Просматриваем ее, ищем следующую. Таким образом можно находить и перечислить все связные компоненты графа. Поиск в ширину применим для нахождения расстояния между вершинами, а также поиска кратчайшего пути между вершинами. **Дано:** произвольный граф $G = \{V, E\}$, $V > 1$. **Требуется** определить расстояние между вершинами $u \in V, v \in V$, а также указать кратчайший из путей, связывающих u и v (если такие пути существуют). **Решение:** стартуя из вершины v , делаем обход графа в ширину, пока не дойдем до u . Если между событиями “ v в очереди” и “ u в очереди” очередь оставалась пустой, пути между v и u нет. Если между указанными событиями очередь ни разу не была пуста, мы нашли путь между v и u , притом этот путь кратчайший. Допустим, ищем кратчайший путь между вершинами v и u в таком графе:



Обходим граф в ширину, начиная с v :

1. В очереди вершина v
2. В очереди a
3. В очереди вершины b и u , останавливаем поиск, так как дошли до вершины u .

Между шагом 1 и шагом 3 очередь не была пуста, значит, u и v принадлежат к одной связной компоненте, и между ними есть путь (или несколько путей). Мы нашли один из таких путей с помощью обхода по ширине. Этот путь $v - a - u$ (смотрим, какие вершины побывали в очереди между v и u). Итак, расстояние между v и u — два ребра или одна вершина. Мы решили правую часть задачи — нашли расстояние между вершинами v и u . **Докажем**, что найденный нами с помощью поиска в ширину путь является кратчайшим путем между v и u .

Доказательство. от противного — пусть мы нашли не кратчайший путь и есть путь еще короче. Тогда вершина на кратчайшем пути встретилась бы раньше при просмотре, поскольку каждый раз мы смотрим на *ближайших* соседей (на смежные вершины) той вершины, то изымается из очереди. Возникло противоречие — утверждение доказано. \square

Упражнение 26. *Напишите псевдокод для поиска связной компоненты.*

Упражнение 27. *Напишите псевдокод для поиска минимального расстояния между двумя вершинами.*

5.3 Обход графа в глубину

Идея этого метода — идти по графу от вершины к вершине вглубь, пока не “упремся”. Затем возвращаемся на этап раньше, снова идем в глубину и т.д. В отличие от обхода в ширину здесь используется не очередь, а стек (first input — last output). При обходе графа будем красить вершины также: белым — вершины, которые еще не обработаны, серым, вершины, которые находятся в процессе обработки, черным — уже обработанные вершины. Итак, алгоритм:

1. Берем произвольную вершину и помещаем в стек (красим в серый цвет)
2. Смотрим на вершину, лежащую на вершине стека. Берем произвольного белого соседа этой вершины и, покрасив в серый цвет, кладем в стек. Если нет белых соседей, то переходим к п. 3., иначе повторяем п.2.
3. Если не осталось белых смежных вершин, то берем вершину из стека и заканчиваем ее обработку (красим в черный) и переходим к 2. Если стек пуст, переходим к п.4.
4. Стек пуст и остались не просмотренные (белые) вершины, то берем из них произвольную, помещаем в стек и переходим к п.2. Иначе завершаем обход.

В таблице 5.3 показана работа алгоритма поиска в глубину. Псевдокод алгоритма обхода в глубину:

```
1. DepthSearch(){
2.   for( each v in V){
3.     if(v.color=white)
4.       VertexDepthSearch(v);
5.   }
6. }
7. VertexDepthSearch(v){
8.   v.color=gray;
9.   for( each w in Neighbours(v) ){
10.    if(w.color=white)
11.      VertexDepthSearch(w);
12.   }
13.   v.color=black;
14. }
```

Пояснения к псевдокоду. Подпрограмма `DepthSearch` просматривает все вершины, и, если находит белую вершину, то начинает просмотр в глубину, начиная с этой вершины (`VertexDepthSearch`). Подпрограмма поиска в глубину с заданной стартовой вершиной работает рекурсивно. Она красит текущую вершину в серый цвет (строка 8), затем находит среди соседей белую вершину и использует ее в качестве стартовой для дальнейшего поиска в

Таблица 5.2: Пример обхода графа в глубину.

Граф	Стек	Действие
	a	Берем произвольную вершину (a) и кладем в стек
	ae	На вершине стека лежит вершина a, поэтому берем ее произвольного белого соседа (вершину e) и кладем в стек.
	aed	На вершине стека лежит вершина e, поэтому берем ее произвольного белого соседа (вершину d) и кладем в стек.
	ae	У вершины d нет белых соседей, поэтому заканчиваем ее обработку (красим в черный)
	aeg	На вершине стека лежит вершина e, поэтому берем ее произвольного белого соседа (вершину g) и кладем в стек
	ae	На вершине стека лежит вершина g, и у нее нет белых соседей. Снимаем ее со стека и красим в черные цвет.
	aec	На вершине стека лежит вершина e, поэтому берем ее произвольного белого соседа (вершину c) и кладем в стек
	ae	На вершине стека лежит вершина c, и у нее нет белых соседей. Снимаем ее со стека и красим в черные цвет.
	a	На вершине стека лежит вершина e, и у нее нет белых соседей. Снимаем ее со стека и красим в черные цвет.
		На вершине стека лежит вершина a, и у нее нет белых соседей. Снимаем ее со стека и красим в черные цвет. Стек пуст.
	h	Находим произвольную белую вершину, кладем в стек.

глубину (строка 11). Когда все соседи обработаны (а каждый из них прошел в глубину до конца), вершину можно покрасить в черный цвет (строка 14). Также как и для поиска в ширину время поиска в глубину составляет $O(|V| + |E|)$. Работа со стеком, как и с очередью, занимает $O(|V|)$. Время на рассмотрение ребер $O(|E|)$.

Упражнение 28. *Перепишите псевдокод поиска в глубину в виде цикла вместо рекурсии. Здесь Вам кроме стандартных методов стека (`push` — положить элемент в стек и `pop` — снять элемент с вершины стека) понадобится метод `peek` — прочитайте элемент с вершины стека не изымая его со стека.*

Разумеется, поиск в ширину и поиск в глубину не являются самостоятельными задачами. Обычно изменение цвета вершины сопровождается дополнительными действиями, например записью вершины в какой-нибудь список, вычислением каких-либо характеристик и т.п. Важно, что при обоих типах поиска мы проходим по соседним вершинам, т.е. анализируем связность графа. Еще раз стоит напомнить, что при решении реальных задач приходится иметь дело с графами достаточно большого размера — тысячи, миллионы и более вершин. Так что нарисовать картинку и на нее посмотреть не удастся.

5.3.1 Типы ребер графа

В результате поиска в глубину получается дерево (или лес из нескольких деревьев) — дерево обхода. Каждая вершина графа v имеет время окончания обработки $t(v)$, т.е. тот шаг алгоритма, при котором вершина стала черной. После завершения обхода графа в глубину все ребра графа оказываются разбитыми на два множества — ребра дерева (леса) обхода и остальные ребра. В ориентированном графе ребра подразделяются на четыре разных типа (рис.5.3):

Ребра деревьев Ребра деревьев поиска в глубину (на рисунке толстые).

Прямые ребра. Вершины u, v принадлежат дереву, но (uv) не является ребром дерева, и $t(u) > t(v)$. Другими словами, это ребра (uv) не являющиеся ребрами дерева и соединяющие вершину u с ее потомком v в дереве поиска в глубину (на рисунке — двойные).

Обратные ребра. Вершины u, v принадлежат дереву и $t(u) < t(v)$. Другими словами, это ребра (uv) , соединяющие вершину u с ее предком v в дереве поиска в глубину (на рисунке — пунктирные).

Перекрестные ребра. Все остальные ребра. Они могут соединять вершины одного и того же дерева поиска в глубину, когда ни одна из вершин не является предком другой, или соединять вершины в разных деревьях (на рисунке — серые).

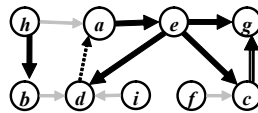


Рис. 5.3: Классификация ребер графа после обхода в глубину. Толстые ребра — ребра дерева обхода, двойные — прямые ребра, пунктирные — обратные ребра, серые — перекрестные ребра.

5.3.2 Построение покрывающего дерева

Покрывающим дерево графа называется совокупность деревьев обхода в глубину, соединенных особыми перекрестными ребрами. Чтобы построить такую структуру, нужно сначала сделать обход графа в глубину, создав лес обхода. Если лес обхода состоит из одного дерева, оно и является покрывающим деревом. Если в лесе обхода несколько деревьев, посмотрим, нет ли перекрестного ребра из вершины одного дерева обхода в корень другого дерева обхода. Такое ребро, направленное из вершины A одного дерева обхода в корень другого, включаем в покрывающее дерево. Новым корнем становится A . Когда все такие ребра исчерпаны, покрывающее дерево готово. Оно состоит из всех деревьев обхода графа, соединенных особыми перекрестными ребрами

5.3.3 Поиск связных компонент

Мы уже обсуждали выше, как искать связные компоненты с помощью обхода графа в ширину. Этот метод подходит для любых графов, тем не менее, на ориентированных графах можно применять другой способ поиска связных компонент. Для нахождения связных компонент в ориентированном графе следует ко всем ребрам графа добавить перекрестные ребра (в отличие от покрывающего дерева это не обязательно ребра, идущие в корень дерева из леса обхода в глубину).

5.3.4 Поиск циклов

Поиск в глубину позволяет найти циклы в графе, причем несколькими способами. 1) Если при обходе в глубину мы положили в стек вершину v и увидели среди ее соседей серую вершину, в графе есть цикл, верно и обратное: если граф циклический, обязательно на каком-либо этапе обхода возникнет встреча с серым соседом. В приведенном примере такая ситуация возникла при обработке вершины d . Впрочем с таким же успехом можно применять 2) Циклы также можно найти, построив покрывающее дерево, как в ориентированном, так и в неориентированном графе. В ориентированном графе для нахождения циклов строим покрывающее дерево. Обратные ребра замыкают циклы. В неориентированном графе также строим покрывающее дерево и ищем все ребра, не принадлежащие к нему. Эти

ребра порождают циклы. Отметим, что задача поиска всех циклов в графе весьма трудоемкая, поскольку количество всех возможных циклов очень велико. Действительно, если в графе есть два цикла с общей вершиной, то можно построить еще один цикл, объединяющий их. Если таких циклов несколько, то возникает большое комбинаторное разнообразие. Тем не менее можно эффективно решать задачу поиска хотя бы некоторых циклов.

Упражнение 29. *Напишите псевдокод для алгоритма поиска циклов в графе. Здесь по-видимому Вам понадобится дополнительная функция в структуре стека, позволяющая просмотреть стек не изымая вершины со стека.*

Упражнение 30. *Двусвязной компонентой неориентированного графа называется такое множество вершин, что через любые две вершины графа проходит цикл такой, в котором вершины проходятся по одному разу. Придумайте алгоритм поиска двусвязных компонент.*

5.3.5 Некоторые приложения к биоинформатике

Есть множество аминокислотных последовательностей (скажем, UNIPROT). Задача состоит в том, чтобы найти белковые семейства. Здесь мы рассмотрим простейший наивный подход к этой задаче. На самом деле процедура кластеризации последовательностей в семейства содержит множество дополнительных деталей. Итак, первым шагом кластеризации является определение расстояний между последовательностями (например, с помощью BLAST или Смита-Ватермана). Теперь можно построить граф: вершины — последовательности, а ребра проводятся, если вес выравнивания превышает заданный порог. Разумеется, в реальности этот порог выбирается из статистического анализа и/или ребро проводится с учетом структуры выравнивания. Здесь есть (и применяется) множество различных вариантов. Теперь мы имеем граф. Можно предположить, что в этом графе белковому семейству отвечает *клика*, т.е. такой подграф, в котором все вершины соединены со всеми (полный подграф). Действительно, в пределах семейства каждый белок должен быть похожим (в смысле выбранной меры сходства) на любой другой белок семейства. Однако такой подход таит в себе ряд проблем. Во-первых, алгоритмическая проблема — поиск клики в графе является сложной задачей, для которой не существует эффективного алгоритма (см. главу 4). Во-вторых, есть содержательная проблема — клика не определяется однозначно — одна вершина (ребро, подграф) может принадлежать сразу нескольким кликам (см. рис. 4.3). Поэтому поиск клик в этой задаче не применяют. Представляется более подходящим подход, основанный на поиске связной компоненты. Такой метод часто применяют. Действительно, связная компонента не обладает недостатками клики — есть эффективные алгоритмы поиска, каждая вершина принадлежит одной и только одной связной компоненте. Для каждой связной компоненты можно оценить ее “качество” — степень насыщенности ребрами. Если количество ребер равно

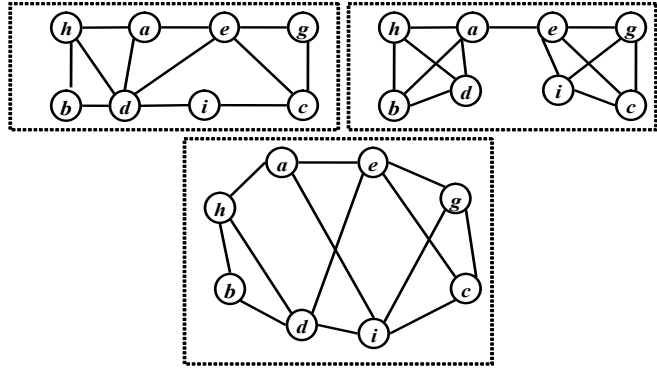


Рис. 5.4: граф для семейства белков. а) Вершина ребро hd принадлежит двум кликам, вершина e принадлежит двум другим кликам. б) связный граф, который очевидно соответствует двум разным семейства. в) Граф, который соответствует одному семейству

$|V| \cdot (|V| - 1) / 2$, то, как не трудно догадаться, эта связная компонента является кликой. Правда, появляется другая содержательная проблема. Посмотрите на рис.5.4б и 5.4в. Эти графы имеют одинаковое количество вершин и ребер, но интуитивно ясно, что граф б) представляет два семейства. Поэтому связная компонента не лучший способ искать семейства. Кстати, этот случай соответствует выравниванию на рис.5.5, т.е. случаю, когда многодоменные белки объединяют в одно семейство группы однодоменных белков. Поэтому представляется более рациональным использовать для этой зада-

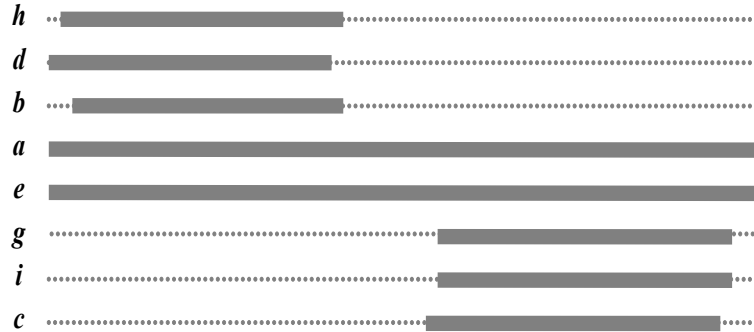


Рис. 5.5: Выравнивание, продуцирующее граф на рис.5.4б

чи поиск двусвязной компоненты. Такой подход допускает ситуацию, когда вершина может принадлежать двум или более двусвязным компонентам. Этот недостаток на самом деле является достоинством, поскольку позволя-

ет находить многодоменные белки.

5.4 Задача Эйлера

Многим известно, что теория графов зародилась в 18 веке после того, как Леонарда Эйлера заинтересовала задача о семи Кенигсбергских мостах: как пройти по всем мостам, не проходя ни по одному из них дважды (рис.5.6 4.5). Отсюда и задача Эйлера, которая формулируется следующим образом: *найти на графе цикл, проходящий через все ребра, причем по каждому ребру один раз*. Эйлеровым циклом называется цикл в графе такой, что

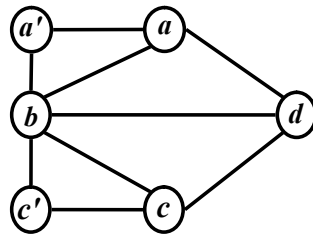


Рис. 5.6: Схема мостов Кенигсберга. Ребра графа — мосты. Вершины a' и c' являются дополнительными, поскольку между a и b и b и c есть по два моста, поэтому нет прямого отображения мостов на граф.

он проходит через все ребра графа по одному разу. Если в графе есть эйлеров цикл, то граф называется эйлеровым, или, говорят, что граф обладает свойством эйлеровости. Понятно, что эйлеровым может быть только связный граф. Для любого связного графа можно сразу сказать, обладает он свойством эйлеровости, или нет. *В связном неориентированном графе есть эйлеров цикл тогда и только тогда, когда индексы всех вершин четные. В связном ориентированном графе есть эйлеров цикл тогда и только тогда, когда индексы всех вершин равны нулю.* Доказательство необходимости лежит на поверхности: рисуя эйлеров граф, мы пройдем через каждое ребро лишь однажды, а значит, если ребро привело нас в некую вершину, сможем выйти из нее только по другому ребру. Из этого следует, что, сколько ребер входит в вершину, столько из нее и выходит — значит, вершина четная. Достаточность четности индекса вершин доказывается ниже. Если мы не требуем условия замкнутости пути, возникает эйлеров путь — путь, содержащий все ребра графа и такой, что по нему можно пройти, не проходя одно и то же ребро дважды. Требования для существования эйлерова пути мягче, чем для эйлерова цикла. В неориентированном графе не более чем две вершины могут быть нечетными (начало и конец), все остальные обязаны быть четными для существования эйлерова пути в графе. Аналогично и с орграфом — все вершины должны иметь индекс 0, не более чем две вершины могут иметь нечетные индексы (отрицательные у начала и

положительные у конца графа). Как искать эйлеров путь и эйлеров цикл в графе, даже если мы знаем, что он там есть? С эйлеровым путем поступим, как любят делать математики (помните задачу про чайник?) — проведем ребро между нечетными вершинами и будем искать эйлеров цикл. Если между этими вершинами уже есть ребро, то вводим дополнительную вершину и соединяем нечетные (ненулевые) вершины ребрами с этой новой вершиной. Таким образом, перед нами стоит задача поиска эйлерова цикла в графе.

5.4.1 Алгоритм поиска Эйлерова цикла в графе

При поиске Эйлерова цикла основывается на наблюдении. Если в эйлеровом графе удалить ребра, образующие цикл, либо останется эйлеровым, либо распадется на эйлеровы графы. Действительно, после удаления такого цикла четность вершин либо не изменится (если эти вершины не участвуют в удаляемом цикле), либо уменьшится на 2, если вершина принадлежит циклу. Поэтому четность вершин после удаления цикла не изменится. Единственное, что может случиться неприятного — потеря связности графа. Нетрудно доказать от противного, что процедуру удаления можно повторять до тех пор, пока в графе не останется ребер. На основе этого наблюдения можно построить алгоритм. Но предварительно докажем лемму.

Лемма 1. *Если в графе все вершины четные, то через любую вершину проходит цикл.*

Доказательство. Начнем просто строить такой цикл. Пометим все ребра как не пройденные. Затем, начиная с выбранной вершины идем в глубину по не помеченным ребрам. При просмотре в глубину (в отличие от стандартного просмотра в глубину) мы помечаем не вершины, а ребра. В конце концов при таком просмотре мы упрямся в ситуацию, когда из очередной вершины не идет не помеченных ребер. Покажем, что эта вершина совпадает со стартовой. Действительно, допустим это некоторая промежуточная вершина. Тогда либо из нее не выходит вообще ребер (эта вершина - тупик), и поэтому она не четная (одно ребро пришло и ни одного не вышло). Либо из этой вершины выходят несколько помеченных ребер. Тогда это значит, что через эту вершину мы проходили, т.е. входили в нее и выходили из нее. Следовательно количество помеченных ребер, выходящих из этой вершины - четное плюс еще одно ребро, по которому мы только что пришли, т.е. количество ребер - нечетное. Таким образом при этом просмотре мы пришли в ту вершину, с которой стартовали. Сам цикл — это просто содержание стека в тот момент, когда мы упрямся в стартовую вершину. \square

1. Помечаем все ребра, как не пройденные.
2. Выбираем произвольную вершину v .
3. Идем в глубину, начиная с вершины v . При этом помечаем не вершины а ребра. Поиск в глубину заканчивается, если из очередной вершины

не выходит белых ребер. Согласно лемме мы получили цикл, который можно записать в список, который можно зациклить, установив ссылку из последнего элемента списка на первый элемент списка.

4. В построенном списке вершин находим вершину v , из которой выходят не помеченные ребра. Если такая вершина не найдется, то найденный цикл является эйлеровым. Действительно, если в построенном цикле нет вершин с не помеченными ребрами, а где-то в графе остались такие, то граф не связный.
5. Повторяем поиск, начиная с вершины v . Найденный цикл записываем в новый связный список и вставляем новый список в основной список на место вершины

Псевдокод алгоритма:

```

1.   Euler(G){
2.     List L;
3.     Vertex v=<any vertex>;
4.     do{
5.       List L1;
6.       SimpleCycle(v,L1);
7.       insert(L,L1,v);
8.       for(v=L.root; v!=null; v=v.next)
9.         if(v has white edges) break;
10.    }while(v!=null)
11.    return L;
12.  }
13.  SimpleCycle(v,L1){
14.    L1.add(v);
15.    for(e in edges_from_v){
16.      if(e.color==white){
17.        e.color=black;
18.        if(e.v1==v)
19.          SimpleCycle(e.v2,L1);
20.        else
21.          SimpleCycle(e.v1,L1);
22.      }
23.    }
24.  }
25.  }

```

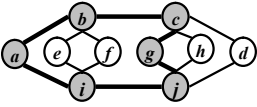
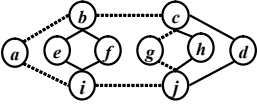
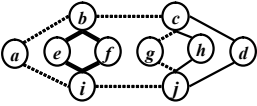
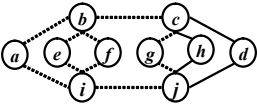
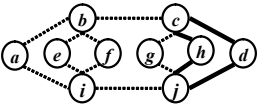
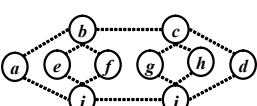
Пояснения: Алгоритм содержит процедуру построения простого цикла (после строки 13). В строке 2 создаем пустой список L. Затем выбираем произвольную вершину v (строка 3). Здесь не описано, как мы это делаем, поскольку это зависит от организации данных. Далее в цикле перебираем

все вершины, имеющие белые ребра (цикл в строках 4-10). Внутри цикла мы организуем новый список L1, куда запишем результаты очередного поиска в глубину. Затем вставляем найденный цикл (список L1) в исходный список L на место v (строка 7). Далее в полном списке L ищем вершину v с белыми ребрами (строки 8-9). Если такая вершина не найдена, то заканчиваем работу (строка 10). Процедура поиска цикла, начиная с заданной вершины (строка 13). Мы добавляем стартовую вершину в список. Затем ищем белое ребро из этой вершины. Как только находим такое ребро (строка 16), красим это ребро в черный цвет и с противоположной вершины этого ребра продолжаем поиск (строки 18-21). Для оптимизации работы можно на вершинах поддерживать количество исходящих белых ребер и список вершин, которые имеют белые ребра — чтобы не искать каждый раз вершину с белыми ребрами. Здесь надо понять, что очередную вершину мы пишем в список на входе в поиск цикла. В таблице 5.4.1 показано как это работает.

5.4.2 Биоинформатика и задача Эйлера: секвенирование геномов

Современная стратегия секвенирования геномов подразумевает массовое прочтение случайных фрагментов генома и затем сборку генома. Характерная длина фрагмента — 300-800 нуклеотидов при методе Сэнгера и 75-150 при секвенировании NGS, например на приборах Illumina. Дальше возникает задача восстановления полной последовательности генома. Поскольку читаются случайные фрагменты, то многие из них пересекаются (имеют общие фрагменты). Эти общие фрагменты и являются ключом для сборки генома. Итак, можно построить граф: вершины — прочитанные фрагменты, ребра проводятся, если фрагменты имеют общий кусок. Задача заключается в том, чтобы провести на этом графе путь, проходящий через все вершины. К сожалению, секвенирование случайных фрагментов может привести к не связному графу, поскольку по случайным причинам некоторые участки генома могут оказаться непокрытыми прочитанными фрагментами. Другая проблема — задача поиска пути, проходящего через все вершины графа — задача Гамильтона, является трудной задачей, для которой не существует в принципе эффективного алгоритма (эта проблема будет обсуждаться в заключительной главе книги). Однако задачу сборки генома можно сформулировать как задачу Эйлера (точнее — как задачу обойти ребра графа). Будем строить граф так, чтобы ребра соответствовали фрагментам, а вершины — суффиксам и префиксам фрагментов. Такой граф называется графом де-Брюина. На рис. 5.7 показан пример графа де-Брюина. Вершины — слова, которые являются префиксами исходящих ребер и одновременно суффиксами входящих ребер. Современная технология сборки геномов существенно использует тот факт, что покрытие генома (сколько раз в среднем прочитан каждый нуклеотид) достаточно велико — 20X и более. Для сборки генома делаем сначала контринтуитивную вещь — разбиваем все прочтения на достаточно короткие слова (размером 15-20 букв). При этом мы теряем существенную информацию, но упрощаем сборку. Потом мы вспомним

Таблица 5.3: Пример поиска Эйлера цикла.

граф	список вершин	
		Стартуя с вершины <i>a</i> проходим в глубину по ребрам. В момент покраски вершины <i>i</i> в черный цвет у нее серый сосед - стартовая вершина <i>a</i> . В списке L1: <i>abcgjia</i>
	<i>abcgjia</i>	Помечаем ребра как пройденные.
	<i>abcgjia</i>	В списке вершин остались вершины с ненулевым числом ребер. Стартуя с вершины, скажем, <i>i</i> , находим цикл: <i>iebf</i>
	<i>abcgj iebfi a</i>	и вставляем цикл <i>iebf</i> на место.
	<i>abcgj iebfi a</i>	В списке вершин остались вершины с ненулевым числом ребер. Стартуя с вершины, скажем, <i>c</i> , находим цикл: <i>cdjhc</i> .
	<i>abcd jhcdj iabei a</i>	Вставляем в список на место вершины <i>c</i> цикл <i>cdjhc</i> . Помечаем ребра как пройденные.

о том, что слова пришли из фрагментов. Теперь из слов строим граф де-Брюина. В этом графе ищем однозначные пути. Если появляется ветвление, то мы обрываем путь. На рис.5.8 показан вариант такой ситуации. Мы выбираем пути *ABCD, DEFG, DHIG, GJK*. Этим путям соответствуют некоторые слова (последовательности). Такие последовательности будем называть *миниконтигами*. Эти объекты во-первых, как правило, существенно длиннее фрагментов, а во-вторых отражают заведомо однозначные прочтения. Такое построение исключает необходимость искать перекрытия прочтений – поскольку перекрытия прочтений – это и есть общие слова фрагментов. Теперь пора вспомнить о том, что миниконтиги пришли из более протяженных фрагментов. Поэтому многие неоднозначности можно разрешить с использованием этого знания. В результате получаются *контиги*. Но в ре-

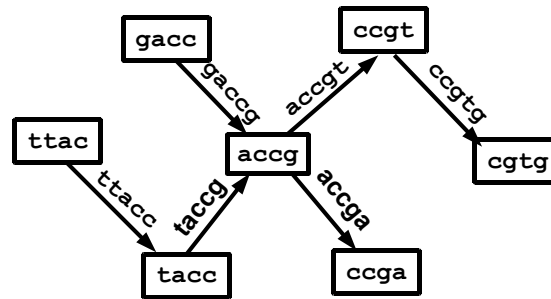


Рис. 5.7: Пример графа де-Брюина – ребра слова, а вершины суффиксы и префиксы. Путь на этом графе соответствует словам. Например, на этом графе можно увидеть слова `ttaccggtg`, `ttaccga`, `gaccgtg`, `gaccga`.

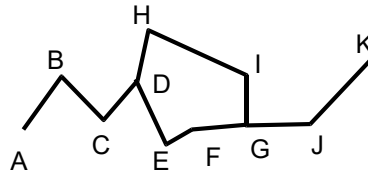


Рис. 5.8: Граф де-Брюина. Слова на этом графе не нарисованы, но они имеются в виду. На этом графе выбираем однозначные пути без ветвлений

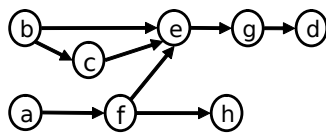
альных геномах часто встречаются повторы, которые длиннее фрагментов. Поэтому не удается разрешить последовательность следования контигов. Для того, чтобы разрешить хотя бы часть оставшихся повторов применяют дополнительные эксперименты. Во-первых можно использовать данные парно-концевых прочтений. Во вторых, можно использовать технику `mate-pair`:

Геномная ДНК подвергается фрагментации (любым методом: с помощью редко щепящих или супер-редкощепящих эндонуклеаз рестрикции, ультразвука, небулизацией). К фрагментам ДНК лигируют адаптеры, содержащие сайты рестрикции специальных других, достаточно часто щепящих эндонуклеаз. Эти адаптеры содержат сайт рестрикции, а также могут быть биотинилированы. Полученные фрагменты ДНК, фланкированные адаптерами закольцовываются. Биотин используют для дальнейшего изолирования целевых спаренных концевых фрагментов на сорбенте со стрептавидином. Затем кольцевая молекула ДНК обрабатывается эндонуклеазой, сайты связывания которой содержатся в адаптерах. Образуются свободные парноконцевые участки. Перед секвенированием, к этим спаренным концевым фрагментам пришивают адаптеры, которые содержат последовательности для отжига праймеров ПЦР. Фрагменты считываются с двух концов. В результате получаем 2 контига, которые можно объединить

в скэффолд.

5.5 Топологическая сортировка

Дан ориентированный граф. **Задача:** перечислить вершины в таком порядке, чтобы из вершины с большим номером не шло ребер в вершины с меньшим номером, иными словами, надо спроецировать вершины на прямую так, чтобы стрелки графа шли только направо. Если две вершины в графе независимы (нельзя сказать, какая из них стоит “позже” другой), их взаимное положение в списке неважно. Однако они могут быть упорядочены опосредовано, например, если между двумя вершинами нет непосредственно ребра, но между ними есть вершина X такая, что из первой вершины есть ребро в X , а из вершины X есть ребро во вторую вершину. Задача топологической сортировки является важной составной частью общей задачи планирования. Если у Вас есть список взаимосвязанных дел, то надо составить порядок их исполнения. Какие-то дела можно делать только после того, как завершатся некоторые другие дела, а выполнение каких-то других дел не связано друг с другом. Эту ситуацию можно представить в виде ориентированного графа. Вершины — это дела, ребра — связи. Будем проводить ребро из вершины A в вершину B , если дело B можно делать только после дела A . Ясно, что связи должны быть устроены так, чтобы в полученном ориентированном графе не было циклов. Высокий пилотаж бюрократии заключается в создании циклов в таком графе дел. Основное правило проецирования: Если из вершины A в вершину B идет ребро, то вершина B обязательно должна стоять в списке позже A . Следовательно, для графа из двух вершин вида $A \rightarrow B$ список будет AB . Если две вершины в графе независимы (нельзя сказать, какая из них стоит “позже” другой), их взаимное положение в списке неважно. Спроецируем согласно приведенным правилам такой граф:



Сортированный список будет:

a b c f e g h d

или

b a f c h e g d

Оба варианта допустимы. Дело в том, что вершины a и b , h и e , h и g и другие независимые вершины можно менять местами. Главное, чтобы вершина f всегда стояла после a , e после c , d после f — вершины в списке должны находиться согласно направлению ребер графа. Можно придумать

несколько вариантов таких записей, и все они будут верными. Таким образом, ациклический ориентированный граф всегда можно спроецировать на список вершин, при этом часто верных проекций может быть несколько. Описанный выше способ проецирования графа на список его вершин называется *топологической сортировкой*. Смотря на граф, на направление ребер в нем, мы сортируем список его вершин надлежащим образом. Рассмотрим топологическую сортировку на наглядном примере. У нас есть неупорядоченный список предметов одежды (список вершин графа). Пусть у нас есть проблема — надо надеть следующие вещи:

- Пиджак
- Брюки
- Галстук
- Рубашка
- Ремень
- Ботинки
- Носки
- Часы

Ясно, что некоторые вещи нельзя надеть раньше других. Например, попробуйте надеть носки после того, как надеты ботинки. Построим ориентированный граф по принципу “что за чем можно надевать”. Если вещь a надевается раньше вещи b , ребро идет из a в b . Получаем такой граф (рис.5.9).

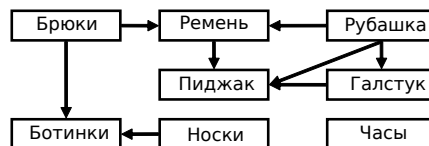


Рис. 5.9: Граф процесса одевания.

5.5.1 Алгоритм топологической сортировки

1. Берем произвольную вершину графа, начинаем с нее поиск в глубину
2. При поиске в глубину, когда вершина чернеет, добавляем ее в список, за ней почерневшую следующей и так далее.
3. Когда все вершины почернели, полученный список разворачиваем.

Реализуем приведенный алгоритм на нашем “одежном” примере. Берем неупорядоченный список вершин и начинаем обход в глубину соответствующего графа с произвольной вершины, к примеру, “брюки”. Обход в глубину дает нам следующий порядок вершин:

Ботинки → Пиджак → Ремень → Брюки

Остались не просмотренные вершины. Стартуем с “рубашки”. Получаем:

Ботинки → Пиджак → Ремень → Брюки → Галстук → Рубашка

Продолжая процедуру получаем наконец:

Ботинки → Пиджак → Ремень → Брюки →
Галстук → Рубашка → Носки → Часы

Разворачивая список находим, наконец порядок надевания вещей:

Часы → Носки → Рубашка → Галстук → Брюки →
Ремень → Пиджак → Ботинки

Топологическая сортировка часто используется при планировании каких-то мероприятий, экспериментов. В частности, она используется в качестве предварительного этапа алгоритма динамического программирования поиска оптимального пути в графе.

Упражнение 31. *Напишите псевдокод процедуры топологической сортировки.*

5.6 Поиск оптимального пути в графе

Есть множество побуждающих примеров для постановки поиска оптимального пути в графе. Это может быть и самый короткий путь (в смысле километров), и самый быстрый — он не обязательно самый короткий — на коротком пути может быть пробка, и самый дешевый — дорога может быть такой, что проезд по ней обойдется большим ремонтом. Этого рода задачи возникают сплошь и рядом при маршрутизации пакетов в компьютерной сети, в частности в Интернете. В любом случае нам важно, что:

1. дан ориентированный граф;
2. на ребрах графа определены веса (т.е. граф взвешенный);
3. Вес пути определяется как сумма весов ребер, через которые проходит этот путь.

Задача состоит в том, чтобы найти путь минимального (или максимального) веса. Рассмотрим один из примеров, возникающих в биоинформатике, когда необходимо найти оптимальный путь. В результате секвенирования ДНК (например, генома) получен ряд последовательностей, которые надо собрать в полный геном. Тогда поступим следующим образом. Построим графа Графа де-Брюина. Каждое ребро графа - слово (фрагмент, прочитанная последовательность) Каждая вершина, из которой выходит ребро - суффикс слова на 1 букву короче слова на ребре. Каждая вершина в которую входит ребро - префикс слова на 1 букву короче слова на ребре. Вес равен длине общей части. Можно составить словарь всех возможных слов длины L из алфавита. При чтении генома каждый ряд характеризуется покрытием - количество его встречаемости. Напротив каждого слова можно поставить его покрытие. Глядя на эти числа, можно выбросить редко встречаемые риды. Тогда задача сборки генома также формулируется как задача поиска пути с наибольшим весом. Любое ветвление означает повтор. Если в таком графе находится путь без ветвлений, то при включении его концевых вершин мы получаем контиг. Итак, для каждого слова в словаре мы можем посмотреть его суффикс и попробовать поискать другое слово с совпадающим префиксом. Если такое соответствие единственное, то мы их объединяем и получаем контиг. Если же таких слов больше, чем один, то возникает неоднозначность или другими словами, повтор. Проблема повторов на практике решается парноконцевыми чтениями.

Задача: дан связный ориентированный взвешенный (т.е. на каждом ребре графа написано неотрицательное число) граф.

Найти: путь между двумя вершинами, имеющий минимальный (максимальный) вес.

5.6.1 Динамическое программирование для поиска оптимального пути

Эта задача решается по-разному в зависимости от некоторых дополнительных условий. Итак, задача. Дан *ациклический* ориентированный взвешенный граф. Задача — найти путь минимального (максимального) веса между двумя вершинами (B и E). Рассмотрим некоторую вершину v . Допустим, мы знаем веса w_i оптимальных путей до всех вершин, непосредственно предшествующих вершине v , т.е. для всех вершин $u_i \in \{prev(v)\}$, из которых есть ребра $e_i = (u_i, v)$. Тогда несложно вычислить вес оптимального пути, ведущего в вершину v . Действительно, вес пути, проходящего через вершину u_i равен $w(B, u_i, v) = w_i + w(e_i)$. Тогда наименьший вес пути из начальной вершины B в вершину v равен:

$$w(B, v) = \min_i \{w(B, u_i, v)\} = \min_i \{w_i + w(e_i)\}$$

Это соображение позволяет построить алгоритм поиска веса оптимального пути из вершины B в вершину E . Припишем каждой вершине v вес $w(v)$, равный весу оптимального пути из начала B в эту вершину. Ясно, что вес, приписанный вершине B равен 0, поскольку путь из B в B не содержит ни одного ребра. Далее мы можем просмотреть всех соседей вершины B и определить вес оптимального пути в них. Когда мы обработаем вершину E , мы найдем вес оптимального пути. *На самом деле это не совсем так.* Рассмотрим фрагмент графа на рис.5.10. Если мы будем подсчитывать вес

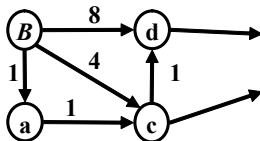


Рис. 5.10: Побудительный пример для необходимости топологической сортировки.

для вершины d , не определив веса у вершины c , то мы положим ее равной 8. Хотя очевидно, что оптимальный путь в вершину проходит через вершины a и c . С другой стороны, если мы будем сначала считать оптимальный вес для вершины c , а потом для вершины d , то мы получим вес 5, что также не является оптимальным весом. А правильный порядок вычисления весов такой: сначала идет вершина a , потом вершина c и только затем вершина d . Таким образом, порядок обхода вершин должен быть таковым, чтобы сначала были обработаны вершины, которые не зависят от последующих вершин, иными словами необходимо провести топологическую сортировку вершин. При этом все вершины, которые в результате топологической сортировки перед вершиной B , равно и вершины, стоящие после вершины E , можно исключить из рассмотрения. Из сказанного ясно, почему для такого алгоритма необходимо, чтобы граф был ациклическим. Итак, алгоритм:

1. Провести топологическую сортировку вершин.
2. Отбрасываем все вершины, предшествующие стартовой вершине B и следующие за вершиной E .
3. В цикле по сортированным вершинам определяем вес вершин, как минимум сумм весов непосредственно предшествующих вершин и ребер. После обработки последней вершины (E) мы определим вес оптимального пути.

Описанный алгоритм находит только вес оптимального пути, но при этом у нас нет оптимального пути (а хотелось бы...). Для того, чтобы найти оптимальный путь надо, чтобы в вершинах графа был еще один атрибут — ссылка на вершину на которой достигается минимум:

$$\pi(v) = \operatorname{argmin}_i (w(u_i) + w(u_k(v)))$$

Обратим внимание на то, что сеть оптимальных переходов образует дерево. Алгоритм, основанный на рекурсивном вычислении целевой функции, называется алгоритмом динамического программирования. Алгоритмы динамического программирования применяются не только при оптимизации на графах, но и при оптимизации других объектов. Для того, чтобы восстановить путь надо провести обратный просмотр. Начиная с вершины E . Ее атрибут $\pi(E)$ скажет нам откуда мы пришли в вершину E . Перехода в соответствии со ссылкой в новую вершину $v_1 = \pi(E)$ определим ее предшественника $v_2 = \pi(v_1)$, и т.д.

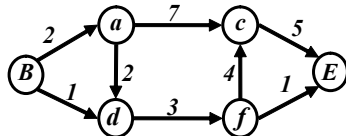
Псевдокод алгоритма:

```

1.  DinProg(B,E){
2.      Vertex v[]=Topol_Sort();
3.      int ib, ie;
4.      for(i=0; i<v.length; i++){
5.          if(v[i]==B) ib=i;
6.          if(v[i]==E) ie=i;
7.      }
8.      if(ie < ib) return “no path”;
9.      v[ib].weight=0;
10.     for(i=ib+1; i <= ie; i++){
11.         v[i].weight=INFINITY;
12.         for(each w in v.previous){
13.             ww=w.weight+weight(w,v[i]);
14.             if(ww < v[i].weight) {
15.                 v[i].weight=ww;
16.                 v[i].pi=w;
17.             }
18.         }
19.     }
20.     Stack path;
21.     w=E;
22.     while(w!=B){
23.         path.push(w);
24.         w=w.pi;
25.     }
26.     return path;
27. }
```

Пояснения. Программа состоит из трех этапов: топологическая сортировка и анализ исходных данных (сроки 2-8). В строке 8 определяется существует ли путь из вершины B в вершину E . В строках 9-19 происходит поиск оптимального веса. Строка 9 означает, что оптимальный вес пути из начала в начало равен 0. Строки 11-17 находят оптимальный вес для вершины $v[i]$, в строке 16 запоминаем оптимальный путь. Строки 20-25 по-

священы обратному просмотру и восстановлению пути. Путь запоминается в стеке. В таблице 5.6.1 приведен пример работы алгоритма динамического программирования для графа:



Заметим, что нам на самом деле неважно, положительные или отрицательные веса заданы на ребрах. Более того, этот алгоритм применим как для минимизации веса пути, так и для его максимизации — стоит только заменить операцию `min` на `max`.

5.6.2 Оценка времени работы динамического программирования

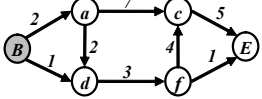
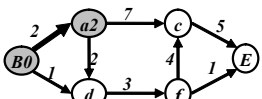
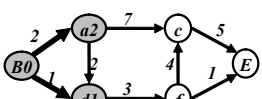
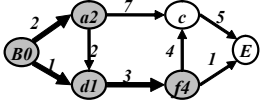
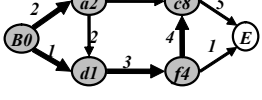
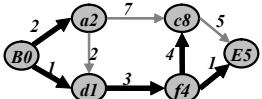
Оценим время работы алгоритма. Топологическая сортировка требует времени порядка $T = O(|V| + |E|)$. Цикл в строках 4-7 обрабатывается за время $T = O(|V|)$. Время работы цикла 10-19 составляет $T = O(|V| + |E|)$, поскольку в этом цикле (при правильной организации данных) просматривается каждая вершина и каждое ребро по одному разу. Время обратного прохода равно длине оптимального пути. В худшем случае $T = O(|V|)$. Итого:

$$T = O(|V| + |E|) + O(|V|) + O(|V| + |E|) + O(|V|) = O(|V| + |E|)$$

5.6.3 Биоинформатические применения поиска оптимального пути в графе. Выравнивание

Задача выравнивания последовательностей возникает во многих задачах биоинформатики и является одной из основных классических задач биоинформатики. На самом деле эта задача возникла намного раньше как задача определения редакционного расстояния. Редакционным событием будем называть одно из действий: удаление символа из слова, вставка символа в слово, замена символа. Редакционным расстоянием между словами называется минимальное количество редакционных событий переводящих одно слово в другое. Задача о редакционном расстоянии возникает во многих контекстах. Например, система исправления опечаток может пытаться построить выравнивание напечатанного слова с одним из слов словаря и тем самым исправить опечатку. Аналогичные алгоритмы используются при поиске плагиата. В биологии задача имеет трактовку: найти минимальное количество элементарных эволюционных событий. Для решения задачи построим граф. Вершины графа соответствуют сопоставлению префиксов слов. Ребра графа — редакционные события, которые удлиняют один или оба префикса, т.е. добавляют символ в одну последовательность (а в другую — добавляют

Таблица 5.4: Пример динамического программирования для поиска оптимального пути в графе.

	$B(0); a(\infty,?);$ $d(\infty,?); c(\infty,?);$ $f(\infty,?); E(\infty,?);$	Весы всех вершин, кроме начальной не определены.
	$B(0); a(2, B);$ $d(\infty,?); c(\infty,?);$ $f(\infty,?); E(\infty,?);$	Обрабатываем вершину а. У нее единственный предшественник — вершина В. $w(a) = w(B) + w(Ba) = 0 + 2 = 2$
	$B(0); a(2, B);$ $d(1, B); c(\infty,?);$ $f(\infty,?); E(\infty,?);$	Вершина d — два предшественника а и В. $w(d) = \min\{(w(B) + w(Bd)), (w(a) + w(ad))\} = \min\{(0 + 1), (2 + 2)\} = 1$
	$B(0); a(2, B);$ $d(1, B); c(\infty,?);$ $f(4, d); E(\infty,?);$	Вершина f — один предшественник d. $w(f) = w(d) + w(df) = 1 + 3 = 4$
	$B(0); a(2, B);$ $d(1, B); c(8, f);$ $f(4, d); E(\infty,?);$	Вершина с — два предшественника а и f. $w(c) = \min\{(w(a) + w(ac)), (w(f) + w(fc))\} = \min\{(4 + 4), (2 + 7)\} = 8$
	$B(0); a(2, B);$ $d(1, B); c(7, f);$ $f(4, d); E(f, 5);$	Вершина E — два предшественника с и f. $w(E) = \min\{(w(c) + w(cE)), (w(f) + w(fE))\} = \min\{(5 + 7), (4 + 1)\} = 5$

пропуск символа) или в другую последовательность, или в обе последовательности. Тогда можно построить граф для определения редакционного расстояния и найти кратчайший путь в этом графе. Итак, определим редакционное расстояние для слов `aacggatcg` и `aaggattcg`. Граф редакционных расстояний организован в виде таблицы (рис.5.11). Горизонтальные переходы порождают символы в горизонтальной последовательности, или (что то же самое) пропуск в вертикальной последовательности. Вертикальные стрелки соответствуют вставке в вертикальной последовательности (пропуск в горизонтальной последовательности). Этим переходам всегда соответствует одно событие редактирования. Диагональные стрелки соответствуют сопоставлению символов (они либо совпадают, либо различаются). В зависимости от того, совпадают соответствующие символы или нет вес этих строк будет либо 0 (нет события редактирования) либо 1 (редактирование — замена). Важно, что в этом графе есть строки и столбцы, которые не соответствуют символам последовательности (отмечены серым). Эти строки (ряды) соответствуют краевым делециям. Любой путь на этом графе

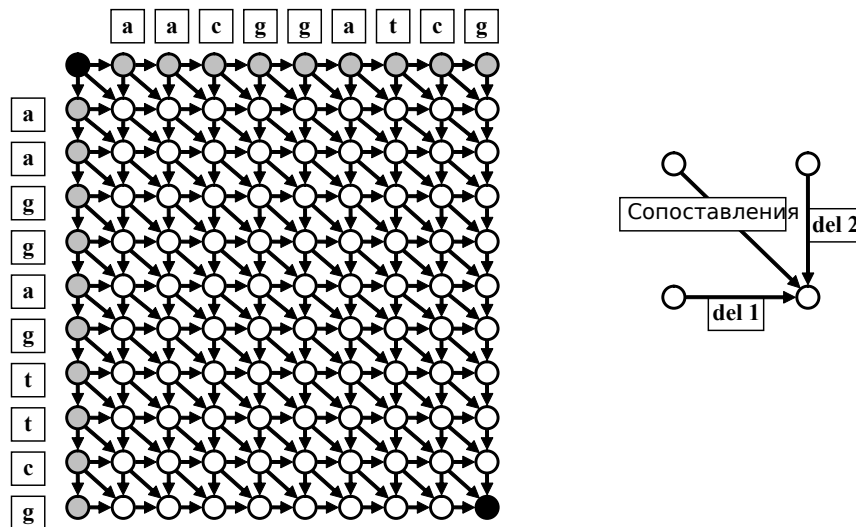


Рис. 5.11: Динамическое программирование для определения редакционного расстояния. Черным отмечены стартовая и финальная вершины.

соответствует некоторой последовательности редакционных событий, или, что то же самое, выравниванию.

Упражнение 32. *Какому выравниванию соответствует путь: сначала по верхней горизонтальной строке до конца, затем по правому вертикальному столбцу до конца?*

Упражнение 33. *Есть ли необходимость проводить для этого графа топологическую сортировку?*

Таблица 5.5: Пример динамического программирования для поиска оптимального пути в графе.

		a	a	c	g	g	a	t	c	g
	0	→1	→2	→3	→4	→5	→6	→7	→8	→9
a	↓1	↘0	→1	→2	→3	→4	→5	→6	→7	→8
a	↓2	→1	↘0	→1	→2	→3	→4	→5	→6	→7
g	↓3	→2	↓1	↘1	↘1	→2	→3	→4	→5	→6
g	↓4	→3	↓2	↓2	↘1	↘1	→2	→3	→4	→5
a	↓5	→4	↘3	↓3	↓2	↘2	→1	→3	→4	→5
t	↓6	→5	↓4	↘4	↘4	→3	↘2	→1	→2	→3
t	↓7	↓6	↓5	↘5	↘5	↓4	↓3	↘2	→3	→4
c	↓8	↓7	↓6	↘5	↘6	↓5	↓4	↓3	↘2	→3
g	↓9	↓8	↓7	↓6	↘5	↘6	↓5	↓4	↓3	↘2

Таблица ?? показывает вычисление редакционного расстояния между последовательностями. Оно оказывается равным 2 — есть вставка в последовательность 1 и вставка в последовательность 2:

aacggat-cg
aa-ggattcg

5.6.4 Полукольцо. Динамическое программирование над полукольцом

Алгоритм динамического программирования использует рекурсию

$$w(B, v) = \min_{i: u_i \in \{prev(v)\}} \{w(B, u_i, v)\} = \min_i \{w_i + w(e_i)\}$$

Операцию минимума по всем вариантам на самом деле есть обобщение бинарной операции минимума: берем первые два элемента, вычисляем минимум, затем вычисляем минимум от полученного значения и следующего элемента т.д.:

$$\min_i \{w_i + w(e_i)\} = \min(v_1 + e_1, \min(v_2 + e_2, \min(\dots)))$$

Минимум на самом деле является бинарной операцией, такой же как, скажем сложение. Если мы введем странные обозначения, а именно вместо + будем писать \otimes , а вместо min будем писать \oplus то можно переписать это выражение в виде:

$$\begin{aligned} & \min(v_1 + e_1, \min(v_2 + e_2, \min(\dots))) = \\ & = (v_1 \otimes e_1) \oplus (v_2 \otimes e_2) \oplus \dots \oplus v_n = \\ & = \bigoplus_i (v_i \otimes e_i) \end{aligned}$$

Таблица 5.6: Примеры полуколец.

Множество	\oplus	\otimes	$\mathbf{0}$
целые неотрицательные числа	+	\times	0
квадратные матрицы	+	\times	нулевая матрица
true, false	or	and	false
$\mathbb{R}, +\infty$	min	+	$+\infty$
$\mathbb{R}, -\infty$	max	+	$-\infty$

Итак, у нас при рекурсии используем две операции — сложение и взятие минимума. При этом нам важны некоторые вещи, а именно хотелось бы, чтобы, например, было свойство $a \oplus (b \oplus c) = (a \oplus b) \oplus c$. Принципиально важно, чтобы вес пути был суммой весов подпутей и, наверное, ряд других свойств. Итак,

Определение 2. *Полукольцом называется множество M на котором определены две бинарные операции \oplus и \otimes обладающие свойствами:*

- $a \oplus b$ — коммутативность по "сложению"
- $a \oplus (b \oplus c) = (a \oplus b) \oplus c$ — ассоциативность по "сложению"
- $a \otimes (b \otimes c) = (a \otimes b) \otimes c$ — ассоциативность по "умножению"
- $\exists 0 : a \oplus 0 = a$ — Существование 0
- $a \otimes (b \oplus c) = a \otimes b \oplus a \otimes c$ — правая дистрибутивность
- $(a \oplus b) \otimes c = a \otimes b \oplus a \otimes c$ — левая дистрибутивность

Таким образом, динамическое программирование можно применять не только для поиска минимума или максимума, но и для решения других задач. Вопрос только каких... Возьмем, к примеру, полукольцо неотрицательных целых чисел с естественными операциями. На ребрах определим значение веса = 1. В стартовой вершине V определим вес 1. Тогда оптимальный вес, подсчитанный в вершине E , даст нам число путей из V в E . Действительно, путь из V в E единственный. Допустим все на всех вершинах u_i , предшествующих вершине v определено количество путей из V . Тогда количество путей, ведущих в вершину v равно сумме чисел путей в предшествующие вершины. На самом деле существуют более изощренные способы применения полуколец при подсчете различных величин. Упомянем подсчет статистических сумм или полных вероятностей. Эти задачи подробно обсуждаются в теории скрытых Марковских моделей.

5.6.5 Поиск минимального пути в графе, содержащем циклы. Алгоритм Дейкстры

Представим себя водителем автомобиля, пытающегося проехать от ФББ до станции метро "Университет" в условиях сильных пробок. Составим план

дорог, соединяющих эти места (рис. 5.12). Кружками на нем показаны вершины “графа объезда” — наиболее важные пункты (ФББ и метро — тоже вершины). Стрелки — это ребра графа, показывающие всевозможные прямые и окольные пути между пунктами-ребрами. На ребрах, вообще говоря, есть веса (варианты: сложность объезда, длина пути, вероятность пробки). Все это вместе образует ориентированный взвешенный граф в котором есть циклы. Чтобы “с комфортом” добраться от ФББ до метро, хорошо было бы найти оптимальный путь в этом графе. Выше мы упоминали, что дина-

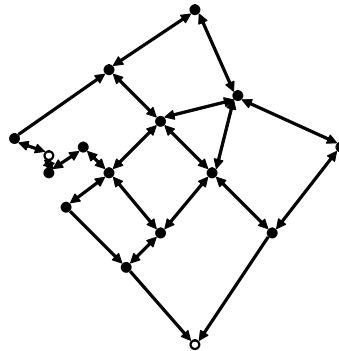


Рис. 5.12: Схема проезда в районе Университета.

мическое программирование на циклическом графе невозможно, поскольку для такого графа нельзя сделать топологическую сортировку, служащую во многом основой динамического программирования. Для построения такого алгоритма нам потребуется новая структура данных — *очередь с приоритетами*. Элементы очереди с приоритетами имеют дополнительный параметр — приоритет. Вход элемента в очередь такой же, как в обычной очереди. Однако из очереди снимается элемент с наивысшим приоритетом. Если есть несколько элементов с одинаковыми приоритетами, то снимется тот элемент, который пришел раньше. Это то, как работает обычная очередь (например, в сбербанке), когда ветераны имеют право внеочередного обслуживания. Если одновременно придет несколько ветеранов, то они образуют свою очередь, которая обслуживается прежде всего. В компьютерной очереди с приоритетами таких приоритетов может быть много. Нетрудно догадаться, что очередь с приоритетами может быть реализована с помощью двоичного дерева поиска (лучше красно-черного). Основой для построения двоичного дерева поиска является операция сравнения, которая позволяет сравнивать два любых элемента и в зависимости от результатов сравнения кладет или ищет элемент в правом или в левом поддереве. Введем понятие времени входа элемента в очередь. Это просто порядковый номер появления элемента в очереди. В начале счетчик равен нулю. Когда приходит очередной элемент, он приобретает номер, равный текущему значению счетчика, а счетчик увеличивается на 1. Теперь определим сравнение. Элемент больше

элемента, если приоритет элемента a меньше приоритета b , а если приоритеты равны, то больше тот элемент, номер входа которого больше.

```

1.     int Compare (a,b){
2.         if(a.pri < b.pri) return 1;
3.         if(a.pri = b.pri) {
4.             if(a.num < b.num) return 1;
5.             if(a.num > b.num) return -1;
6.             return 0;
7.         }
8.         return -1;
9.     }

```

Пояснение. Эта подпрограмма возвращает 1, если $a > b$, -1, если $a < b$ и 0, если $a = b$. Отметим, что последнее условие никогда не должно выполняться, поскольку все элементы имеют заведомо разные номера. Теперь можно построить двоичное дерево поиска, основанное на этом сравнении элементов. Если мы будем добавлять в это дерево элементы с использованием описанного сравнения, а снимать с дерева наименьший элемент, то это и будет очередь с приоритетами. Так же, как и в динамическом программировании, в каждой вершине будем записывать вес оптимального (или какого-либо) пути из начала в эту вершину. В начальный момент все вершины, кроме начальной имеют вес, равный ∞ . Стартовая вершина имеет вес, равный 0. Основной промежуточный шаг алгоритма поиска оптимального пути — это релаксация одной вершины относительно другой. Если вес вершины v больше суммы веса вершины u и веса ребра $e(uv)$, то вес вершины v приравнивается этой сумме, и в вершине v запоминается переход на вершину u .

```

1.     Relax(v,u){
2.         if(u.d+w(uv) < v.d){
3.             v.d= u.d+w(uv);
4.             v.pi=u;
5.         }
6.     }

```

Теперь опишем алгоритм Дейкстры поиска кратчайшего пути. Он очень прост и несколько напоминает алгоритм обхода графа в ширину (за исключением того, что используется очередь с приоритетами). В качестве приоритета используется текущий вес вершины. Алгоритм Дейкстры:

1. Все вершины не обработаны (белые, вес ∞), кроме последней (серая, стоит в очереди, вес 0).
2. Снимаем вершину из очереди, красим в черный цвет.
3. Для всех соседей этой вершины (в том числе и серых) делаем Relax, и, если они не серые, то помещаем в очередь.

4. Если очередь пуста, заканчиваем обработку, иначе переходим к п.2

Отметим, что в процессе работы алгоритма вершины, стоящие в очереди могут менять приоритеты. Поэтому после релаксации следует вершину переместить в правильное место очереди. На самом деле для реализации очереди с приоритетами используется другая, более удобная структура данных, но ее изучение выходит за рамки курса.

Псевдокод:

```

1. Deikstra(){
2.   PriQueue q;
3.   q.put(E); E.color=gray; E.d=0;
4.   do{
5.     v=q.get();
6.     for(each u in v.neighbours){
7.       if(u.color != black){
8.         Relax(u,v);
9.         if(u.color == gray) q.correct(u);
10.        if(u.color == white){
11.          q.put(u);
12.          u.color=gray;
13.        }
14.      }
15.    }
16.    v.color=black;
17.  }while(v!=B);
18. }
```

Пояснение. Этот код предполагает, что есть как-то организованная очередь с приоритетами. Эта структура данных имеет стандартные для очереди операции `put` и `get`. Кроме того есть операция перестройки очереди `correct` — она необходима, если элемент, уже стоящий в очереди изменил свой уровень приоритета. В начале все вершины белые. Кладем в очередь конец нашего пути (строка 3). Потом снимаем из очереди очередной элемент (строка 5) и релаксируем относительно него всех не черных соседей — черные уже приобрели свой конечный вес (строки 6-12). После релаксации не забываем поправить очередь (строка 9) и положить белых соседей в очередь (строки 11, 12). Процедура повторяется до тех пор, пока из очереди не снимем начальную вершину пути. Поскольку процедура `Relax` запоминает оптимальный переход, то восстановление пути проблемы не составит.

Корректность работы алгоритма Дейкстры

Теорема 1. По завершении обработки алгоритмом Дейкстры взвешенного ориентированного графа $G = V, E$ с неотрицательной весовой функцией w и истоком s для всех вершин $u \in V$ выполняется равенство $d(u) =$

Таблица 5.7: Пример работы алгоритма Дейкстры.

Граф	Очередь и веса вершин	Действия																					
	<table border="1"> <tr><td>E</td><td></td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>∞</td><td>∞</td><td>∞</td><td>∞</td><td>∞</td><td>∞</td><td>0</td></tr> </table>	E							B	a	b	c	d	e	E	∞	∞	∞	∞	∞	∞	0	В очереди стоит вершина E
E																							
B	a	b	c	d	e	E																	
∞	∞	∞	∞	∞	∞	0																	
	<table border="1"> <tr><td>e(1)</td><td>d(8)</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>∞</td><td>∞</td><td>∞</td><td>∞</td><td>8</td><td>1</td><td>0</td></tr> </table>	e(1)	d(8)						B	a	b	c	d	e	E	∞	∞	∞	∞	8	1	0	Снимаем E из очереди и релаксируем всех соседей и помещаем в очередь (если они не там)
e(1)	d(8)																						
B	a	b	c	d	e	E																	
∞	∞	∞	∞	8	1	0																	
	<table border="1"> <tr><td>c(3)</td><td>d(5)</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>∞</td><td>∞</td><td>∞</td><td>3</td><td>5</td><td>1</td><td>0</td></tr> </table>	c(3)	d(5)						B	a	b	c	d	e	E	∞	∞	∞	3	5	1	0	Снимаем e из очереди и релаксируем соседей относительно нее
c(3)	d(5)																						
B	a	b	c	d	e	E																	
∞	∞	∞	3	5	1	0																	
	<table border="1"> <tr><td>d(4)</td><td>a(5)</td><td>b(5)</td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>∞</td><td>5</td><td>5</td><td>3</td><td>4</td><td>1</td><td>0</td></tr> </table>	d(4)	a(5)	b(5)					B	a	b	c	d	e	E	∞	5	5	3	4	1	0	Снимем из очереди c.
d(4)	a(5)	b(5)																					
B	a	b	c	d	e	E																	
∞	5	5	3	4	1	0																	
	<table border="1"> <tr><td>a(5)</td><td>b(5)</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>∞</td><td>5</td><td>5</td><td>3</td><td>4</td><td>1</td><td>0</td></tr> </table>	a(5)	b(5)						B	a	b	c	d	e	E	∞	5	5	3	4	1	0	Снимем из очереди d. Ни одна вершина не релаксировала
a(5)	b(5)																						
B	a	b	c	d	e	E																	
∞	5	5	3	4	1	0																	
	<table border="1"> <tr><td>b(5)</td><td>B(13)</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>13</td><td>5</td><td>5</td><td>3</td><td>4</td><td>1</td><td>0</td></tr> </table>	b(5)	B(13)						B	a	b	c	d	e	E	13	5	5	3	4	1	0	Снимаем вершину a. Вершина B приобрела вес 13.
b(5)	B(13)																						
B	a	b	c	d	e	E																	
13	5	5	3	4	1	0																	
	<table border="1"> <tr><td>B(6)</td><td></td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>6</td><td>5</td><td>5</td><td>3</td><td>4</td><td>1</td><td>0</td></tr> </table>	B(6)							B	a	b	c	d	e	E	6	5	5	3	4	1	0	Снимаем вершину b. Вершина B релаксировала до 6
B(6)																							
B	a	b	c	d	e	E																	
6	5	5	3	4	1	0																	
	<table border="1"> <tr><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td>a</td><td>b</td><td>c</td><td>d</td><td>e</td><td>E</td></tr> <tr><td>6</td><td>5</td><td>5</td><td>3</td><td>4</td><td>1</td><td>0</td></tr> </table>								B	a	b	c	d	e	E	6	5	5	3	4	1	0	Снимаем вершину B. Вершина B обработана, оптимальный путь найден!
B	a	b	c	d	e	E																	
6	5	5	3	4	1	0																	

$\delta(s, u)$, где $d(u)$ — вес, подсчитанный в соответствии с алгоритмом Дейкстры, $d(u) = \delta(s, u)$ — вес кратчайшего пути.

Доказательство. Для доказательства этого утверждения используют следующие соображения (см. рис. 5.13). Есть множество обработанных вершин (черные). Есть множество “пограничных” вершин (серые), которые имеют

ребра к черным вершинам, есть множество белых вершин, которые не обрабатывались.

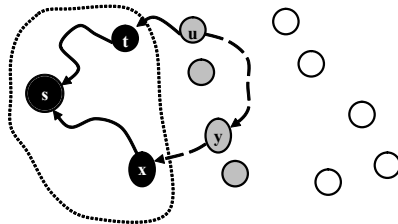


Рис. 5.13: К доказательству правильности алгоритма Дейкстры.

1. Никакое ребро в черную вершину не лучше отмеченных ребер, (поскольку серые — релаксированы).
2. Выберем вершину u , которая стоит первой в очереди. Она имеет минимальное значение d среди всех серых вершин (по принципу построения очереди). Допустим есть путь из u в s , лучше, чем путь $u \rightarrow t \rightarrow s$. Тогда этот путь обязательно пройдет через серую вершину y : $u \rightarrow y \rightarrow x \rightarrow s$, поскольку путь в конце концов войдет в черную область, а соседями черной области — серые. Но:

$$w(u \rightarrow y \rightarrow x \rightarrow s) = w(u \rightarrow y) + w(x \rightarrow s) \text{ — вес пути есть сумма весов подпутьей}$$

$$w(y \rightarrow x \rightarrow s) = d(y) > d(u) = w(u \rightarrow t \rightarrow s) \text{ — вес } u \text{ минимален по построению очереди с приоритетами.}$$

$$w(u \rightarrow y) > 0 \text{ — веса положительны!}$$

$$w(u \rightarrow y \rightarrow x \rightarrow s) > w(u \rightarrow t \rightarrow s) \text{ — противоречие!}$$

□

Заметим, что поскольку допустимы циклы, то алгоритм Дейкстры можно применять и для неориентированных графов — неориентированный граф можно трактовать как ориентированный, где каждое неориентированное ребро (vi) эквивалентно двум ориентированным ребрам (vi) и (iv) . Кстати пример, приведенный в таблице 5.6.5 является неориентированным графом.

Упражнение 34. Приведите пример графа с некоторыми отрицательными ребрами, когда алгоритм Дейкстры не даст оптимального решения.

Время работы алгоритма Дейкстры

Время работы алгоритма Дейкстры на графе немного больше, чем время работы алгоритма динамического программирования. Каждая вершина

проходит одну релаксацию относительно всех соседей, поэтому надо потратить время $O(|E|)$. Но организация очереди стоит времени, причем перестройка очереди происходит при каждой релаксации (в худшем случае на каждом ребре). В худшем случае алгоритм работает за $T = O(|E| \cdot \log(|V|))$. Логарифм в выражении появляется потому, что мы пользуемся бинарным деревом для реализации очереди с приоритетами

5.6.6 Алгоритм Беллмана-Форда

Описанные ранее алгоритмы поиска оптимального пути имели ограничения — либо ациклический граф, либо только неотрицательные веса. Для случая, когда веса как положительны, так и отрицательны, применим следующий алгоритм.

Алгоритм Беллмана-Форда основывается на повторе релаксации всех ребер (в произвольном порядке) $|V| - 1$ раз. Вообще говоря, не очевидно, что такой алгоритм находит оптимальный путь. Докажем это.

Теорема 2. *Алгоритм Беллмана-Форда находит оптимальный путь в графе.*

Доказательство. Если есть оптимальный путь $\{v_1, v_2, \dots, v_n\}$, то на i -шаге будет определена длина пути от i -вершины до конца. Доказательство по индукции. Инициация тривиальна — в начале мы имеем пустой список (тривиальный оптимальный путь из end в end нам известен заранее, это 0). Пусть на i -шаге $d(v_i)$ равна длине кратчайшего пути до вершины E . Тогда на следующем шаге произойдет релаксация ребра (v_i, v_{i+1}) и таким образом определится вес оптимального пути от v_{i+1} до E . После первой релаксации первая вершина на оптимальном пути приобретет правильный вес, потом вторая вершина . . . В конце концов все вершины оптимального пути приобретут нужный вес. Самый длинный путь, который можно представить себе, состоит из $|V|$ вершин, поэтому в худшем случае нам необходимо пройти $|V| - 1$ шагов. \square

Заметим, что алгоритм Беллмана-Форда позволяет очень просто определить, существует ли в графе отрицательный цикл, достижимый из вершины end. Достаточно произвести релаксацию ровно $|V|$ раз. Если при исполнении последней релаксации длина кратчайшего пути от какой-либо вершины строго уменьшилась, то в графе есть отрицательный цикл, достижимый из E . Отметим, что наличие отрицательного цикла делает задачу вырожденной — мы можем сколько угодно раз крутиться по этому циклу и уменьшать длину пути до $-\infty$.

Таблица 5.8: Сводка по алгоритмам поиска оптимального пути.

Алгоритм	Область применимости	Время работы
Динамическое программирование	Минимум и Максимум. Ориентированный ациклический граф. На веса нет ограничений	$O(V + E)$
Дейкстра	Минимум. Ориентированный и неориентированный граф, веса строго положительны	$O(E \cdot \log V)$
Беллман-Форд	Минимум. Ориентированный граф, веса любые, нет циклов отрицательного веса.	$O(V \times E)$

5.6.7 Сводка основных алгоритмов, ищущих оптимальный путь на графе

В таблице 5.6.7 приведена сводка по рассмотренным алгоритмам поиска оптимального пути в графе.

5.7 Сети и потоки в них

5.7.1 Основные понятия

Представим себя владельцем нефтеперерабатывающего завода. Пусть есть нефтяная вышка S , от которой нефть течет к нашему заводу по некому множеству труб (рис.5.7.1). Каждая труба обладает пропускающей способностью. Кроме труб есть узлы, в которых трубы сливаются, раздваиваются ... Допустим, по какой-то причине поменять трубы на трубы с большей пропускающей способностью мы не можем (запрещает экологическая инспекция). Тогда использовать ресурсы нашего источника максимально мы можем, только искусно распределив поток нефти по имеющимся трубам. Перед нами стоит задача: распределить поток нефти по трубам с известной

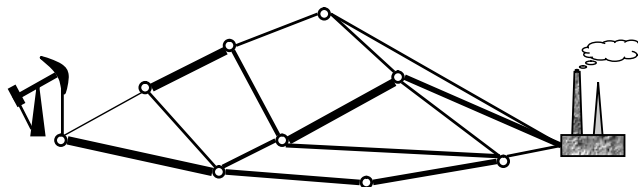


Рис. 5.14: Сеть труб

пропускной способностью так, чтобы была максимальная прибыль. Задача, на самом деле далеко не простая. Если пустить слишком мало нефти по трубе — прибыль понизится, если слишком много — труба лопнет и прибыль

понижится тоже. Аналогичные задачи решаются при передаче информации через Интернет. Нужно оптимально распределить пакеты информации так, чтобы они максимально быстро и без потерь пересылались через существующие каналы, которые имеют разную пропускную способность. Условия перечисленных задач очень сходны: есть сеть, есть источник (данных, веществ и др.), есть место, куда эти данные и др. нужно доставить. Необходимо найти наилучший поток в этой сети. Дадим формальное определение сети и потока:

Определение 3. Сеть — это

1. Ориентированный граф $G = \{V, E\}$;
2. На каждом ребре (uv) определен строго положительный вес $c(uv) > 0$ (пропускная способность).
3. На каждом несуществующем ребре (на каждой паре вершин, где нет ребра), определен нулевой вес
4. Выделены две вершины s — источник, t — сток.
5. Для любой вершины v существует путь от s через v к t : $s \rightarrow v \rightarrow t$.

Условия 2, 3 означают, что для любой пары вершин определено неотрицательное число $c : V \times V \rightarrow R^+$. Это число называется *пропускной способностью ребра*. Если ребра нет, то пропускная способность равна 0.

Определение 4. Поток в сети называется функцией, определенной на паре вершин $f : V \times V \rightarrow R$ и обладающая свойствами:

1. $f(uv) \leq c(uv)$ для любого ребра (в том числе и для несуществующего ребра) поток ограничен весом ребра — “пропускной способностью”.
2. $f(uv) = -f(vu)$. Это свойство называется *кососимметричностью*.
3. для любой вершины $v \in V \setminus \{s, t\}$ $\sum_{u \in V} f(uv) = 0$, это свойство можно назвать *законом сохранения* — сколько в вершину “втекает”, столько и “вытекает”. Обратим внимание на то, что закон сохранения выполняется для всех вершин, кроме источника и стока;
4. $|f| = f(s) = \sum_{v \in V} f(s, v)$ — суммарный поток (величина потока) равен потоку в источнике и сумме потоков, текущих из источника в каждую вершину.

Упражнение 35. Из этих свойств потока следует, что $f(t) = -f(s)$, то есть сколько вытекло из источника, столько втекло в сток. Докажите это утверждение.

Задача о максимальном потоке. Дана сеть $G = \{V, E, c\}$ (см. определение). Найти максимальный поток, т.е. $|f| = f(s) = \sum_{v \in V} f(s, v) \rightarrow \max$. Попробуем найти хоть один поток. Для этого найдем на нашем ориентированном графе путь P (любой) из вершины s в вершину t , проходящий по настоящим ребрам (там, где пропускная способность строго больше нуля). Это можно сделать с помощью поиска в ширину или в глубину. На этом пути можно найти ребро с минимальной пропускной способностью $c_{\min} = \min_{e \in P} c(e)$. Теперь определим поток:

$$f(uv) = \begin{cases} c_{\min} & uv \in P \\ -c_{\min} & vu \in P \\ 0 & uv \notin P \end{cases} \quad \text{где } u \text{ и } v \text{ — смежные вершины}$$

Теперь, наверное, можно поток улучшить. Метод поиска максимального потока, основанный на последовательном увеличении потока (если оно возможно) был придуман Фордом и Фалкерсоном и носит их имя.

5.7.2 Метод Форда-Фалкерсона

Пусть у нас дан поток f . Определим остаточную пропускную способность следующим образом $c_f(uv) = c(uv) - f(uv)$. Поскольку по определению потока он не может превышать пропускную способность, остаточная пропускная способность неотрицательна. Кроме того остаточная пропускная способность может превратить некоторые виртуальные ребра (для которых пропускная способность была равна 0), в реальные ребра с положительной пропускной способностью. Например, если есть ребро (uv) , нет ребра (vu) и $f(uv) = f_{uv} > 0$. Тогда по определению потока $f(vu) = -f(uv)$. Стало быть, остаточная пропускная способность на ребре (vu) определится как $c_f(vu) = c(vu) - f(vu) = 0 - (-f_{uv}) = f_{uv} > 0$, т.е. станет положительной. Смысл возникновения таких несуществующих ребер заключается в следующем. Допустим, мы пропустили по ребру слишком большой поток, а в оптимальном потоке эта величина должна быть меньше. Тогда такие дополнительные ребра позволяют нам исправить ошибки, пропуская поток в обратном направлении, что эквивалентно сокращению потока в прямом направлении.

Определение 5. *Остаточной сетью называется сеть, в которой ребра есть между теми вершинами, где остаточная пропускная способность положительна $c_f(uv) > 0$:*

$$c_f = \{V, E_f = \{(uv) \in V \times V : c_f(uv) > 0\}, c_f\}$$

В основе метода Форда-Фалкерсона лежит следующая

Лемма 1. *Лемма. Пусть G сеть, f — поток G' — остаточная сеть и f' — поток в остаточной сети. Тогда $f + f'$ является потоком в сети G и его величина равна $|f + f'| = |f| + |f'|$.*

Доказательство. Для доказательства надо проверить, что для $f + f'$ выполняются свойства потока:

- Кососимметричность. Поскольку $f(vu) = -f(uv)$ и $f'(vu) = -f'(uv)$, то $f(vu) + f'(vu) = -f(uv) + f'(uv)$
- Закон сохранения. Поскольку $\sum_v f(uv) = 0$ и $\sum_v f'(uv) = 0$, то $\sum_v (f(uv) + f'(uv)) = 0$
- Ограничение на пропускную способность. Действительно $f'(uv) < c_f(uv)$. Поскольку по определению $c_f(vu) = c(vu) - f(vu)$, то $f'(uv) \leq c(vu) - f(vu)$, или $f(uv) + f'(uv) \leq c(uv)$
- Величина потока $|f + f'| = \sum_v (f(sv) + f'(sv)) = |f| + |f'|$

□

Если остаточная сеть имеет путь из источника s в сток t , то по сети можно пропустить дополнительный поток так же, как мы это делали в первый раз. Процедуру можно повторять до тех пор, пока в остаточной сети не останется путей из источника в сток. На каждом шаге такой процедуры мы увеличиваем поток. Всего возможных способов увеличения потока конечно. Поэтому процедура сходится за конечное число шагов. Правда число шагов может оказаться очень большим. В таблице 5.7.2 показан пример работы алгоритма. Вообще говоря, из того, что в дополнительной сети нет пути из источника в сток, не следует, что мы нашли максимальный поток. Может быть, если бы мы выбирали другие пути из s в t для пропуска дополнительного потока, мы смогли достичь в конце концов большего потока? Есть теорема, доказывающая, что этот алгоритм дает в конце максимальный поток. Она будет доказана позже.

5.7.3 Разрезы в сетях

Определение 6. *Разрезом сети $G = \{V, E, c\}$ называется разбиение множества вершин на два подмножества S и T , таких, что:*

1. $V = S \cup T$
2. $S \cap T = \emptyset$
3. $s \in S, \quad t \in T$

Говорят, что ребро $e = (vu)$ принадлежит разрезу, если вершина u принадлежит подмножеству S , а вершина v — подмножеству T . Просто говоря, ребро принадлежит разрезу, если мы его “разрезали”, разделив вершины между S и T . Принадлежность ребра e к разрезу обозначается $\in E_c$. Отметим, что ребра, идущие из T в S разрезу не принадлежат (рис. 5.15).

Таблица 5.9: Поиск максимального потока (числа в скобках — величина потока).

сеть	поток	пояснения
	1	Находим путь из источника в сток (выделено)
	1	Строим дополнительную сеть, ребро et исчезает.
	2	Находим путь по дополнительной сети
	2	исчезает ребро cd
	3	Находим путь
	3	Строим остаточную сеть исчезло ребро sb
	6	Поток увеличился на 3
		В дополнительной сети нет пути из s в t. Поток нельзя увеличить.
		Результирующий поток

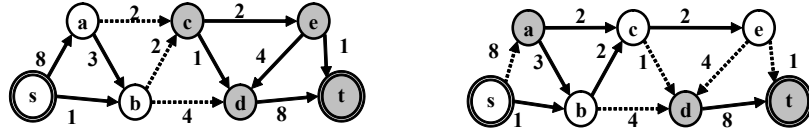


Рис. 5.15: Примеры разрезв в сетях. Белым показаны вершины, принадлежащие множеству S , серым — множество T . Пунктиром показаны ребра, принадлежащие разрезу.

Определение 7. *Пропускной способностью разреза называется сумма пропускных способностей всех ребер разреза: $C(S, T) = \sum_{u \in S, v \in T} c(uv)$. Разрез минимальной пропускной способности называют минимальным разрезом сети. Поток из S в T равен: $f(S, T) = \sum_{u \in S, v \in T} f(uv)$. Суммирование ведется по всем вершинам u , принадлежащим S и по всем вершинам v , принадлежащим T . Максимальный поток через разрез равен пропускной способности разреза.*

Свойства потока между множествами вершин сети $G = \{V, E, c\}$: 1. $f(X, X) = 0$; 2. $f(X, Y) = -f(Y, X)$; 3. Пусть $X, Y, Z \subseteq V$ и $X \cap Y = \emptyset$, тогда $f(X \cup Y, Z) = f(X, Z) + f(Y, Z)$; $f(Z, X \cup Y) = f(Z, X) + f(Z, Y)$, т.е. поток дистрибутивен. 4. Пусть $X \subseteq V \setminus \{s, t\}$ (в подмножество X могут входить любые вершины сети, кроме источника и стока). Тогда $f(X, V) = 0$. Поток из любого множества вершин, не включающего s и t , в множество вершин, равен нулю.

Упражнение 36. *Докажите свойства 1-4.*

Лемма 2. *Имеется некая сеть $G = \{V, E, c\}$, некий ее разрез $S \cup T$ и некий поток через эту сеть f . Тогда для любого разреза выполняется $f(S, T) = |f|$, т.е. поток через любой разрез равен потоку между источником и стоком.*

Доказательство. Очевидно, что любой поток обязательно проходит через ребра разреза (поскольку любая вершина из V принадлежит либо S , либо T). Более строгое доказательство выглядит так:

$$f(S, T) = f(S, S) + f(S, T) = f(S, S \cup t) = f(S, V) = f(S \setminus s) + f(s, V) = |f|$$

Здесь использовано свойство 4: множество $S \setminus s$ не содержит ни источника, ни стока, поэтому $f(S \setminus s) = 0$. \square

Лемма 3. *для любого разреза $S \cup T = V$ и любого потока f выполнено $|f| \leq c(ST)$.*

Доказательство. следует непосредственно из первого свойства потока и леммы 1. \square

5.7.4 Теорема о максимальном потоке и минимальном разрезе

Теорема 3 (Фалкерсона). *Следующие утверждения эквивалентны.*

1. Поток f — максимальный
2. Остаточная сеть не содержит путей из s в t
3. Существует разрез $S \cup T = V$, такой, что $|f| = c(S, T)$, в частности, это значит, что если вершины $u \in S$, $v \in T$ то $f(uv) = c(uv)$, то есть существует такой разрез, поток через который равен пропускной способности разреза. Иными словами, существует хотя бы один разрез сети, поток через который максимален и равен пропускной способности разреза.

Доказательство. **1** \rightarrow **2** очевидно. От противного: если бы остались пути из s в t , поток можно было бы еще увеличить.

2 \rightarrow **3** Пусть в остаточной сети нет путей из s в t . Рассмотрим в качестве S' множество вершин, достижимых из источника s по остаточной сети, в качестве T' — остальные вершины. Мы получили разрез, поскольку $S' \cap T' = \emptyset$, $S' \cup T' = V$, $s \in S'$, $t \in T'$. По построению ни одно ребро остаточной сети не принадлежит разрезу (мы собрали в множестве все ребра, достижимые по остаточной сети из источника). Поэтому остаточный вес ребер разреза равен нулю. Формула для получения остаточного веса ребра (uv) $c_f(uv) = c(uv) - f(uv)$, следовательно, для ребер, принадлежащих разрезу $c(uv) = f(uv)$, и поток через разрез равен пропускной способности разреза, что и требовалось доказать.

3 \rightarrow **1** Для любого разреза $S \cup T = V$ выполнено $|f| \leq c(S, T)$, поэтому из равенства $|f| = c(S, T)$ следует максимальность потока. □

5.7.5 Случай многих источников и стоков

Разумеется, Вы владеете не одной качалкой для нефти и заводов у Вас несколько. Как найти оптимальное распределение нефти по трубопроводам? Очень просто. Надо ввести дополнительные вершины — общий источник и общий сток и провести ребра из универсального источника в Ваши качалки и от Ваших заводов к универсальному стоку. Если Ваше месторождение имеет ограничение по дебету (мощности), то эти дополнительные ребра должны иметь пропускную способность, равную дебету скважин. Аналогично для ограничений мощности заводов (рис.5.16).

5.7.6 Паросочетания в двудольном графе

Двудольный граф позволяет разбить множество вершин на два подмножества, так, что ребер внутри этих подмножеств нет. Паросочетания — мно-

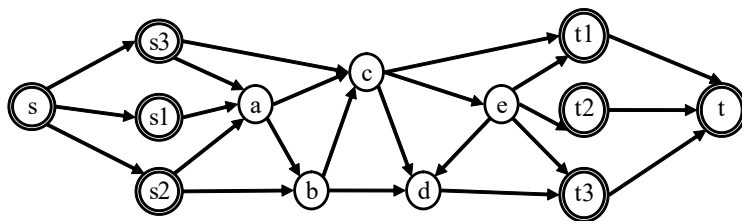


Рис. 5.16: Случай многих источников и стоков.

жество ребер в графе, такое, что нет ребер, имеющих общую вершину. Добавляем к двудольному графу источник и сток, пропускные способности устанавливаем на всех ребрах графа $=1$. Ищем максимальный целочисленный поток. Ребра, которые будут иметь поток, равный 1 обеспечат максимальное паросочетание. Таким способом может решиться задача о браке: есть n юношей и n девушек, нужно составить максимальное число супружеских пар. Пусть юноши составляют одну “долю” графа, девушки другую (все они являются вершинами графа). Между двумя вершинами есть ребро, если между соответствующими юношей и девушкой есть симпатия. Найдя максимальное паросочетание в таком графе, мы решим задачу о максимальном количестве браков. В биоинформатике поиск потока на двудольном графе — сети с вершинами — последовательностями можно использовать для поиска ортологов. Вершины графа одной доли — гены в одном геноме, другая доля графа соответствует второму геному. Задача о паросочетаниях, а также другие задачи, связанные с потоками встречаются в качестве вспомогательных задач весьма часто, в том числе и в биоинформатике.

5.8 Понятие о линейном программировании

Есть набор из некоторого (обычно большого) числа переменных x_1, \dots, x_n . Эти переменные неотрицательны $x_i \geq 0$. На переменные еще есть некоторое количество линейных ограничений $\sum_i d_i^k x_i \leq f_k$. Надо найти максимум некоторой линейной формы $\sum a_i x_i \rightarrow \max$. Эти задачи называются задачами линейного программирования. На самом деле это никакого отношения к программированию не имеет — просто использован неудачный перевод. Правильнее было бы называть этот класс задач задачами линейного планирования. Итак,

- Есть набор переменных $\{x_1 \dots x_n\}$
- Переменные неотрицательны $x_i \geq 0$
- Есть ограничения типа неравенства

$$\sum_j d_{ij} x_j \geq f_i$$

- Есть ограничения типа равенства

$$\sum_j b_{ij}x_j = c_i$$

- Надо максимизировать линейную форму

$$\sum_i a_i x_i \rightarrow \max$$

Примеры таких задач:

- Задача о максимальном потоке:
 - переменные $x_i = f(uv)$ — значения потока
 - *линейные* ограничения $\sum_v f(vu) = 0$
 - *линейные* ограничения типа неравенства $f(uv) \leq c(uv)$
 - максимизируем $\sum_v f(sv)$
- Транспортная задача — оптимизировать транспортные потоки
- Анализ метаболизма

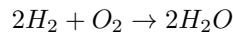
Транспортная задача Имеется некий однородный груз, который нужно перевезти с n складов на m заводов. Для каждого склада i известно, сколько в нчм находится груза a_i , а для каждого завода известна его потребность b_j в грузе. Стоимость перевозки пропорциональна расстоянию от склада до завода (все расстояния c_{ij} от i -го склада до j -го завода известны). Требуется составить наиболее дешевый план перевозки. Решающими переменными в данном случае являются x_{ij} — количества груза, перевезенного из i -го склада на j -й завод. Они удовлетворяют ограничениям:

$$\begin{aligned} x_{i1} + x_{i2} + \dots + x_{im} &\leq a_i, \\ x_{1j} + x_{2j} + \dots + x_{nj} &\geq b_j. \end{aligned}$$

Целевая функция имеет вид: $c_{11}x_{11} + c_{12}x_{12} + \dots + c_{nm}x_{nm}$, которую надо минимизировать.

Анализ метаболизма Часто возникают задачи изучения метаболизма. Традиционный подход основан на изучении динамики изменения концентраций метаболитов с описанием ферментативных реакций. Этот подход имеет вполне достойное применение, например, при моделировании кинетики действия и распада лекарств, при изучении биосинтеза различных веществ и много, где еще. Для описания динамики пишут большие и сложные дифференциальные уравнения, которые потом решают. Однако, при составлении этих уравнений надо использовать кинетические константы, которые, в большинстве своем, не известны, или известны для других, но

похожих ферментов, в условиях, часто отличных от физиологических. При решении этих уравнений надо задавать начальные условия, которые не известны. Так что такой подход, несмотря на привлекательность, имеет целый ряд серьезных проблем. Но есть другой подход к анализу метаболизма, который использует гораздо меньше информации. Более того эта информация в большинстве случаев известна достоверно. Это стехиометрия химических реакций. Например, мы знаем, что из 2 молекул водорода и одной молекулы кислорода получается две молекулы воды. Можно написать:



Можно ввести понятие потока – количество элементарных реакций, происходящих в единицу времени. Или, количество того или иного вещества убывающего или прибывающего в единицу времени в результате реакции. Стехиометрия связывает эти потоки. Например, можно написать:

$$F(H_2O) = -F(O_2), \quad F(H_2O) = -1/2F(H_2)$$

Можно обзреть всю метаболическую карту какого-нибудь организма и написать все стехиометрические ограничения, которые следуют из этой карты. Заметим, что потоки неотрицательны. Между ними есть линейные ограничения типа равенства. Можно поставить задачу о максимизации чего-нибудь. Часто в качестве функционала используют биомассу, при этом кроме стехиометрии используют соотношение основных метаболитов из реальных данных, например аминокислотный состав, липидный состав и пр. В результате получается большая задача линейного программирования, которую можно решить и оценить, насколько активны те или иные реакции. Другое применение этого подхода – оптимизация продуцентов целевых продуктов. Например, Вы разрабатываете штамм-продуцент какого-то витамина. Из анализа метаболической карты и решения задачи линейного программирования Вы можете понять сколько грамм целевого продукта в принципе можно получить из грамма глюкозы.

5.8.1 Геометрическая интерпретация

Итак, у нас есть линейные ограничения. Ограничения типа равенства можно заменить парой ограничений типа неравенства:

$$\sum_j b_{ij}x_j = c_i \Rightarrow \sum_j b_{ij}x_j \leq c_i ; \sum_j b_{ij}x_j \geq c_i;$$

Каждое линейное неравенство определяет полупространство (рис.5.17). Совокупность всех ограничений выделяет выпуклый многогранник. Линейная форма, которую оптимизируем также является гиперплоскостью. Значение функционала определяет положение этой гиперплоскости.

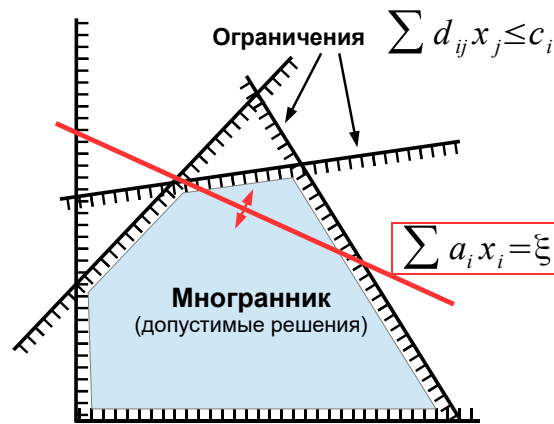


Рис. 5.17: Графическое представление задачи линейного программирования

5.8.2 Немного истории

В 1938 году Леонид Витальевич Канторович в порядке научной консультации приступил к изучению чисто практической задачи по составлению наилучшего плана загрузки лущильных станков (фанерный трест). Эта задача не поддавалась обычным методам. Стало ясно, что задача не случайная. В 1939 году Леонид Канторович опубликовал работу «Математические методы организации и планирования производства», в которой сформулировал новый класс экстремальных задач с ограничениями и разработал эффективный метод их решения, таким образом были заложены основы линейного программирования. Изучение подобных задач привело к созданию новой научной дисциплины линейного программирования и открыло новый этап в развитии экономико-математических методов. В 1949 году американский математик Джордж Бернард Данциг разработал эффективный метод решения задач линейного программирования (ЗЛП) — симплекс-метод. В 1975 году Леонид Витальевич Канторович стал лауреатом Нобелевской премии по экономике (совместно с Тьяллингом Купмансом за вклад в теорию оптимального распределения ресурсов). Это единственный отечественный исследователь, получивший Нобелевскую премию по экономике. Термин «программирование» нужно понимать в смысле «планирования» (один из переводов англ. programming). Он был предложен в середине 1940-х годов Джорджем Данцигом, одним из основателей линейного программирования, еще до того, как компьютеры были использованы для решения линейных задач оптимизации.

5.8.3 Симплекс метод

Заметим, что каждое из линейных неравенств на переменные ограничивает полупространство в соответствующем линейном пространстве. В результате

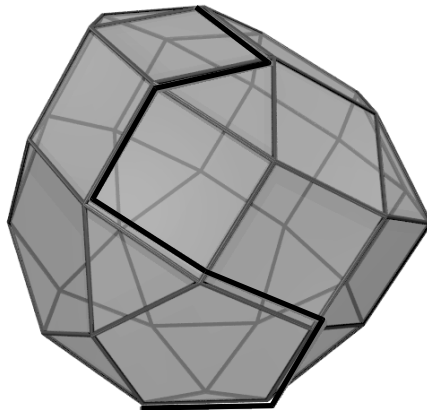


Рис. 5.18: Поиск оптимальной вершины симплекс-методом

все неравенства ограничивают некоторый многогранник (возможно, бесконечный). Уравнение $W(x) = c$, где $W(x) = \sum_i a_i x_i$ — максимизируемый (или минимизируемый) линейный функционал, порождает гиперплоскость $L(c)$. Зависимость от c порождает семейство параллельных гиперплоскостей. Тогда экстремальная задача приобретает следующую формулировку: требуется найти такое наибольшее c , что гиперплоскость $L(c)$ пересекает многогранник хотя бы в одной точке. Заметим, что пересечение оптимальной гиперплоскости и многогранника будет содержать хотя бы одну вершину, причём, их будет более одной, если пересечение содержит ребро или k -мерную грань. Поэтому максимум функционала можно искать в вершинах многогранника. Принцип симплекс-метода состоит в том, что выбирается одна из вершин многогранника, после чего начинается движение по его рёбрам от вершины к вершине в сторону увеличения значения функционала (рис.5.18). Когда переход по ребру из текущей вершины в другую вершину с более высоким значением функционала невозможен, считается, что оптимальное значение c найдено. Последовательность вычислений симплекс-методом можно разделить на две основные фазы:

- нахождение исходной вершины множества допустимых решений,
- последовательный переход от одной вершины к другой, ведущий к увеличению (уменьшению) значения целевой функции.

При этом в некоторых случаях исходное решение очевидно или его определение не требует сложных вычислений, например, когда все ограничения представлены неравенствами вида $x_i \leq b_i$ (тогда нулевой вектор совершенно точно является допустимым решением, хотя и, скорее всего, далеко не самым оптимальным).

Глава 6

Свойства задач и NP-полные задачи.

6.1 Нижняя оценка времени сортировки массива

До сих пор мы рассматривали разные задачи и для них строили алгоритмы, которые их решают. При этом мы делали оценку времени работы алгоритма. Например, когда мы рассматривали задачу сортировки массива, мы сначала построили наивный алгоритм, который сортировал массив за время $T = O(N^2)$. Потом мы построили алгоритм, который выполнял сортировку за время $T = O(N \log N)$. Может быть если еще подумать, то можно построить алгоритм сортировки за время, скажем, $T = O(N)$, или даже $T = O(\log N)$?

Итак, задача сортировки массива. Сортировка массива предполагает, что элементы массива Θ можно сравнивать, т.е. для любых двух элементов массива можно сказать верно одно из следующих утверждений $\forall a, b \in \Theta$ $a > b$ или $b > a$ или $a = b$. При этом выполняется условие транзитивности: если $a > b$ и $b > c$, то $a > c$. Аналогично для равенства. Задача: переставить элементы массива так, чтобы $\forall 0 < i < N$ выполнялось $a[i] > a[i - 1]$.

Теорема 1. *Не существует алгоритма, который решал бы задачу сортировки сравнениями массива быстрее, чем за $O(N \log N)$, иначе говоря, быстрее, чем за $c \cdot N \log N$, где c — некая константа.*

Доказательство. Любому алгоритму сортировки соответствует дерево решений: надо сделать сравнение и в зависимости от результата что-то сделать (например, переставить два элемента) и затем сделать следующее сравнение. В результате алгоритму соответствует дерево решений. Листья этого дерева соответствуют отсортированному массиву. Если зафиксировать длину массива, то есть взаимно-однозначное соответствие дерева решений

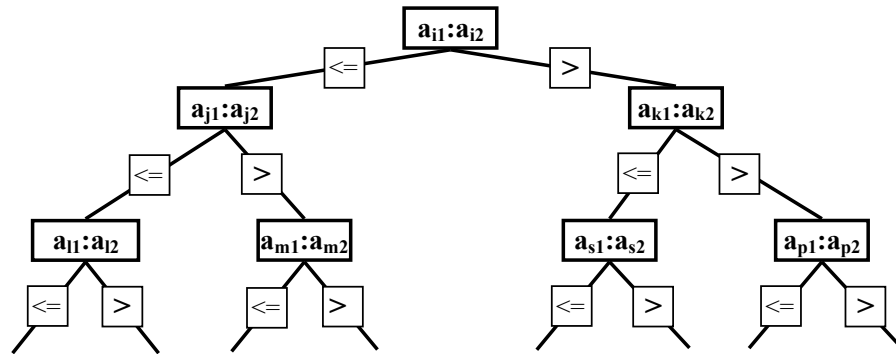


Рис. 6.1: Дерево решений в задаче сортировки. Узлы — сравнения. В зависимости от результата сравнения происходят какие-то действия, определяемые алгоритмом и происходят следующие сравнения.

и алгоритма (если не учитывать такие “мелочи”, как действия, связанные с перестановкой элементов). Количество листьев дерева решений равно числу возможных перестановок n элементов. Действительно, каждому отсортированному массиву соответствует набор индексов исходного массива в том порядке, в каком элементы будут лежать в отсортированном массиве.

Время работы алгоритма сортировки равно высоте h дерева решений. Максимальное число листьев в дереве высоты h равно 2^h . Число листьев в дереве решений равно $n!$, поскольку каждому листу соответствует одна перестановка элементов массива. Поэтому для любого алгоритма должно выполняться

$$N! \leq 2^h \text{ или } h \geq \log_2 N$$

Применяя оценку для факториала $N! \geq (N/e)^N$, получаем окончательно:

$$T(N) = c \cdot h \geq c \cdot N \log N - N \cdot \log e = O(N \cdot \log N)$$

□

Таким образом, $O(N \cdot \log N)$ — нижняя оценка времени работы любого алгоритма сортировки массива. Отметим, что эта оценка является оценкой для наихудшего случая. Не исключено, что дерево решений может иметь достаточно короткие пути, и для лучшего случая время работы алгоритма будет меньше. На самом деле есть алгоритмы сортировки, которые это делают за линейное время (!). Это алгоритмы целочисленной сортировки, которые учитывают не только результат сравнения $a > b$, но и то насколько один элемент больше другого. Это позволяет сделать ветвление в дереве решений не на две ветки, а на количество ветвей, сравнимое с длиной массива. Для целых чисел такое сравнение возможно, но, например для строк, это сделать затруднительно, а для многих других случаев просто невозможно. Из этой теоремы следует важный вывод. Оказывается можно до

создания алгоритма оценить минимально достижимое время работы, и это время является *свойством задачи*, а не алгоритма, которого мы не знаем.

6.2 Формальные языки. Допускающий и распознающий алгоритмы

Задачи бывают разные. В том числе, бывают и те, в которых ответ надо дать в виде “да” или “нет”. Для решения задач такого плана (а их еще иногда называют задачами разрешения) удобно использовать терминологию теории формальных языков. Понятие “формальный язык” определяет множество конечных слов (строк) над конечным алфавитом. Вход любого алгоритма сводится к строке (любая программа имеет на входе двоичный файл). Поэтому время работы алгоритма, вообще говоря, правильно оценивать как функцию длины входной строки. Итак:

Алфавит — некое (конечное) множество Σ .

Слово длины n — последовательность из n символов алфавита Σ .

Σ^* — Множество всех слов всех длин

Язык L над алфавитом Σ — подмножество Σ^* : $L \subseteq \Sigma^*$ — (некоторое множество строк символов из алфавита).

Поскольку алфавит конечен, то есть отображение множества всех слов в алфавите $\Sigma = \{a_i\}$ на множество всех слов в алфавите $\Sigma_{0,1} = \{0, 1\}$. Далее будем рассматривать языки над алфавитом $\{0, 1\}$. Так, например, можно рассматривать алфавит $\Sigma_{0,1} = \{0, 1\}$ и язык $L = \{10, 11, 101, 111, 1011, 1101, 10001 \dots\}$, состоящий из двоичных записей простых чисел.

Мы будем рассматривать класс алгоритмов, результатом которых является ответ либо “да” либо “нет”. Т.е. алгоритмы, которые реализует отображение $\{0, 1\}^* \rightarrow \{0, 1\}$. Эти алгоритмы называются решающими алгоритмами. Таким образом, алгоритм A определяет подмножество $L \subseteq \Sigma^*$, такое, что $A(x \in L) = 1, A(x \notin L) = 0$. Иными словами, алгоритм определяет некоторый язык. Теперь становится интуитивно ясным, что есть соответствие между формальными языками и задачами (ведь алгоритм решает задачу). Разумеется, класс решающих алгоритмов гораздо уже, чем все мыслимые алгоритмы. В частности, алгоритм сортировки затруднительно представить в виде решающего алгоритма. Однако, алгоритмы для задач оптимизации (например, задача поиска кратчайшего пути в графе) можно представить как решающие алгоритмы. Соответствующий решающий алгоритм будет отвечать на вопрос “имеет ли данный граф путь, короче, чем заданная величина?”.

Введем несколько важных определений.

Определение 1. Алгоритм A допускает слово x , если $A(x) = 1$.

Определение 2. Алгоритм A отвергает слово x , если $A(x) = 0$.

Заметим, что алгоритм может не допускать и не отвергать некоторые слова. Случается и так, что алгоритм не дает никакого ответа (например, заикливается или аварийно завершается).

Определение 3. Алгоритм A допускает язык L , если он допускает те и только те слова, которые принадлежат L .

Понимать следует и то, что алгоритм A , допускающий некоторый язык L , вовсе не обязан отвергать всякое слово, не входящее в L . Пусть, например, язык L это все Эйлеровы графы. Ясно, что граф можно описать в виде двоичного кода. Однако есть (в зависимости от двоичного представления) двоичные слова, которые не представляют никакого графа. В случае, если алгоритму представлено недопустимое слово, алгоритм может аварийно завершить работу, т.е. допускающий алгоритм для языка L может не дать ответа.

Определение 4. Алгоритм A распознает язык L , если он допускает все слова из L и отвергает все другие слова.

Определение 5. Язык L допускается за полиномиальное время, если существует допускающий алгоритм A , причем всякое слово x допускается или отвергается алгоритмом за время порядка $O(n^k)$ (т.е. время работы алгоритма A порядка $O(n^k)$), где n — длина слова, а k — некоторое не зависящее от x число. (Оно универсально для для всех слов).

Определение 6. Язык распознается за полиномиальное время, если существует распознающий алгоритм A и число k такие, что его время работы порядка $O(n^k)$, где n — длина слова.

Определение 7. Класс языков, для которых существуют полиномиальные распознающие алгоритмы называется классом P .

6.3 Типы задач

Разрешимые и неразрешимые задачи. Задача останова Задачи бывают двух типов. Неразрешимые и разрешимые. К неразрешимым задачам относится задача останова. Ставится задача определится ли программа, на заданных данных. В частном случае, если в программе нет сложных циклов или сложных переходов, то определить, что она остановится вполне возможно, однако в общем случае эта задача неразрешима.

Доказательство. Доказательство.

Программа, также как и входные данные написаны в виде последовательности из 0 и 1, следовательно отличить их в самом компьютере не представляется возможным. Тогда можно преобразовать эти последовательности в пусть и большие, но натуральные числа. Переведем последовательность соответствующую коду программы в N , а входным данным в M . Пусть

существует другая программа - Анализатор, который принимает на вход натуральные числа (N, M) .

Базовое предположение (Этап 1): Анализатор выдает 1, если программа N останавливается при входных данных M , 0 - иначе.

Этап 2: Создадим Анализатор', который выдает 1, если программа N не останавливается. Иначе он уйдет в цикл (в случае, если программа N останавливается). Замечание: Анализатор' можно создать из Анализатора

Этап 3: Введем программу Диагонализатор, который вызывает Анализатор' от N и N . Его код тоже можно зашифровать, например, K . Он останавливается только тогда, когда программа N не останавливается при входных данных M .

Вызовем Диагонализатор сам от себя. 1) Допустим он остановился. Это значит, что Анализатор' не остановился, N остановился, и тогда Диагонализатор не остановился - противоречие.

2) Допустим, что Диагонализатор не остановился, тогда Анализатор' остановился, N не остановился, и тогда Диагонализатор остановился, снова противоречие.

Получаем ошибку! Следовательно где-то было не верное предположение, а единственное предположение было в существовании Анализатора. Итог - его не существует. \square

6.4 Проверяющие алгоритмы. Класс языков NP

Определение 8. Назовем проверяющим алгоритмом A алгоритм с двумя аргументами. Первый аргумент будем называть входной строкой (слово $x \in L$), а второй — сертификатом. Проверяющий алгоритм допускает слово x , если существует сертификат $y \in \{0, 1\}^*$, такой, что $A(x, y) = 1$ для $x \in L$ и не существует сертификата для других слов, не принадлежащих L . Отметим, что сертификат может не иметь никакого отношения к языку L .

Приведем пример. Гамильтоновым циклом в графе является такой цикл, который проходит через все вершины графа ровно по одному разу. Графы, имеющие Гамильтонов цикл называются Гамильтоновыми графами. Пусть язык Ham — множество Гамильтоновых графов. Мы не умеем проверить непосредственно, является ли предъявленный граф гамильтоновым, т.е. у нас нет распознающего алгоритма, такого, что $A(x) = 1$, если слово x принадлежит языку Ham и $A(x) = 0$ в противном случае. Однако, если кто-то нам предъявит последовательность вершин и заявит, что эта последовательность является Гамильтоновым циклом, мы можем это проверить, т.е. можем проверить, что предъявленное двоичное слово действительно представляет последовательность вершин исследуемого графа, что эта последовательность является циклом и что она содержит все вершины графа. Эта последовательность может использоваться в качестве сертификата. Отметим, что последовательность вершин данного графа не является

вообще каким-либо графом. Тем не менее, эта последовательность вершин представляется в виде двоичного слова, т.е. $y \in \{0, 1\}^*$. Алгоритм, который проверяет, действительно ли y представляет гамильтонов цикл в графе x довольно легко реализовать. Ясно, что этот алгоритм должен зависеть от двух аргументов — собственно слова x , представляющего граф, и слова y , представляющего последовательность вершин — $A(x, y)$.

Определение 9. *Язык принадлежит классу NP, если существует проверяющий алгоритм, работающий за полиномиальное время $A(x, y)$, причем:*

- *Существуют такие константы C_p и p_2 , что длина сертификата ограничена полиномом: $|y| \leq C_p \cdot |x|^{p_1}$;*
- *Время работы проверяющего алгоритма ограничено полиномом от длины входного слова: $T(A(x, y)) \leq C_A \cdot |x|^{p_1+p_2}$, где C_A, p_1, p_2 — константы, не зависящие от входного слова.*

NP (nondeterministic polynomial (time)), что расшифровывается как “неопределенное полиномиальное (время)”.

Язык (или задача) тогда принадлежит классу NP, когда он легко проверяем проверяющим алгоритмом. Заметим, что мы не говорим, что умеем решать такую задачу — умеем быстро проверять, является ли представленное решение правильным. Отметим, что мы здесь нигде не конкретизируем, а как задачу представить в виде двоичных слов. Для записи тех же графов можно придумать разнообразные способы записи, т.е. разнообразные отображения $ЗАДАЧА \rightarrow \{0, 1\}^*$. Однако можно считать, что все мыслимые отображения в некотором смысле эквивалентны, т.е. длины слов в разных представлениях одной и той же задачи отличаются друг от друга не более, чем полиномиально.

Упражнение 37. *Докажите, что все языки, для которых существует полиномиальный распознающий алгоритм, принадлежат классу NP.*

6.5 Сводимость языков. NP-полные задачи (NPC)

Говоря о сводимости задач вспоминается старый анекдот. **Задача 1.** Дано: чайник, кран, газовая плита и спички. Требуется вскипятить воду. Решение: наливаем в чайник воду, зажигаем газ, ставим чайник на газ, ждем 10 минут. **Задача 2.** Дано чайник с водой, кран, газовая плита и спички. Требуется вскипятить воду. Решение. Выливаем воду из чайника и сводим задачу к предыдущей. В этом анекдоте заложена идея сводимости одних задач к другим.

Определение 10. *Язык L_1 сводится к языку L_2 , если существует отображение $f : \{0, 1\}^* \rightarrow \{0, 1\}^*$ такое, что $x \in L_1 \Leftrightarrow f(x) \in L_2$ и $x \notin L_1 \Leftrightarrow f(x) \notin L_2$. Отметим, что отображение f не является взаимно-однозначным, т.е. может не существовать обратного отображения (рис. 6.2).*

Частным случаем такого отображения является допускающий алгоритм (если существует). В этом случае язык L_2 состоит только из одного слова $\{1\}$, а отображение есть алгоритм: $f(x) = A(x)$. Если сводящее отображение $f(x)$ вычисляется за полиномиальное время от длины слова x , то говорят, что язык L_1 сводится к языку L_2 за полиномиальное время, и обозначается $L_1 \xrightarrow{P} L_2$.

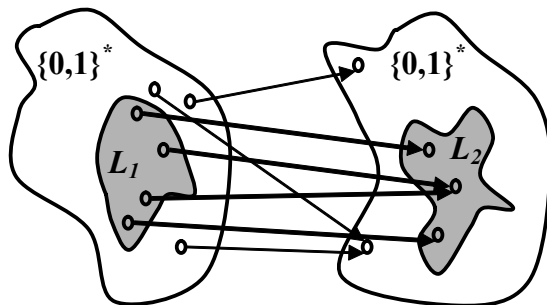


Рис. 6.2: К понятию сводимости.

Определение 11. Язык L принадлежит классу NP-полных (NP-complete, NPC), если:

1. $L \in NP$ – язык принадлежит класс NP (полиномиально проверяемые задачи)
2. $L' \xrightarrow{P} L$ для любого языка $L' \in NP$, любой такой язык сводим к L ,

Если выполняется только второе, то задача принадлежит к NP-hard.

Иными словами, NP-полные задачи это полиномиально проверяемые задачи, которые полиномиально не менее сложные, чем любая NP-задача.

Благодаря второму свойству (свойству сводимости) существует ОДНА И ТОЛЬКО ОДНА NP-полная задача. Если для хотя бы одной NP-полной задачи найдется полиномиальный алгоритм, то ВСЕ NP-полные задачи можно будет решить за полиномиальное время.

Есть гипотеза: для NP-полных задач не существует полиномиального алгоритма.

Некоторое количество NP-полных задач надо знать наизусть. Это полезно, поскольку если вы сидите и решаете какую-то задачу, стараясь найти полиномиальный алгоритм, а к этой задаче сводится NP-полная задача, не надо биться — полиномиального алгоритма решения нет. И если Вы смогли построить полиномиальный алгоритм, то, скорее всего, Вы придумали неправильный алгоритм, который, по крайней мере, в некоторых случаях не

решает Вашу задачу. Надо понимать, что для NP-полных задач нет эффективных полиномиальных алгоритмов, а есть, по-видимому, экспоненциальные. Но если Ваша практическая задача небольшая (например граф имеет немного — порядка 10 вершин), то вполне может устроить и неэффективный алгоритм полного перебора. Впрочем даже для достаточно большого размера задачи есть выход — искать приближенное решение, или, опираясь на особенности задачи решать не универсальную задачу, а суженную задачу. Эта ситуация будет рассмотрена позже.

Впрочем, если использовать вычислительные средства, в которых количество процессоров *сравнимо* с размером задачи, то можно строить эффективные алгоритмы решения таких задач. Кроме того, квантовые вычислители также могут решать сложные задачи (правда, в некотором *вероятностном* смысле). Впрочем рассмотрение этих вычислителей выходит далеко за рамки курса.

6.5.1 Иерархия типов задач

NP-hard (nondeterministic polynomial hard) относят к разрешимым задачам, они включают в себя часть NP ((nondeterministic polynomial). NP(nondeterministic polynomial) включают 2 класса: P (polynomial) и NP-complete (nondeterministic polynomial complete).

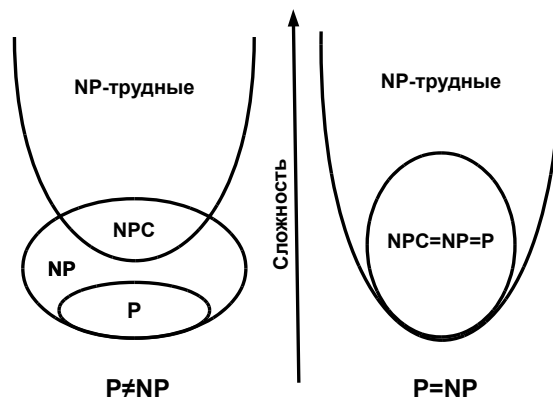


Рис. 6.3: Иерархия типов задач.

6.6 Некоторые NP-полные задачи

Поиск максимальной клики в графе. Поиск максимальной клики в графе — это поиск максимального полного подграфа (где все ребра связаны со всеми). Поиск любой клики тривиален — кликой мощности 2 является любое ребро. Задача поиска клики в графе часто встречается, в том числе и

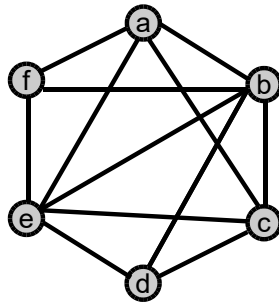


Рис. 6.4: Пример графа с 6 вершинами и с кликами в 4 вершины (abef, bcde, abce). Видно, что эти клики имеют общие вершины и ребра

в биоинформатике. К примеру, программа CluD, используемая для поиска гидрофобных ядер в белках, основана на поиске клики в специальном графе “гидрофобных групп” белка, или поиск ортологов — по большому счету это поиск клик в многодольных графах. Мы уже упоминали, что задача на максимум заменяется задачей о том есть ли в графе клика больше заданного размера. При формализации прикладной задачи (например, из области биоинформатики) надо иметь в виду, что часто максимальная клика не единственна!

Поиск гамильтонова пути. Поиск гамильтонова пути — в отличие от поиска эйлерова пути так просто не решается. Выше мы уже говорили (см. раздел про применение эйлеровости графа в биоинформатике) о возможном применении поиска гамильтонова пути и замене его поиском эйлерова пути на другом графе. Надо понимать, что такая замена возможна не всегда, иначе мы бы решили все NP-полные задачи.

Задача коммивояжера. Задача коммивояжера основана на предыдущей задаче о гамильтоновом пути. Это интересная прикладная задача. Коммивояжер (бродячий торговец) ездит по городам и предлагает свои товары. Понятно, что чем больше городов он объедет, тем больше денег заработает, при этом мы учитываем стоимость переезда из одного города (вершины графа) в другой — не каждая дорога (ребро графа) коммивояжеру выгодна. Сложность задачи еще и в том, что коммивояжер нещадно обманывает и обсчитывает своих покупателей, поэтому появляться дважды в одном городе ему нельзя. Наша задача — помочь наглomu коммивояжеру заработать как можно больше денег.

Задача о рюкзаке. Задача о рюкзаке знакома каждому туристу, путешественнику и студенту ФББ (как завсегда таю ЗБС и ББС). Есть большой рюкзак и много вещей разного заданного размера. Нужно положить в

рюкзак максимально количество вещей. Заметим, что жадный алгоритм — хватать самое большое и запихивать в рюкзак — здесь не подходит. Эта задача (несколько модифицированная) имеет приложения в биоинформатике. Речь идет о расшифровке аминокислотных последовательностей по данным масс-спектрометрии. В весьма упрощенной постановке задача выглядит следующим образом. Есть белок с неизвестной последовательностью. Его случайным образом порезали на фрагменты. Обычно это делают с использованием различных протеаз. Далее, с помощью масс-спектрометра определяют массы фрагментов. Мы знаем массы аминокислотных остатков. Вопрос: из каких остатков могут состоять фрагменты. Дальнейшая, весьма сложная работа заключается в восстановлении аминокислотной последовательности. К сожалению, этот подход непосредственно применить очень трудно еще из-за того, что есть неоднозначность в определении массы фрагментов (некоторые фрагменты могут при себе нести дополнительные химические группы, например, OH и поэтому эти модификации дают дополнительные пики в спектре).

Решение задачи о рюкзаке.

Метод "Branch and Bound" (Метод ветвей и границ)

Задача набрать предметы с максимальной суммарной стоимостью

Предполагаем, что изначально в рюкзаке лежит ϵ . Далее мы можем положить в него любой объект из имеющихся - вершина дерева. На следующем шаге, находясь в первой вершине мы можем положить уже предметы со второго до N-го. И так для каждого предмета. (Ветвление).

После построения дерева мы можем выбрать некоторое множество вариантов предметов - одна из ветвей. Выберем ее и посчитаем его вес (стоимость всех предметов). Все остальные варианты будем оценивать относительно него. Если стоимость всех продолжений этих вариантов заведомо меньше уже набранной суммы меньше, то не будем их рассматривать, а ветку, которая соответствует данному набору предметов обрежем.

Существует 2 подхода просмотра дерева:

1) Exploration (обход в ширину, при этом на каждом шаге мы просматриваем все возможные предметы, которые можно положить в рюкзак). 2) Exploitation (обход в глубину, при этом мы каждый раз выбираем предмет из возможных вариантов по уже выбранной ветке - углубляемся в одном направлении).

На самом деле задача о рюкзаке может быть решена динамическим программированием за полиномиальное время! Нет ли здесь противоречия? Дело в том, что за полиномиальное время задача решается для рюкзака с целочисленной емкостью и с целочисленными размерами предметов. Если предметы имеют не целочисленный размер, то размер матрицы динамического программирования возрастает. Поэтому, когда говорят о размере задачи (например, задачи о рюкзаке) надо понимать что размер задачи это не число предметов, а длина записи задачи. Даже если предметов немного, но все они имеют длинное число, описывающее их размер, то размер задачи становится большим. Например, 10,5,4,3,1 — достаточно короткая запись, а 9.89654, 5.109875, 3.987650981, 3.00098126745, 1.12095674 — длинная запись.

И если в первом случае задача решается быстро, то во-втором решение будет более длинным. Правда, если мы округлим размеры предметов, то быстро получим приближенное решение.

Задача о выполнимости булевой формулы (SAT от satisfaction). Задача о выполнении булевой формулы. Имеется некая (любая) булева формула от любого конечного числа булевых переменных, к примеру, такая:

$$(x_1 \& \neg x_2) | ((x_2 \& (x_3 | \neg x_1)) \& x_4) | (x_2 \& (\neg x_3 | x_4))$$

Задача заключается в том, чтобы убедиться, что существует набор значений переменных, при которых эта формула принимает значение истина.

Среди приведенных задач есть задачи оптимизации (например, максимальная клика), которые, вообще говоря не являются задачами распознавания. Как уже говорилось, их можно переформулировать в виде задач распознавания. Например задачу о максимальной клике можно сформулировать, как задачу о существовании клики заданного размера. Для того, чтобы доказать, что какая-либо задача принадлежит классу NP-полных задач необходимо, во-первых, показать существование полиномиального проверяющего алгоритма с сертификатом полиномиальной длины. Для приведенных задач очевидно, что существуют быстрые проверочные алгоритмы — их очень просто написать. А во вторых, свести к этой задаче какую-нибудь NP-полную задачу. Обратим внимание, на то что к чему надо сводить. Чтобы доказать NP-полноту задачи L надо к ней свести известную NP-полную задачу, а не наоборот.

Для задачи SAT доказана NP-полнота полностью, т.е. доказано, что любая задача из класса NP сводится к ней. Эта теорема достаточно сложна и мы не будем ее здесь доказывать. Однако мы покажем, как производится сводимость одних NP-полных задач к другим.

6.6.1 Задачи выполнимости SAT и 3CNF.

3-Conjunctive Normal Form или 3-конъюнктивная нормальная форма, сокращенно 3CNF — это комбинация из троек переменных, внутри каждой тройки между переменными стоит ИЛИ, а перед некоторыми переменными стоит отрицание \neg , а между тройками стоит И. Пример 3CNF:

$$(x_1 | x_2 | \neg x_3) \& (x_1 | \neg x_2 | x_3).$$

Комбинацию вида x или $\neg x$ называется термом, Совокупность трех термов, разделенных знаком $|$, называется группой. Следовательно, 3CNF — это конечное число групп, разделенных знаком $\&$. 3CNF сводима к SAT. Это очевидно, поскольку 3CNF является частным случаем SAT, и если у нас есть алгоритм решения задачи SAT, то автоматически мы умеем решать задачу 3CNF. Но для доказательства NP-полноты 3CNF надо наоборот свести задачу SAT к 3CNF!

Докажем, что SAT сводится к 3CNF. Для этого покажем, что для любой SAT можно написать эквивалентную 3CNF, быть может, для этого нам придется ввести дополнительные переменные. Рассмотрим пример. Пусть нам дана булева формула:

$$\varphi = ((x_1 \& (x_2 | x_3)) \& (x_2 | x_4)) | x_1 \quad (6.1)$$

Построим для этого выражения так называемое дерево разбора. Чтобы вычислить это выражение нам надо вычислить выражения в скобках, и сделать соответствующие подстановки. Введем дополнительную переменную $y_1 = x_2 | x_3$. Тогда мы сможем вычислить следующую скобку $y_2 = x_1 \& y_1$. Продолжая вводим еще переменные $y_3 = x_2 | x_4$, $y_4 = y_2 \& y_3$, $y_5 = y_4 | x_1$. Эту процедуру можно представить в виде дерева разбора (см. рис.6.5). Листья дерева — переменные. Узлы — бинарные операции. Теперь (6.1) можно

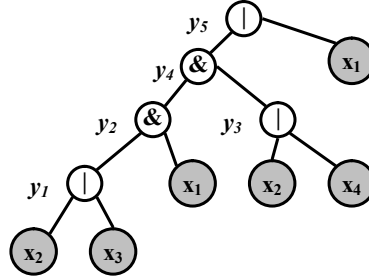


Рис. 6.5: Дерево разбора для (6.1).

переписать в виде:

$$\varphi = y_5 \& (y_5 \equiv y_4 | x_1) \& (y_4 \equiv y_2 \& y_3) \& (y_3 \equiv x_2 | x_4) \& (y_2 \equiv y_1 \& x_1) \& (y_1 \equiv x_2 | x_3)$$

Отметим, что, во-первых, эта формула эквивалентна исходной формуле. Во-вторых, она почти 3CNF, поскольку она является комбинацией из групп по 3 терма. Единственное, вместо необходимых операций “или” там встречаются разные другие операции — проверка тождественности, операции “и”. Кроме того, есть группа, состоящая из одного терма — y_5 . Самое простое — представить эту группу в виде тройки: $(y_5 | y_5 | y_5)$, которая вполне является 3-группой.

Упражнение 38. Докажите, что приведенная формула эквивалентна исходной, т.е. она принимает значение true на тех и только на тех наборах переменных x , на которых принимает значение true исходная формула.

Приведем полученные тройки к 3CNF. Для каждой тройки вида $x \equiv y | z$ либо $x \equiv y \& z$ существует всего 8 разных наборов булевых переменных, обращающих тройку в “true” либо “false”. К примеру, если все три переменные равны нулю, тройка обращается в равенство $(0=0\&0)$. Это равенство верное, ему соответствует “true” (1).

Перечислим все возможные тройки булевых переменных и соответствующие им элементы решения тройки ($x \equiv y \& z$):

тройка	Значение
0 0 0	1
0 0 1	1
0 1 0	1
0 1 1	0
1 0 0	0
1 0 1	0
1 1 0	0
1 1 1	1

Рассмотрим сочетание (0, 1, 1). Попробуем представить это значение в виде комбинации термов, объединенных операцией “|” так, чтобы оно было равно 0 при значениях переменных (0, 1, 1) и только при них. Эта комбинация будет такой: $(x \equiv y \& z) = (x | \neg y | \neg z)$. Действительно, чтобы получить 0 нам надо, чтобы все термы были равны 0. Ясно, что там, где в исходных данных была 1 нам надо применить отрицание, а там, где был 0 надо принять переменную как есть. Заметим, что тройка $(x | \neg y | \neg z)$ является группой 3CNF. Аналогично для любой тройки, результатом вычисления которой есть 0 мы определяем свой терм. Это процедура называется приведением термов.

Для тех, комбинаций, которые в результате дают 1 приведение термов не возможно. Но это не важно. Просто мы в финальную 3CNF не будем писать соответствующие группы. Эти группы равны 1, а операция &, которая объединяет группу, нейтральна по отношению к этому значению.

В результате мы каждую комбинацию вида $x \equiv y \diamond z$ заменяем некоторым количеством групп (не более 8), объединенных операцией “&”. А эти комбинации, в свою очередь, также объединены операцией “&”. Итак, мы получили 3CNF, причем ее длина не более, чем в 16 раз больше, чем длина исходной формулы.

Алгоритм приведения произвольной SAT к 3CNF:

1. Строим дерево разбора
2. Расширяем термы (вводим новые переменные)
3. Приводим термы и строим 3CNF

Итак, мы показали (не доказали), как свести SAT к 3CNF. Если рассуждения, приведенные в примере обобщить, то будет доказано, что задача SAT не сложнее задачи 3CNF. Но задача SAT является NP-полной, поэтому и задача 3CNF также является NP-полной. Обратите внимание на то, в какую сторону мы делали сводимость. Мы сводили NP-полную задачу к нашей (не наоборот!).

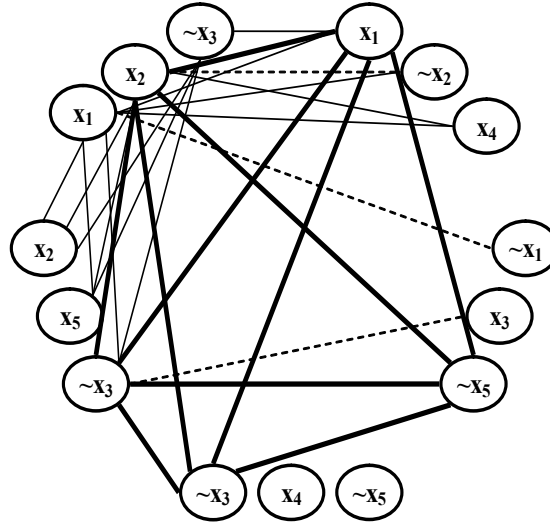


Рис. 6.6: Многодольный граф, соответствующий 3CNF. Показаны не все ребра. Пунктиром отмечены пары несовместимых вершин. Жирным выделены ребра, принадлежащие клике.

6.6.2 Задача о клике (CLIQUE).

Дан граф $G = \{V, E\}$. Выяснить есть ли в нем клика, содержащая более k вершин. Проверка решения задачи тривиальна (достаточно предъявить список из вершин), поэтому это задача класса NP. Оценим ее полноту. Для этого сведем к этой задаче другую NP-полную задачу. Покажем, что $3CNF-SAT <_p CLIQUE$.

Пусть дана некая 3CNF:

$$(x_1|x_2|\neg x_3) \& (x_1|\neg x_2|x_4) \& (x_2|x_5|\neg x_3) \& (\neg x_3|x_4|\neg x_5) \& (x_3|\neg x_3|x_3)$$

Построим по ней многодольный граф, причем такой, что если в нем есть k -клика, то эта формула выполнима. Для каждой группы рисуем 3 вершины (по одной на терм). Вершины u и v соединяем ребром, если они принадлежат разным тройкам и совместимы. Несовместимыми считаются вершины, несущие — одна некий терм, другая — его отрицание. К примеру, несовместимы вершины, несущие термы x_3 и $\neg x_3$. Граф получается очень запутанный (на рис.6.6 показаны не все ребра).

Лемма 1. Если в графе, построенном по некой 3CNF по указанному правилу, есть клика размером в число групп 3CNF, эта 3CNF разрешима.

Доказательство. Если формула разрешима, то существует набор переменных, когда все группы истинны. Тогда в каждой группе есть истинный терм,

и, следовательно, группа тоже истинна. Соответствующие вершины образуют клику. Обратное тоже верно — если есть клика размера k , то формула разрешима. \square

В нашем случае для разрешимости 3CNF нужно требовать существования клики мощностью 5. Такая клика в графе есть (выделено жирным). Для получения набора переменных, удовлетворяющего 3CNF, взяли термы, которые образуют клику в графе, и присвоили им значение 1 (“true”), значения остальных переменных не важно. При существовании показанной клики наша 3CNF разрешима с таким набором переменных: $x_1 = true, x_2 = true, x_3 = false$, так как $\neg x_3 = true, x_4$ — любое $x_5 = false$.

Итак, NP-полная задача 3CNF сводится к задаче о клике в многодольном графе. Поэтому задача о клике в многодольном графе также является NP-полной. С другой стороны, задача о клике в многодольном графе является частным случаем (следовательно, тривиально сводится) к задаче о клике в произвольном графе. Следовательно, задача о клике в произвольном графе также является NP-полной. Опять обратите внимание на направление сводимости.

6.6.3 Решение NP-полных задач

Для NP-полных задач нет полиномиальных алгоритмов. Это в настоящее время не доказано, но является предметом почти повсеместной *веры*. Это не значит, что для этих задач вообще нет алгоритмов. Алгоритмы есть! только они требуют полного перебора вариантов, и поэтому их время работы очень быстро растет с размером задачи. Однако, не все реальные задачи большие. Часто на практике встречаются задачи достаточно небольшой размерности (десятки вершин в графе). Поэтому вопрос об алгоритмах для NP-полных задач не является бессмысленным. Итак, решать NP-полные задачи надо! Мы упоминали выше, какими важными являются, к примеру, задачи о клике и гамильтоновом пути, как много применений имеют эти задачи.

Во-первых, не все NP-полные задачи для общего случая являются таковыми для некоторых специальных случаев. Например, задача о максимальной клике является NP-полной для общего вида графов, но является тривиальной для деревьев. Дерево по определению ациклический граф, а в любой клике мощностью больше трех всегда есть циклы. Поэтому, даже не глядя на дерево, можно сказать, что мощность максимальной клики равна двум.

Во вторых, для NP-полных задач существуют эвристические алгоритмы решения NP-полных задач. Эвристическими называются методы, сокращающие полный перебор при решении задачи. Такие методы дают не точные, а приближенные решения. Например, ищут не максимальную клику, а достаточно большую, или ищут не совсем клику, а очень плотный подграф. Для эвристических алгоритмов нужно делать оценку качества их работы.

Часто можно сделать некоторые предварительные действия для того, чтобы сократить перебор. Например, в задаче о поиске клики хорошо бы

сначала выделить связные компоненты, поскольку клики по определению являются связными компонентами. Далее, можно попробовать в графе выделить в графе Другие структуры. Ясно, что клика размером более двух вершин обязательно содержит циклы. При декомпозиции графа мы ищем в графе определенные точки и составные части, по которым граф можно было бы разделить на подграфы. Такое деление часто облегчает задачу, к примеру, задачу поиска клики. Неориентированный граф можно разбить на узлы, мосты, деревья и двусвязные компоненты (см.рис.6.7):

Двусвязная компонента. Подграфы, в которых две любые вершины принадлежат общему циклу

Узел. Вершина, соединяющая две двусвязные компоненты графа, или, что то же самое связные компоненты, между которыми нет ребер.

Мост. Цепочки ребер, соединяющие две двусвязные компоненты, между которыми больше нет ребер.

Дерево. Ациклический подграф.

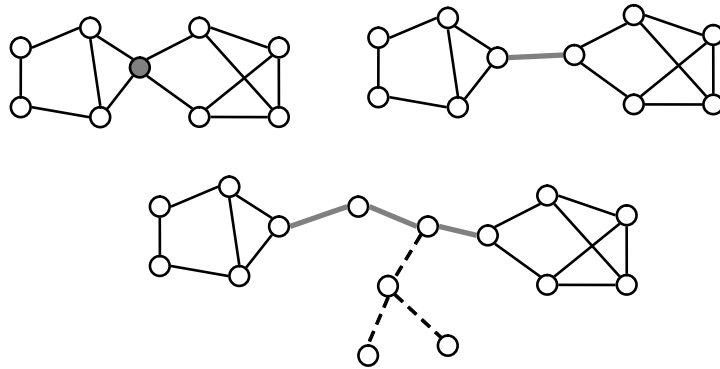


Рис. 6.7: Двусвязные компоненты (показаны обычными линиями), узлы (серые вершины), мосты (толстые серые ребра), деревья (показаны штриховой линией).

Если в графе несколько двусвязных компонент, соединенных ребрами, мостами и т.п., гораздо удобнее искать клику в каждой компоненте по отдельности, поскольку клики мощностью больше, чем количество вершин в самой большой двусвязной компоненте, в таком графе нет. Кстати, в отличие от клик, для двусвязных компонент существует эффективный алгоритм поиска.

6.6.4 Переборные алгоритмы

Обычно перебор ведется с помощью обхода некоторого дерева. Рассмотрим задачу суммы чисел. Дан набор положительных чисел $S = \{x_1, x_2, \dots, x_n\}$

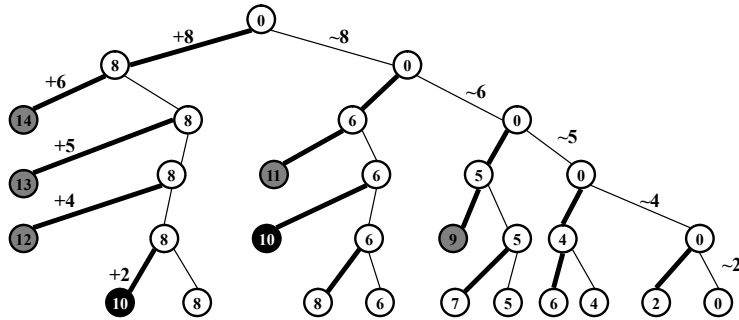


Рис. 6.8: Дерево поиска оптимальной загрузки рюкзака. Множество весов: 8,6,5,4,2, порог равен 10. толстые ребра отвечают включению элемента в список, а тонкие — исключению из списка. Серые вершины отвечают бесперспективным путям, черные вершины — найденным решениям.

и некоторое положительное число t . Надо найти такое подмножество $T = \{x_{i_1}, x_{i_2}, \dots, x_{i_k}\} \subseteq S$, чтобы их сумма $s = \sum_j x_{i_j}$ была максимальной, но не превышала t :

$$s \rightarrow \max |s \leq t$$

Пусть множество S упорядочено по убыванию (пусть первым будет самый тяжелый элемент). На первом шаге у нас есть выбор: брать первый элемент или нет. Если он тяжелее ограничения, то мы заведомо его не берем (и не рассматриваем вариантов, содержащих его). На каждом следующем шаге мы стоим перед выбором: или брать очередной элемент или не брать. При этом проверяем перспективность дальнейшего. Простейшее правило остановки — если добавление самого легкого элемента переводит за порог, то дальнейшее продвижение не перспективно. Если он превышает ранее зафиксированный результат, то запоминаем его в качестве максимума.

Здесь возможна еще одна эвристика. Если мы в каждой позиции будем помнить сумму весов последующих элементов, то мы сможем отсекал больше неперспективных путей. В этом алгоритме мы ветвимся и отсекаем неперспективные пути. Поэтому метод называется методом ветвей и границ.

Упражнение 39. Как можно использовать эвристику суммы остатков?

6.7 Стохастические алгоритмы

Допустим, нам нужно посчитать такой замечательный интеграл:

$$\int_0^1 \int_0^1 \int_0^1 \int_0^1 \int_0^1 \int_0^1 \int_0^1 f(x, y, z, t, u, v, w) dx dy dz dt du dv dw$$

Функция определена в семимерном пространстве, и считать точное значение интеграла очень долго. При любом регулярном методе подсчета нам нужно семь вложенных циклов. Для разумной точности надо в каждом цикле сделать порядка ста шагов. Итого получаем примерно 10^{14} операций. Можно предложить более эффективный способ подсчета. Воспользуемся знанием математической статистики и попробуем посчитать значение интеграла как среднее по большой выборке частных значений. Для этого берем какие-то конкретные случайные значения всех семи переменных из равномерного распределения на отрезке $[0, 1]$: $rx, ry \dots rw$ и подставляем в функцию f , получаем некое значение этой S . Так делаем много — 10000 уже вполне достаточно — раз. Берем все 10^4 значений S и считаем среднее. Получаем значение интеграла с хорошим приближением. При подсчете этого интеграла на сетке с шагом 10^{-2} объем вычислений составляет 10^{14} , что на много порядков больше, чем стохастический подсчет.

Идея использования случайных чисел широко применяется для решения многих сложных задач оптимизации. Эта группа алгоритмов имеет общее название методы Монте-Карло.

6.7.1 Искусственный отжиг

Пусть нам надо оптимизировать (допустим, минимизировать) некую функцию на некотором множестве. Дано множество S (конечное или бесконечное), функция $f : S \rightarrow R$. Задача: найти элемент множества $s^* \in S$ такой, что $f(s^*) \rightarrow \min$. Пусть, например, мы хотим найти Гамильтонов цикл в графе. На самом деле, мы заранее не знаем есть ли такой цикл в графе, но понимаем, что хотелось бы найти такой цикл, в котором вершины не слишком часто повторялись и чтобы поменьше было не пройденных вершин. Получаем задачу оптимизации, в которой S — множество циклов в графе, а функция f — число “неправильных” вершин.

Множество S мы в дальнейшем будем называть множеством состояний, а функцию f энергией. На множестве состояний S строится неориентированный граф: вершины соединяются ребром, если эти состояния (в некотором смысле) не сильно отличаются. Например, в случае поиска Гамильтонова цикла, множество S — это множество циклов, а два цикла называются соседними, если они имеют пересекающиеся множества ребер. Важно, чтобы граф состояний был связным, т.е., чтобы из любого состояния в другое состояние был бы путь через некоторые промежуточные состояния. Надо понимать, что граф состояний очень большой и не описан явно. Теперь мы начинаем блуждать по графу, просматривая соседние состояния. Для этого нам не нужно полное описание графа состояний, но для каждого состояния надо уметь находить множество соседей $\{s_i\}_{neighbor}$. Важно, чтобы множество соседей было не очень большим.

Находясь в состоянии s выбираем произвольное состояние s_{new} , принадлежащее множеству соседних состояний $\{s_i\}_{neighbor}$. Рассмотрим, насколько изменится энергия, если мы перейдем в новое состояние s_{new} : $\Delta f = f(s_{new}) - f(s)$. Чем меньше это значение, тем полезнее будет пе-

реход системы в состояние s_{new} с точки зрения минимизации энергии f . Если переход в какое-то соседнее состояние оказался энергетически выгоден ($\Delta f < 0$), то мы переходим в новое состояние. Казалось бы мы так будем всегда уменьшать энергию, пока не дойдем до минимума. Однако этот минимум будет локальным, а движение только в сторону уменьшения энергии исключает возможность нащупать глобальный минимум. Кроме того, если несколько соседних состояний лучше текущего, то, выбрав неправильное направление, мы можем попасть в тупик. Поэтому нам надо оставить возможность выбираться из локального минимума и переваливать через высокие барьеры. Граф состояний не ориентирован, следовательно, если состояние s_{new} является соседом состояния s , то и состояние s является соседом состояния s_{new} . Поэтому, введя обратимость переходов, мы всегда можем вернуться в исходное состояние.

Вспоминая статистическую физику, мы знаем, что отношение вероятности перехода из состояния a в состояние b к вероятности обратного перехода определяется законом:

$$\frac{p(a \rightarrow b)}{p(b \rightarrow a)} = \exp\left(-\frac{G(b) - G(a)}{kT}\right)$$

Поскольку одна из вероятностей равна 1 (всегда переходим в сторону уменьшения энергии), то другая вероятность должна определяться из этого соотношения. Получаем формулу для определения вероятностей переходов:

$$p(s \rightarrow s_{new}) = \begin{cases} 1 & \Delta f < 0 \\ \exp\left(-\frac{\Delta f}{T}\right) & \Delta f \geq 0 \end{cases}$$

Теперь можно описать весь алгоритм.

1. Выбираем произвольное допустимое состояние $s \in S$
2. Для текущего состояния выбираем произвольное соседнее состояние $s_{new} \in S$
3. Вычисляем вероятность перехода $p(s \rightarrow s_{new})$
4. Разыгрываем случайную величину r , распределенную равномерно на $[0, 1]$. Если $p(s \rightarrow s_{new}) < r$, то переходим в новое состояние.
5. Переходим к п. 2.

Здесь надо заметить, что если температура T очень велика по сравнению с Δf , то мы будем просто свободно блуждать по пространству состояний. Если же $T \ll \Delta f$, то вероятности невыгодных переходов будут очень малы и мы будем быстро приходить в локальный минимум. Поэтому в этом алгоритме используют процедуру понижения температуры. Обычно температуру полагают равной

$$T = \frac{1}{\alpha\sqrt{i}},$$

где α — некоторый параметр, i — номер итерации.

Существует теорема, что при достаточно плавном понижении температуры процедура сходится к глобальному минимуму. Однако скорость сходимости сильно зависит от стратегии понижения температуры, от определения функции энергии, от определения пространства состояний и множества соседних состояний. Поскольку в процессе работы алгоритма мы постоянно понижаем температуру, то эта процедура напоминает отжиг. Поэтому алгоритм называется алгоритмом искусственного отжига (Annealing). Авторство такого подхода принадлежит Метрополису и Гастингсу (1954). Другой алгоритм, имеющий аналогию в статистической физике, устроен точно также, но без понижения температуры. В процессе блуждания по пространству мы в каких-то состояниях бываем чаще, а в каких-то реже. Наиболее посещаемые состояния на самом деле являются оптимальными. Этот алгоритм называется Гиббс-сэмплер. Оба подхода имеют достаточно широкое применение в биоинформатике.

6.7.2 Генетические алгоритмы

Генетические алгоритмы — стохастические, эвристические оптимизационные методы, впервые предложенные Холландом (1975). Предложенная Ч. Дарвином в 1859г. эволюционная теория сильно повлияла на мировоззрение людей с самого ее появления. Дарвин выявил главный механизм развития живого: отбор в сочетании с изменчивостью. В общем-то, неудивительно, что ученые, занимающиеся компьютерными исследованиями, обратились к теории эволюции в поисках вдохновения. Итак, идея генетического алгоритма заимствована у живой природы и состоит в организации эволюционного процесса, конечной целью которого является нахождение оптимального решения в сложной комбинаторной задаче. Разработчик генетических алгоритмов выступает, в данном случае, как “создатель”, устанавливающий законы “эволюции”, для наиболее быстрого достижения желаемой цели.

Введем основные понятия, важные для генетических алгоритмов. Допустимое решение описывается как массив значений (например, битов или байтов) и называется **генотипом** или **особью** — в этой науке это одно и то же. Значение каждого признака (бита, байта) называется **аллельным состоянием**. Допустимое решение должно обладать свойством: Если мы возьмем часть значений из одного допустимого решения, а остальные значения из другого, то получим снова допустимое решение.

Функция приспособленности f (от слова fitness) — вычисляется по генотипу и характеризует качество соответствующего допустимого решения.

Мутация — изменение аллельного состояния одного из признаков

Кроссинговер — порождение новой особи в результате скрещивания двух родительских особей. Новая особь получает часть генов от одного ро-

дителя, остальные — от другого. Обычно используется схема случайного и независимого выбора (сцепленности генов здесь нет)

Популяция — совокупность из некоторого количества особей (обычно размер популяций составляет сотни особей)

Генетический алгоритм:

1. Создается популяция случайных особей.
2. Выбирается несколько произвольных особей (это количество — параметр алгоритма) и проводится мутагенез по некоторому количеству случайных аллелей.
3. Выбирается произвольная пара особей и осуществляется скрещивание. В результате появляется новая особь.
4. С вероятностями, зависящими от функции приспособленности f (например $p_i = \exp(-f_i)/Z, Z = \sum_i \exp(-f_i)$) выбирается жертва естественного отбора, и соответствующая особь удаляется из популяции.
5. Переход к п.2.

Эволюция популяции приводит к появлению оптимальных особей (набор аллелей становится примерно одинаковым у каждой). О сходимости можно судить по двум признакам — по среднему по популяции значению функции приспособленности и однородности популяции, т.е. по дисперсии функции приспособленности. Успех или не успех применения генетического алгоритма в большой степени зависит от искусства. Надо правильно построить генотипы, правильно сформировать функцию приспособленности, правильно подобрать параметры мутирования. Обычно хорошие результаты можно получить уже через 40 “поколений” популяции (каждое поколение означает выполнение действий, описанных в 1-3 пунктах алгоритма).

Правильно запрограммированные генетические алгоритмы могут быть очень эффективными. Применение их весьма широко: это и разнообразные задачи на графах (задача коммивояжера, например), и оптимизация запросов в БД, и биоинформатика (свертывание белков, докинг) и многое другое.

Упражнение 40. Постройте схему генетического алгоритма для поиска максимальной клики в графе.

Заключение

Описанные в этой книге структуры, идеи, алгоритмы являются лишь небольшой частью айсберга, имя которому теоретическая информатика. Существует и используется гораздо более широкий арсенал подходов, и методов решения разного рода алгоритмических задач. Искусство применения теоретической информатики заключается, прежде всего, в том, чтобы для прикладной задачи найти правильную формализацию и разглядеть черты уже известных задач теоретической информатики. В большинстве удается новую прикладную задачу свести к уже известной алгоритмической задаче. Однако, одна и та же прикладная задача может быть формализована разными путями, что может привести к разным не эквивалентным алгоритмическим задачам.

Хочется предупредить читателя, что изобретательство велосипеда дело достойное и интересное, но мало продуктивное. Поэтому прежде, чем пытаться построить новый алгоритм для решения какой-либо прикладной задачи, хорошо бы посмотреть, а нет ли уже известных подходящих алгоритмов. Если же не удалось найти подходящего алгоритма решения задачи, и Вы построили новый алгоритм, то надо озаботиться тем, чтобы попробовать доказать теоретически, что предложенный алгоритм действительно решает поставленную задачу. Строить собственные алгоритмы необходимо не с нуля, а опираясь на уже существующие алгоритмы.

В заключении можно порекомендовать следующую литературу:

- Томас Х. Кормен, Чарльз И. Лейзерсон, Рональд Л. Ривест, Клиффорд Штайн. Алгоритмы, построение и анализ
- Дональд Кнут. Искусство программирования
- Дэн Гасфилд. Строки, деревья и последовательности в алгоритмах
- <http://rain.ifmo.ru/cat/view.php/aboutus/vis-authors>
- Google

Предметный указатель

Предметный указатель

- BLAST, 66
- byte, 12
- Merge_Sort, 36
- Merge_Sort, время работы, 38
- NP, 176
- Partition, 40
- QSort, 39
- RLE, 104
- автомат Ахо-Корасик, 91
- алгоритм, 10
 - Ахо-Корасик, 87
 - Беллмана-Форда, 158
 - Дейкстры, 151
 - время работы, 156
 - Кнута-Морриса-Пратта, 74
 - Рабина-Карпа, 72
 - Форда-Фалкерсона, 161
 - генетический, 190
 - допускающий, 174
 - искусственного отжига, 188
 - переборный, 186
 - проверяющий, 175
 - распознающий, 174
 - решающий, 173
 - стохастический, 187
- аллельное состояние, 190
- алфавит, 79, 87, 173
- архивирование, 100, 102, 104, 105
- байт, 12
- бинарное дерево поиска, 51
 - добавление, 52
 - максимальный элемент, 54
 - минимальный элемент, 54
 - поиск, 52
 - преобразования, 58
 - сбалансированное, 58
 - следующий элемент, 54
 - удаление элемента, 56
- блок схема, 16
- булева алгебра, 18
- быстрая сортировка, 39
- время работы алгоритма, 11
- выравнивание, 125, 147
- генотип, 190
- граф, 121
 - к-дольный, 122
 - Гамильтонов, 175, 179
 - Эйлеров, 135
 - Эйлеров путь, 135
 - Эйлеров цикл, 135
 - ациклический, 121
 - взвешенный, 121
 - двудольный, 122
 - паросочетания, 165
 - двусвязная компонента, 125, 134, 186
 - задача Эйлера, 135
 - клика, 122, 125, 133
 - максимальная, 178
 - матрица смежности, 124
 - минимальный путь, 151
 - многодольный, 184
 - обход в глубину, 129, 143
 - обход в ширину, 125
 - время работы, 126
 - оптимальный путь, 143

- динамическое программирование, 144
- покрывающее дерево, 132
- полный, 122
- расстояние между вершинами, 128
- связная компонента, 126, 132, 133
- связный, 122
- сети
 - поток, 159
 - сильно связный, 122
 - смешанный, 122
 - типы ребер, 131
 - топологическая сортировка, 141
 - цикл, 132
 - Гамильтонов, 179
- граф де-Брюина, 138
- граф неориентированный, 121
- граф ориентированный, 121
- дерево, 122, 186
 - суффиксное
 - алгоритм построения, 95
 - Ахо-Корасик, 88
 - красно-черное, 59
 - поиска, 34, 51
 - добавление, 52
 - красно-черное, 59
 - максимальный элемент, 54
 - минимальный элемент, 54
 - поиск, 52
 - преобразования, 58
 - сбалансированное, 58
 - следующий элемент, 54
 - удаление элемента, 56
 - покрывающее, 132
 - суффиксное, 94
 - наибольшее общее подслово, 98
 - неявное, 95
 - поиск образца, 98
 - поиск повторов, 98
- динамическое программирование, 151
- время работы, 147
- жадный алгоритм, 102
- задача
 - SAT, 181
 - Эйлера, 135
 - класс
 - NP полная, 177
 - коммивояжера, 179
 - о Гамильтоновом пути, 179
 - о выполнимости 3-конъюнктивной нормальной формы, 181
 - о выполнимости булевой формулы, 181
 - о максимальной клике, 178, 184
 - о рюкзаке, 179
- задача останова, 174
- индекс вершины, 122
- источник, 160
- картирование на геном, 109
- клика, 122, 125, 133
- код Хаффмана, 102
- коллизия, 63
- компактное
 - хранение, 102
- компонента
 - двусвязная, 122
 - связная, 121
 - сильно связная, 121
- компьютер, 10
- конечный автомат, 79
 - граф переходов, 80
 - детерминированный, 79
 - недетерминированный, 85
 - построение, 85
 - поиск образца, 81
 - построение, 82, 85
 - состояние, 79
 - функция перехода, 79
- красно-черное дерево, 59
 - вставка, 61
 - поиск, 61
 - удаление, 61, 62
 - черная высота, 60
- кроссинговер, 190

- лексикографический порядок, 32
лес, 122
линейное программирование, 166
- массив, 13
 добавление элемента, 31
 поиск, 29
 слияние, 37
 сортированный, 41
 добавление элемента, 35
 поиск, 32
 удаление элемента, 35
 сортировка, 36, 171
 оценка времени, 171
 разделением, 39
 слиянием, 36
 суффиксный, 93
 удаление элемента, 31
 упорядоченный, 32, 41
- машина Тьюринга, 21, 22
- метод
 ветвей и границ, 187
- минимальный разрез, 164
- многопроцессорные вычисления, 11
- мост, 186
- мутация, 190
- наибольшее общее подслово, 98
- образец, 69
 структура, 74
- описатель универсальный, 101
- ортологи
 поиск, 125
- особь, 190
- оценка времени
 в среднем, 12, 40, 71
 в худшем, 12, 40, 70
- очередь, 48, 125
- очередь с приоритетами, 152
- память, 10
 распределение, 47
- паросочетания, 165
- паттерн, 69
 структура, 74
- повтор, 75, 98
- подграф, 121
- подстрока, 69
- поиск
 многих образцов, 87, 92
 повторов, 98
 подстроки
 алгоритм Кнута-Морриса-Пратта, 74
 алгоритм Рабина-Карпа, 72
 конечный автомат, 79
 наивный алгоритм, 70
- поиск образца, 75, 98
- покрывающее дерево, 132
- полукольцо, 150
- популяция, 191
- посимвольное сжатие, 102
- преобразование Барроуза-Уиллера, 109
- препроцессинг Кнута-Морриса-Пратта, 77
- префикс, 69, 75, 147
 собственный, 69
- префикс-функция, 74
- программа, 11
- пропускная способность, 160
- псевдограф, 121
- псевдокод, 16
- путь, 121
 простой, 121
 элементарный, 121
- разбор формулы, 49
- разделяй и властвуй, 36, 39
- разрез, 162
- расстояние
 редакционное, 147
- регулярное выражение, 84
- редакционное расстояние, 147
- рекурсия, 33, 36, 129
- сводимость задач, 177
- сводимость языков, 177
- связная компонента, 133
- секвенирование геномов, 138
- сертификат, 175

- сети
 - поток, 159
- сеть
 - источник, 160
 - остаточная, 161
 - пропускная способность, 160
 - разрез, 162
 - минимальный, 164
 - пропускная способность, 164
 - сток, 160
- сжатие
 - метод RLE, 104
 - по Лемпелю-Зиву, 104
 - по Лемпелю-Зиву-Вельчу, 105
 - текста, 100, 104
- сжатие со словарями, 105
- симплекс метод, 169
- система описания текста, 100
- слияние массивов, 37
- сложность текста
 - относительно программы, 101
 - по Колмогорову, 100
- случайный текст по Колмогорову, 101
- сортировка
 - всплывающий пузырек, 17
 - массива, 36
 - наивная, 16
 - разделением, 39
 - слиянием
 - время работы, 38
- состояние, 188
- список
 - двусвязный, 45
 - добавление, 46
 - поиск, 46
 - сортировка, 46
 - удаление, 47
 - односвязный, 42
 - вставка элемента, 43
 - поиск, 42
 - удаление, 44
 - циклический, 46
- статистическая физика, 189
- стек, 49, 129
- сток, 160
- строка, 69
- структура, 14
- суффикс, 69
 - собственный, 69
- суффикс-префикс, 75, 77
- суффикс-функция, 82
- суффиксное дерево, 94
 - алгоритм построения, 95
 - наибольшее общее подслово, 98
 - неявное, 95
 - поиск образца, 98
 - поиск повторов, 98
- суффиксный массив, 93
- текст, 69
- температура, 189
- теорема
 - Фалкersona, 165
- точное вхождение, 69
- узел, 186
- указатель, 14
- универсальный описатель, 101
- функция неудач, 88
- функция приспособленности, 190
- хеш
 - таблица, 63
 - коэффициент заполненности, 65
 - функция, 63
- хранение
 - текста, 102
- целое число без знака, 12
- целое число со знаком, 12
- цепь, 121
- цикл, 121, 132
 - Гамильтонов, 175
 - Эйлера, 136
- цикл Эйлера, 135
- черная высота, 60
- число с плавающей точкой, 13

энергия, 188

язык, 87, 173

класс

NP, 176

NP полный, 177