

Матосновы алгоритмов

Кабашный Иван (@keba4ok)
на основе лекций А. С. Охотина и А. В. Тискина

22 января 2020 г.

https://users.math-cs.spbu.ru/~okhotin/teaching/algorithms_2020/ - источник с конспектами А. С. Охотина, где всё в тысячу раз подробнее и грамотнее.

Основные моменты.

Содержание

1 Лекция 2.	3
1.1 Быстрая сортировка.	3
1.2 Сортировка кучей.	3
1.3 Скорость сортировки.	3
1.4 Нахождение i -го по величине элемента массива.	4
1.5 Метод динамического программирования.	4
2 Лекция 3.	4
2.1 Продолжение динамического программирования.	4
2.2 Нахождение наибольшей общей подпоследовательности.	5
2.3 Поиск в ориентированном графе.	6
3 Лекция 4.	7
3.1 Окончание поисков в орграфе.	7
3.2 Поиск в алгоритме с весами.	7
3.3 Окончание поиска с весами.	8
3.4 Нахождение минимального остовного дерева.	9
3.5 Структура данных для непересекающихся множеств.	10
4 Лекция 5.	11
4.1 Пути между всеми парами вершин.	11
4.2 Быстрое умножение матриц.	11
5 Лекция 6.	13
5.1 Двоичные деревья поиска	13
5.2 В-деревья.	14
6 Лекция 7.	14
6.1 Полиномиальное хэширование.	14
6.2 Алгоритм Кнута-Морриса-Пратта.	15
6.3 Поиск с помощью конечных автоматов.	16
6.4 Представления множеств строк.	16
6.5 Поиск нескольких строк одновременно: алгоритм Ахо-Корасик.	17
6.6 Суффиксные деревья.	17
7 Лекция 9.	17
7.1 Сжатие данных, основанное на повторении строк.	17

1 Лекция 2.

1.1 Быстрая сортировка.

Алгоритм 1. Быстрая сортировка. Выбираем *опорный элемент*, с которым сравниваем все остальные элементы (на это уходит линейное время). Затем рекурсивно работаем с тем, что справа от него и слева от него.

Теорема 1. Если все элементы массива различны и опорный элемент выбирается случайно, то среднее время работы алгоритма - $\Theta(n \log n)$.

Доказательство. Время работы пропорционально числу сравнений между элементами. Рассматриваем два элемента y_i и y_j , $i < j$, тогда они сравниваются только, если выбран один из них в качестве опорного. Если будет выбран какой-то y_k , $i < k < j$, то они никогда больше не будут сравнены, если что-то на отрезке не между ними - плевать, относительно отрезка между ними ничего не поменялось. Тогда среднее количество сравнений между этими элементами:

$$\frac{2}{j - i + 1}.$$

Тогда всего среднее количество сравнений:

$$\sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j - i + 1} = \sum_{i=1}^{n-1} \sum_{k=1}^{n-1} \frac{2}{k + 1} = O(n \log n).$$

Последняя оценка получается из

$$\sum_{k=1}^n \frac{1}{k} \approx \int_1^n \frac{1}{x} dx = \ln n.$$

□

1.2 Сортировка кучей.

Алгоритм 2. Сортировка кучей. Для начала, мы строим дерево: записываем по порядку все вершины так, что у вершины x_i потомки - x_{2i} и x_{2i+1} . Затем начинаем на все вершины, кроме висячих смотреть и делать вот что: если она меньше потомка, то меняем с ним (если меньше обоих, то с меньшим). Так доведём её до куда сможем, и продолжим рассмотрение для оставшихся невисячих вершин (в изначальном дереве). Так сверху окажется наименьшая вершина, вынесем её, затем - по индукции.

Утверждение 1. Работает за $O(n \log n)$ (построение дерева - $O(\log n)$, вынесение вершин - $O(n)$), можно разогнать оценку до $2n$.

1.3 Скорость сортировки.

Теорема 2. Всякий алгоритм сортировки, основанный на сравнении, требует $\Omega(n \log n)$ операций сравнения.

Доказательство. Построим дерево того, как мы спускаемся к определению последовательности, его высота ограничивается $\log_2 n!$, оценим факториал $\left(\frac{n}{e}\right)^n < n! < n^n$, а после - оценим через это логарифм $n \log_2 n - O(n)$. □

Алгоритм 3. Сортировка подсчётом. Если у нас есть массив из конечного обозримого количества типов элементов, можно сначала посчитать количество первого, затем количество второго, и так далее. Время работы - $O(n + k)$, где k - количество типов, n - количество переменных.

Алгоритм 4. Поразрядная сортировка. Сортируем числа сначала по первому разряду, затем по второму, и так далее... Время работы: $O(l(n + k))$, где сравниваются строки длины l , алфавит из k символов.

1.4 Нахождение i -го по величине элемента массива.

Алгоритм 5. Нахождение i -го элемента. Делим массив на пятёрки подряд идущих элементов (возможно, последняя пятёрка будет неполной). Теперь в каждой пятёрке выделяем медианы, и смотрим на медиану медиан. Сделаем её опорным элементом и как в быстрой сортировке, раскидаем всё по сторонам. Если этот элемент под номером i , то мы его нашли, иначе - действуем рекурсивно с одной из сторон. Время работы - линейное.

1.5 Метод динамического программирования.

Задача 1. Имеется стержень длины n . Продав стержень длины i , можно выручить p_i денежных единиц. Как выгоднее всего распилить имеющийся стержень?

Алгоритм 6. Начинаем с первого, и делаем полный перебор. Говнище

Алгоритм 7. Жадный алгоритм. Отпиливаем самый дорогой кусок, затем опять самый дорогой из возможных, и так далее. Не самый оптимальный.

Алгоритм 8. Метод динамического программирования. Суть этого метода такова. Пусть на каждом шаге надо сделать выбор (принять решение). Известно, что какой-то выбор приводит к оптимальному результату. Этому выбору соответствует некий набор подзадач. Тогда сперва находятся ответы для всех подзадач данной задачи, возникающих при различном выборе, после чего, имея все эти ответы перед глазами, можно будет в каждом случае сделать наилучший выбор.

Пример(ы) 1. Пусть стержень длины 0 не стоит ничего. Для j от 1 до n пока не найдено никаких способов продать стержень, для всякой длины отрезаемого куска, сложим его цену с выручкой за остаток. Если так можно выручить больше известного, то цена стержня длины j улучшается, и так рекурсивно мы дойдём до получения цены за весь стержень.

2 Лекция 3.

2.1 Продолжение динамического программирования.

Задача 2. Пусть нужно умножить n матриц $M_1 \times \dots \times M_n$. В силу ассоциативности, скобки можно расставить как угодно. От их расстановки зависит общее число операций, и, чтобы умножить матрицы быстрее, надо заранее определить наилучший порядок их умножения.

Пример(ы) 2. Строим верхнедиагональную матрицу T , в которой $T_{i,j}$ - наименьшее число действий, необходимых для вычисления $M_{i+1} \times \dots \times M_n$.

Внешний цикл по длине куска $l = j - i$, второй - по i , во внутреннем перебираются все разбиения произведения на два, и вычисляется следующее значение:

$$T_{i,j} = \min_{k=i+1}^{j-1} (T_{i,k} + T_{k,j} + m_i m_k m_j).$$

Разбираем так все по порядку и вычисляем наилучший способ. Время раюоты: $O(n^3)$ - строим таблицу, далее - $2n - 1$ вызовов процедуры перемножить (i, j) , в каждом - $O(n)$ итераций цикла. И ещё само умножение матриц.

Примечание 1. Для простого понимания - простой пример с кузнечиком, который прыгает на 1 или 2, и ему нужно пропрыгать n , сколькими способами это можно сделать? Мы заводим массив $dp[i]$ длины n (кол-во способов добраться до i), тогда $dp[i] = dp[i-1] + dp[i-2]$, и так насчитываем все значения, находим ответ для n .

2.2 Нахождение наибольшей общей подпоследовательности.

Определение 1. Строкой над алфавитом Σ называется всякая конечная последовательность $w = a_1 \dots a_l$, где $l \geq 0$, и $a_1, \dots, a_l \in \Sigma$ - символы.

Алгоритм 9. Народный алгоритм. Динамический способ нахождения наибольшей общей подпоследовательности. Заводим таблицу T и в ячейке $T_{i,j}$ записываем длину наибольшей общей подпоследовательности на префиксах длины i и j первого и второго слова соответственно. Заполняем таблицу последовательно от более коротких мар до самых длинных, и в итоге получим ответ в задаче.

Строим таблицу так: берём $T_{i,j}$. Если у них одинаковые последние элементы, то получим $T_{i-1,j-1} + 1$. Если они разные, то $\max(T_{i-1,j}, T_{i,j-1})$.

Саму последовательность элементов потом восстанавливаем с конца понятно как. Недостаток в том, что чтобы найти подпоследовательность, нужно хранить всю таблицу, а это $O(mn)$, и это много. Однако, если нужна только длина, то можно ограничиться лишь двумя столбцами (или двумя строчками).

Алгоритм 10. Алгоритм Хиршенберга. Построение наибольшей общей подпоследовательности за время $O(mn)$, используя память $O(\min(m, n))$. Пусть $u = u'u''$ - некоторое разбиение u . Тогда оптимальное совмещение u и v совмещает u' с каким-то начальным куском v - пусть это v' , и u'' - с остатком v'' . Нужно найти это разбиение $v = v'v''$, чтобы потом отдельно запустить совмещение двух соответствующих пар кусков.

Алгоритм делит u на две подстроки примерно равной длины. Сперва динамическим программированием находится последняя строчка таблицы $T^{u',v}$, как в базовом алгоритме. Её j -ый элемент содержит длину наибольшей общей подпоследовательности u' и u_j - префикса длины j . Аналогично находится последняя строка таблицы $T^{(u'')^R, v^R}$ (R - reverse). Далее складываем таблицы поэлементно и посмотрим, где достигается максимум - это и есть искомое разбиение $v = v'v''$. Для этого вычисления алгоритм использовал $O(|v|)$ ячеек памяти, которые теперь можно освободить.

Далее алгоритм вызывается рекурсивно, чтобы вычислить лучшее совмещение u' и v' , и u'' и v'' . Полученные совмещения последовательно приписываются друг к другу.

Теорема 3. Алгоритм Хиршберга работает за время $O(mn)$.

Доказательство. Принимая за единицу времени время, затрачиваемое на вычисление значения одного элемента $T_{i,j}$ в "народном" алгоритме, утверждается, что в общей сложности будет выполнено не более, чем $2mn$ шагов.

Пусть $f(m, n)$ - время работы в наихудшем случае. Тогда индукцией по m и n доказывается неравенство $f(m, n) \leq 2mn$. При запуске на строках u и v , где их мощности соответственно равны m и n , вычисление таблицы $T_{u,v}$ займёт $\frac{1}{2}mn$ шагов, и за столько же шагов будет вычисляться таблица $T_{(u'')^R, v^R}$. После этого проводятся два рекурсивных вызова, один из которых занимает $f(\frac{m}{2}, k)$ шагов, а другой - $f(\frac{m}{2}, n - k)$ шагов, для некоторого k . Время работы рекурсивных вызовов оценивается по предположению индукции, откуда получается оценка того же вида для $f(m, n)$.

$$f(m, n) = 2 \cdot \frac{1}{2}mn + \max_k (f(\frac{m}{2}, k) + f(\frac{m}{2}, n - k)) \leq mn + \max_k (mk + m(n - k)) = 2mn$$

□

2.3 Поиск в ориентированном графе.

Алгоритм 11. Поиск в ширину (BFS). В каждый момент времени вершина графа может быть помечена или не помечена. Если вершина уже помечена, значит алгоритм нашёл путь из корня в неё. Кроме пометок на вершинах, алгоритм хранит очередь, в которой находятся все те помеченные вершины, для которых ещё не обработаны исходящие дуги. Таким образом, в каждый момент времени вершина может быть не помеченной, помечанной и обработанной, и помечанной и необработанной. Идём из корня и последовательно отмечаем и заносим в очередь тех, к кому пришли.

Утверждение 2. В каждый момент времени очередь состоит из некоторых вершин, находящихся на расстоянии l от s , вслед за которыми идут некоторые вершины, находящиеся на расстоянии $l + 1$ от s . При этом все вершины на расстоянии, меньшем, чем l , уже обработаны, ровно как и все вершины на расстоянии l , не вошедшие в очередь. Из вершин на расстоянии $l + 1$ в очереди есть ровно все потомки обработанных вершин.

Утверждение 3.

- алгоритм помечает вершину v тогда и только тогда, когда есть путь из s в v ;
- если алгоритм находит v по дуге (u, v) , то один из кратчайших путей из s в v идёт через u ;
- все пройденные дуги (u, v) образуют дерево.

Алгоритм 12. Поиск в глубину (DFS). Идём в глубину до конца, отмечаем вершины в чёрный, если из них начали идти вниз, серым, если они нам просто встретились на пути. После того, как дошли до конца, идём вверх до первой вершины. Время работы: $O(|V| + |E|)$.

Задача 3. Топологическая сортировка. Нужно найти остовные деревья в орграфе (естественно, он должен быть ациклическим). Решается через DFS из следующих утверждений.

Утверждение 4. Граф ациклический тогда и только тогда, когда при поиске в глубину никогда не рассматривается дуга, ведущая в вершину, находящуюся в стеке возврата (дуга из серой в серую).

Утверждение 5. Если в ациклическом графе есть дуга (u, v) , то время завершения v меньше, чем время завершения u .

3 Лекция 4.

3.1 Окончание поисков в орграфе.

Задача 4. Найти в данном графе его компоненты сильной связности.

Алгоритм 13. Алгоритм Косараджу-Шарира. Спервая запускается поиск в глубину для G , а затем запускается поиск в глубину для обращённого графа G^R , в котором направления всех дуг изменены на обратные (однако, компоненты связности те же). При поиске в глубину в обращённом графе, во внешнем цикле вершины рассматриваются в порядке их завершения при первом поиске в глубину, от конца к началу. После этого оказывается, что каждый запуск процедуры *DFS* во внешнем цикле будет находить очередной сильно связный компонент исходного графа. Время работы: $O(|V| + |E|)$, корректность обосновывается следующим:

Утверждение 6. Пусть в графе G есть сильно связанные компоненты C и D , и есть дуга $(u, v) \in E$ из C в D . Тогда при поиске в глубину в графе G самое позднее время завершения вершины в C превосходит таковое в D .

Доказательство. Рассматриваются два случая: самое ранне обнаружение в C меньше, чем в D , тогда рассматриваем первую обнаруженную вершину $x \in C$, у всех вершин из D время окончания меньше, чем у неё. Если же это не так, то рассмотрим самую раннюю $y \in D$. Все остальные из этой компоненты будут обнаружены на рекурсивных вызовах, а из C на этом этапе не обнаружатся, поэтому время окончания всех вершин из C больше. \square

Теорема 4. Вершины каждого дерева, найденного алгоритмом Косараджу-Шарира при втором поиске в глубину - это и есть сильно связанные компоненты исходного графа.

Доказательство. Индукция по количеству найденных компонент связности. Надо доказать, что если первые k найденных компонент связности действительно таковы, то и следующая также обладает этим свойством. \square

3.2 Поиск в алгоритме с весами.

Задача 5. Пусть в орграфе для каждой дуги задан вес. Нужно найти пуи наименьшего веса из данной вершины $s \in V$ во все вершины графа.

Алгоритм 14. Алгоритм Беллмана-Форда. Для каждой вершины вычисляются значения d_v - наименьший вес пути из s в v и π_v - предыдущая вершина на пути наименьшего веса из s в v . Изначально полагается, что $d_v = \infty$ и $\pi_v = \text{NULL}$ для всех вершин, и $d_s = 0$. Далее алгоритм постепенно находит пути меньшего веса в другие вершины, запоминая веса лучших из найденных путей в этих переменных. Значения уменьшаются с помощью элементарной операции улучшения пути, используя некоторую дугу $(u, v) \in E$. Если $d_u + w_{u,v} < d_v$, то $d_v = d_u + w_{u,v}$, а $\pi_v = u$. Эта операция применяется, пока можно что-то улучшить. Как будет показано, для этого достаточно рассмотреть все дуги $|V| - 1$ раз.

Утверждение 7. После i -ой итерации внешнего цикла алгоритм Беллмана-Форда находит все пути наименьшего веса длины не более чем i .

Доказательство. Индукция по i . \square

Теорема 5. Алгоритм Беллмана-Форда за $|V| \cdot |E|$ шагов или правильно вычисляет пути наименьшего веса из вершины s во все вершины, или сообщает о наличии достижимого цикла отрицательного веса.

Доказательство. Рассмотрим систему после $|V|-1$ итераций, и согласно утв. 7, найдены все пути наименьшего веса, состоящие не более чем из $|V|-1$ дуг. Если какой-то путь ещё можно улучшить, то в этом пути какая-то вершина встретилась дважды, следовательно, найден цикл отрицательного веса. Если же ничего нельзя улучшить, то любой достижимый цикл будет иметь неотрицательный вес, так как для последовательный вершин цикла можно записать не условие улучшения и сложить по все парам. \square

Алгоритм 15. Алгоритм Дейкстры. Этот алгоритм решает ту же задачу, однако на каждом шаге находит очередную вершину u , путь наименьшего веса в которую уже известен. Тогда алгоритм рассматривает все дуги, исходящие из вершины u . Это позволяет ему рассматривать каждую дугу графа лишь однажды.

Теорема 6. Алгоритм Дейкстры работает правильно.

Доказательство. Индукцией по длине вычисления доказываем, что для всякой вершины $u \notin Q$, путь наименьшего пути уже построен. Для этого рассматриваются две вершины на кратчайшем пути - одна в среди рассмотренных, другая - среди не рассмотренных (Q), и для последней проводятся вычисления, связанные с длиной самого короткого пути для неё. \square

Примечание 2. Для представления множества Q алгоритм использует особую структуру данных: очередь с приоритетами. Каждый элемент Q находится там вместе со своим текущим значением d_v . Также заданы операции:

- $insert(x)$ - вставить новый элемент;
- $min()$ - выдать минимальный элемент;
- $extract_{min}()$ - выдать минимальный и удалить;
- $decrease(x, k)$ - изменить значение элемента $x \in Q$ на k , если k меньше текущего значения x_i .

Скорость работы алгоритма: $|V|$ раз $extract_{min}$ и $|E|$ раз $decrease$, поскольку каждая дуга обрабатывается лишь однажды.

3.3 Окончание поиска с весами.

Сложность алгоритма Дейкстры зависит от того, как реализована очередь с приоритетами. Тупая реализация: хранить массив x_v , индексированный по $v \in V$. Тогда $decrease$ работает за $O(1)$, но $extract_{min}$ требует время $|V|$. Отсюда общее время работы - $|V|^2 + |E| = O(|V|^2)$.

Алгоритм 16. Улучшение Дейкстры кучей. Используется куча, в которой значение в каждой внутренней вершине не больше, чем значение в любом из её потомков ($min\text{-heap}$). Операции выполняются следующими:

- $insert(x)$ - вставить новый элемент (он становится листом), после чего дать ему всплыть до его законного места;

- $\min()$ - выдать минимальный элемент (просто вернуть x_1);
- $\text{extract}_{\min}()$ - переместить x_n в x_1 , убрав его в конце, а затем запустить исправление кучи из корня, то есть, $\text{heapify}(1)$;
- $\text{decrease}(i, k)$ - изменить значение элемента x_i на k , после чего дать элементу x_i всплыть наверх, пока возможно.

Примечание 3. Дать всплыть - значит, если x_i меньше своего родителя, то он обменивается так до тех пор, пока не займёт положенное место.

3.4 Нахождение минимального остовного дерева.

Задача 6. Дан неориентированный связный граф, рёбрам которого сопоставлены числа - веса. Нужно найти одно из остовных деревьев с наименьшим весом.

Алгоритм 17. Общий принцип действия алгоритмов. Одно за другим присоединяются рёбра к поддереву какого-то минимального, чтобы опять получилось поддерево минимального.

Определение 2. Сечение графа - разбиение множества вершин на два дизъюнктных множества. Ребро пересекает сечение, если один из его концов лежит в одной части, а другой - во второй.

Утверждение 8. Пусть T - одно из минимальных остовных деревьев графа, а F - его подмножество. Пусть также имеется какое-то сечение, что никакое ребро из F его не пересекает. Тогда ребро с наименьшим весом, пересекающее это сечение, принадлежит некоторому минимальному остовному дереву T' , которое также содержит F .

Доказательство. Если (u, v) - такое ребро, входит в T , то нам подойдёт T . Если его там нет, то рассмотрим цикл, который с ним получается и поменяем его на другое ребро из цикла, которое пересекает сечение. \square

Алгоритм 18. Алгоритм Прима. Начинаем с произвольной вершины и на каждом шаге добавляем ребро из одной из уже имеющихся вершин в некоторую незадействованную (естественно, из всевозможных рёбер выбирается то, у которого минимальный вес).

Время работы: $|E| \log |V|$, так как выполняется $|V|$ операций внешнего цикла. Операции над очередью с приоритетами - $\log |V|$. Внутренний цикл по v : за всё время работы алгоритма каждое ребро рассматривается дважды: с одного конца и с другого. Поэтому тело цикла в общей сложности выполняется $2|E|$ раз, и всякий раз выполняется операция над очередью с приоритетами за время $O(\log |V|)$.

Примечание 4. Структура данных алгоритма - очередь с приоритетами, в которой хранятся все вершины, ещё не попавшие в дерево. Значение d_v каждой вершины v - это наименьший вес ребра, соединяющий её с деревом. Также для v запоминается вершина π_v , через которую v соединена с деревом ребром наименьшего веса.

Всякий раз, когда в дерево добавляется новая вершина u , для любой вершины v не из дерева может оказаться, что ребро (u, v) легче, чем ранее известное ребро наименьшего веса (π_v, v) , соединяющее её с деревом, и это так, если $w_{u,v} < d_v$. В этом случае значения d_v и π_v обновляются, переключая v на соединение с деревом через u .

Алгоритм 19. Алгоритм Крускала. Текущее подмножество остовного дерева - лес, соединяем компоненты самыми лёгкими путями.

Теорема 7. На каждом шаге работы алгоритма Крускала переменная T содержит подмножество одного из остовных деревьев минимального веса.

Доказательство. Индукция по длине вычисления. □

Утверждение 9. При использовании леса непересекающихся множеств алгоритм работает за время $|E| \log |E|$ - и, стало быть, $|E| \log |V|$.

Доказательство. Сортировка займёт время $|E| \log |E|$. Далее алгоритм выполняет $3|E|$ операций над структурой данных, каждая из которых, при реализации через лес непересекающихся множеств, выполняется за время $O(\log n)$, и в общей сложности получается $|E| \log |V|$, что уже быстрее сортировки. □

3.5 Структура данных для непересекающихся множеств.

Задача 7. Абстрактная структура данных для представителя разбиения множества на непересекающиеся подмножества. У каждого множества есть выделенный элемент - *представитель*. Заданы операции:

- $make_set(x)$ - создать одноэлементное множество;
- $find_set(x)$ - найти представитель множества, содержащего данный элемент;
- $set_union(x, y)$ - объединение двух множеств.

Перед нами стоит задача: как эффективно реализовать эту абстрактную структуру данных?

Алгоритм 20. Лес непересекающихся множеств. У нас есть лес, если одноэлементно - просто вершинка, если представитель - корень дерева, если объединяем, то корень одного будет указывать на корень другого.

Примечание 5. Используются следующие улучшения:

В каждой вершине хранится её условная сложность - *ранг*. Процедура $make_set(x)$ устанавливает ранг в нуль. При объединении двух деревьев корень дерева меньшего ранга станет указывать на корень дерева большего ранга. Если оба корня имеют одинаковый ранг, то их объединение получит ранг, больший на единицу.

При выполнении каждой операции поиска все встреченные на пути вершины пренаправляются в корень для ускорения последующих операций цикла.

Утверждение 10. Пусть всего элементов - n , и к ним применяется m операций. Тогда, с применением первого улучшения, они выполняются за время $O(m \log n)$.

Доказательство. Смотрим на дерево, через логарифмы вычисляем. □

Теорема 8. (Тарьян.) Пусть всего элементов - n , и над ними выполняется m операций. Тогда, с применением первого и второго улучшения, эти операции выполняются за совокупное время $m\alpha(n)$, где $\alpha(n)$ - обратная функция к $A_n(n)$, а $A(k, n) = A_k(n)$ - функция Аккермана.

Задача 8. Дан неориентированный граф - возможно, несвязный. Надо определить его компоненты связности и научиться быстро отвечать на вопрос о том, лежат ли две данные вершины в одной компоненте.

Алгоритм 21. Алгоритм пишется через представление компонент связности в виде структуры данных для непересекающихся множеств.

4 Лекция 5.

4.1 Пути между всеми парами вершин.

Задача 9. Дан ориентированный граф $G = (V, E)$, где $V = \{1, \dots, n\}$, а $E \subseteq V \times V$. Ставится задача проверить существование пути из каждой вершины в каждую - то есть: для каждой пары вершин (i, j) определить, есть ли путь из i в j . После решения мы получим *транзитивное замыкание* графа.

Алгоритм 22. Очевидный алгоритм. Перебираем все пары вершин (i, j) и для каждой пары все промежуточные вершины k . Если она соединена с обоими, то добавляем ребро (i, j) . Делаем так пока можем.

Утверждение 11. После каждой t -ой итерации внешнего цикла будут просчитаны все пути длины не более 2^t , и для каждого из них в граф будет добавлена дуга. Откуда время работы - не более, чем $n^3 \log_2 n$ шагов.

Алгоритм 23. Алгоритм Варшалла. Можно построить транзитивное замыкание за время n^3 , обойдясь без внешнего цикла (без пока можем). Для этого достаточно пометить циклы местами и теперь рассматривать каждую вершину как смежную, соединяя все пары, смежные с ней.

Утверждение 12. После k -ой итерации внешнего цикла найдены все пути, проходящие только через промежуточные вершины из множества $\{1, \dots, k\}$.

Доказательство. Доказательство индукцией по количеству итераций. □

Задача 10. Пусть граф теперь не только ориентированный, но и рёбра имеют вес. Нужно найти кратчайшие пути между вершинами.

Алгоритм 24. Алгоритм Флойда-Варшалла. То же самое, что и алгоритм варшала, но мы ещё и храним вес рёбер.

Примечание 6. О достижимости можно говорить в рамках матриц смежности. Умножение рассматривается как

$$c_{i,j} = \bigvee_{k=1}^l a_{i,k} \wedge b_{k,l}.$$

Тогда если A - матрица смежности, то A^l - матрица достижимости по путям длины ровно l .

4.2 Быстрое умножение матриц.

Ну, какие-то способы усложнения - ебанина какая-то, писать не буду.

Теорема 9. (Основная теорема о времени работы рекурсивных алгоритмов). Пусть задача размера n решается путём разбиения её на a подзадач того же типа, размера $\frac{n}{b}$ каждая, и процесс разбиения, а также соединения результатов занимает время $O(n^c)$.

$$T(n) = aT\left(\frac{n}{b}\right) + O(n^c)$$

Тогда время работы алгоритма оценивается как

- Если $c < \log_b a$, то $O(n^{\log_b a})$;
- Если $c = \log_b a$, то $O(n^c \log_b n)$;

- Если $c > \log_b a$, то (n^c) .

Доказательство. В дереве рекурсии $\lceil \log_b n \rceil$ уровней, на нулевом - одна подзадача размера n , расчёты занимают время n^c , на первом - a подзадач размера $\lceil \frac{n}{b} \rceil$, расчёты для каждой занимают время $\lceil \frac{n}{b} \rceil^c$, и так далее, на уровне i - a^i подзадач размера $\lceil \frac{n}{b^i} \rceil$ каждая, расчёты каждой занимают время $\lceil \frac{n}{b^i} \rceil^c$. Общее время работы - сумма по всем уровням рекурсии.

$$T(n) = \sum_{i=0}^{\lceil \log_b n \rceil} a^i \left(\frac{n}{b^i} \right)^c = n^c \sum_{i=0}^{\lceil \log_b n \rceil} \left(\frac{a}{b^c} \right)^i$$

Таким образом, если $a = b^c$, то сумма геометрической прогрессии - \log_b^n , и отсюда $O(n^c \log_b n)$. Если $a < b^c$, то сумма ограничена сверху константой, отсюда $O(n^c)$. И если $a > b^c$, то это возрастающая геометрическая прогрессия, и применив формулу для вычисления получим $O(n^{\log_b a})$. \square

Задача 11. Умножение булевых матриц.

Алгоритм 25. Можно записать матрицу в виде кучи нулей и единиц, после чего перемножить в кольце вычетов по модулю $n+1$ или по любому удобному модулю, превосходящему n . Тогда вместо дизъюнкции конъюнкций будет вычислена сумма произведений по модулю $n+1$ (в точности равная количеству истинных конъюнкций в дизъюнкции конъюнкций, и потому она лежит в диапазоне от 0 до n , откуда следует, что по модулю $n+1$ она вычислится точно). Чтобы узнать значение дизъюнкции конъюнкций, достаточно будет проверить вычисленную сумму на равенство нулю.

Алгоритм 26. Метод четырёх русских. Чисто комбинаторный метод умножения булевых матриц за время $O(n^3)$. Пусть A и B - две булевых матрицы размера $n \times n$, цель - вычислить их произведение $C = AB$. Пусть k - число, намного меньшее, чем $\log_2 n$, и пусть n делится на k (если не делится - немного увеличим n до делимости). Каждая строка матрицы A делится на $\frac{n}{k}$ векторов размера $1 \times k$, называемых кусочками. Кусочки в i -ой строке обозначаются через $A_{i,1}, \dots, A_{i,\frac{n}{k}} \in \mathbb{B}^{1 \times k}$. Матрица B разделяется на $\frac{n}{k}$ подматриц размера $k \times n$, называемых полосами и обозначаемыми через $B_1, \dots, B_{\frac{n}{k}} \in \mathbb{B}^{k \times n}$. Тогда всякая i -ая строка C , обозначаемая через $C_i \in \mathbb{B}^{1 \times n}$, представима в виде следующей поэлементной дизъюнкции строк:

$$C_i = \bigvee_{r=1}^{\frac{n}{k}} A_{i,r} B_r.$$

Произведение кусочка $A_{i,r}$ на соответствующую полосу B_r - это строка размера $1 \times n$, и знак дизъюнкции в формуле для C_i означает поразрядную дизъюнкции таких строк,

Примечание 7. Казалось бы, никакого принципиального улучшения не получилось, однако с программистской точки зрения, это в чём-то удобнее.

Главная идея метода же состоит в том, что так как k невелико, возможных кусочков всего 2^k , и кусочки будут часто повторяться. Таким образом, если запомнить произведения таких кусочков на полосы B , в плане вычислений это займёт много меньше времени. Общее время работы будет $O\left(\frac{n^3}{\log n}\right)$.

5 Лекция 6.

5.1 Двоичные деревья поиска

Речь сейчас пойдёт о структурах данных для представления *множества* различными способами, с разными временем выполнения операций (строго говоря, мы будем говорить о *мультимножестве*, так как может быть несколько элементов с одинаковым значением). Элементы могут быть любого вида, на них определено отношение порядка " \leq ". Операции: найти элемент с данным значением; найти наибольший (наименьший) элемент; найти элемент, предшествующий или следующий за данным; вставить / удалить элемент.

Алгоритм 27. Бинарное дерево поиска. Представление множества в виде дерева так, что вершина содержит одно значение и три указателя: на предка (если корень, то *NULL*), на левое поддерево и на правое (если их нет, то *NULL*). При этом все вершины в её левом поддереве содержат не меньшие значения, а все вершины в правом поддереве - не большие. Понятно, как на таком дереве задаются требуемые операции.

Примечание 8. Сложность каждой операции - не более, чем высота дерева. Если дерево сбалансированное, то есть, все пути примерно одинаковой длины, то все операции выполняются за логарифмическое время.

Задача 12. Даже если дерево сбалансированное, то операции над ним могут привести его к дисбалансу, надо как-то сложить структуру, чтобы избежать критических изменений.

Алгоритм 28. AVL-деревья. Усложнённая разновидность двоичных деревьев поиска. Такую разновидность называют почти сбалансированной, то есть, высота деревьев-потомков одной вершины отличается не более, чем на 1. Этого ограничения достаточно, чтобы дерево имело логарифмическую высоту.

Утверждение 13. AVL-дерево высоты h содержит не менее, чем $F_{h+3} - 1$ вершин, где F_n - n -ое число Фибоначчи.

Доказательство. Индукция по h . База очевидна, переход таков. Если дерево имело высоту h , то один из его потомков имеет высоту $h - 1$, а другой - высоту не менее, чем $h - 2$, то есть, не менее чем $F_{h+2} - 1$ и $F_{h+1} - 1$ вершин соответственно. Всего, с учётом корня и по определению чисел Фибоначчи, как раз и получится то, что нужно. \square

Утверждение 14. Высота AVL-дерева с n вершинами не превосходит $\log_{\varphi} n$, где φ - золотое сечение.

Примечание 9. Естественно, просто так AVL-дерево не получить, поэтому придумана операция *вращение*, которая состоит в том, что если у нас есть корень y , из которого направо идёт поддерево t_3 , а налево - поддерево с корнем x и поддеревьями t_1 и t_2 налево и направо соответственно. Тогда мы подвешиваем дерево за x и перебрасываем его поддерево t_2 в левое поддерево y . При помощи такой операции можно исправить балансировку при различных операциях, подробно расписывать я это пока что не буду.

Задача 13. Двоичные деревья рассчитаны на то, чтобы храниться в оперативной памяти компьютера, а для того, чтобы хранить деревья во внешней, медленной памяти, требуется некая адаптация.

5.2 В-деревья.

Алгоритм 29. В-деревья. Адаптация деревьев поиска для хранения во внешней памяти. Основная мысль - использовать вершины большой степени - с тем, чтобы каждая вершина занимала один блок, а высота дерева уменьшилась бы.

Пусть вершины в двоичном дереве - это 2-вершины, поскольку у каждой из них 2 потомка и 1 значение, по которому эти потомки разделяются. У m -вершины - m потомков (деревья t_1, \dots, t_m - возможно, пустые), и в ней находится $m - 1$ значение: x_1, \dots, x_{m-1} , где значения упорядочены по неубыванию. Все значения в каждом поддереве t_i , больше или равны x_i , и меньше или равны x_{i+1} .

В В-дереве могут одновременно соежаться вершины различных степеней: выбирается некоторое число $k \geq 2$, после чего корень может иметь степень от 0 до $2k$, а все остальные вершины - любые степени от k до $2k$. При этом, дерево сбалансировано.

Примечание 10. То, как производятся требуемые операции - лучше один раз увидеть и немного прочесть в конспекте А. С. Охотина, чем читать мои заметки.

Алгоритм 30. В-дерево для $k = 2$ называется 2-3-4 деревом, и такое дерево очень удобно хранить в оперативной памяти.

Алгоритм 31. Красно-чёрное дерево - это представление 2-3-4 дерева в виде двоичного дерева, в котором каждая вершина представлена в виде одной или нескольких связанных между собою двоичных вершин, каждая из которых покрашена в один из двух цветом. Каждая 2-вершина остаётся собой и считается чёрной. Каждая 3-вершина разбивается на две двоичных: чёрную - корень поддерева, и красную вершину - правого или левого потомка чёрной. Далее - к этой паре вершин, так же как и к исходной 3-вершине, присоединяются три поддерева. Наконец, каждая 4-вершина разбивается на три двоичных - чёрную и два красных потомка, к которым присоединены четыре поддерева.

6 Лекция 7.

6.1 Полиномиальное хэширование.

Задача 14. Дана длинная строка $w = a_1 \dots a_n$ и короткая искомая строка $x = b_1 \dots b_m$. Требуется найти все вхождения x в w в качестве подстроки.

Алгоритм 32. Хэш-функции. Каждой строке ставится в соответствие некоторое число $w \mapsto h(w)$, и это число хранится вместе со строкой. Тогда, если нужно сравнить две строки, то мы сравниваем сперва соответствующие им числа. Если числа разные, то и строки точно разные. Если соответствующие числа одинаковые, то надо сравнить строки и получить итоговый результат. Самое тупое - сложить коды всех символов, но тогда распределение неравномерное, и вообще это не очень кайфовый вариант. Намного лучше полиномиальное хэширование. Берём некоторое основание степени p , пусть $w = a_1 \dots a_l$ - строка, тогда используется сумма $\sum_{i=1}^l a_i \cdot p^{l-i}$, взятая по некоторому модулю M .

Алгоритм 33. Алгоритм Рабина-Карна. Алгоритм поиска подстроки $x = b_1 \dots b_m$ в строке $w = a_1 \dots a_n$, основанный на полиномиальном хэшировании степени $m - 1$: сперва вычисляется значение хэш-функции для искомой строки, а затем для всех m -символьных подстрок данного текста последовательно вычисляется значение их хэш-функции. Когда значение для подстроки и искомой строки совпадают, производится прямое сравнение. Значение хэш-функции для искомой строки: $X = \sum_{i=1}^m b_i \cdot p^{m-i}$ по модулю q , а

значение хэш-функции для подстроки со смещением s : $W_s = \sum_{i=1}^m a_{s+i} \cdot p^{m-1}$ по модулю q . Алгоритм сперва вычисляет X , а затем - последовательно все W_s . Из соотношения $W_{s+1} = pW_s - a_{s+1}p^m + a_{s+m+1} \pmod{q}$. Сложность: $\Theta(m)$ - на подготовку, и затем в худшем случае $O(mn)$, если хэ-функция выдаст одинаковые значения для всех подстрок. В среднем случае получается время работы $\Theta(n)$.

Задача 15. (Наибольшая общая подстрока). Даны две строки: $u = a_1 \dots a_m$ и $v = b_1 \dots b_n$. Надо найти самую длинную их общую подстроку x .

Алгоритм 34. Тупое решение: для каждой пары позиций найти самую длинную подстроку, время $((m+n)^3)$.

Алгоритм 35. Улучшение через полиномиальное хэширование - научиться отвечать на вопрос "Есть ли общая подстрока длины l ?" за время $O((m+n)\log(m+n))$. Для этого находятся значения хэш-функции для всех подстрок и длины l , размещаются в двоичном дереве поиска - время $m \log m$. То же самое - для v за время $n \log n$.

Алгоритм 36. Улучшение через двоичный поиск по l . Его можно записать в виде рекурсивной процедуры, отвечающей на вопрос "Найти длину наибольшей общей подстроки, если известно, что её длина не меньше, чем l_1 , и строго больше, чем l_2 ?". Процедура $f(l_1, l_2)$: если $l_2 - l_1 = 1$, то возвращаем l_1 , $k = \lfloor \frac{l_1 + l_2}{2} \rfloor$. Теперь, если есть общая подстрока длины k , то возвращаем $f(k, l_2)$, иначе возвращаем $f(l_1, k)$.

Задача 16. Найти в строке самую длинную подстроку-палиндром.

Алгоритм 37. Самый тупой алгоритм: искать вокруг каждого символа на длину вплоть до $\frac{n}{2}$, время $O(n^2)$.

Алгоритм 38. Улучшение основано на решении подзадачи: "Есть ли подстрока-палиндром данной длины l ?". После того, как будет построен алгоритм для решения подзадачи, останется использовать двоичный поиск по l .

Алгоритм 39. Решение через небольшую переформулировку: разворачиваем строку и для таковой и оригинальной рассматриваем хэш-функции для всех подстрок длины l . После чего, естественно, надо выполнить посимвольную проверку. Двоичный поиск по l даст время $O(n \log n)$.

6.2 Алгоритм Кнута-Морриса-Пратта.

Алгоритм 40. Строка w читается слева направо, и после чтения i символов w он помнит длину j самого длинного префикса строки x , на которой заканчивается $a_1 \dots a_i$ - то есть, наибольшее число j , для которого верно $a_{i-j+1} \dots a_i = b_1 \dots b_j$.

При чтении очередного символа a_{i+1} алгоритм сравнивает его с b_{j+1} , и если эти символы равны, это значит, что префикс x длины j продлевается на один символ. Если же $a_{i+1} \neq b_{j+1}$, то самый длинный префикс x , на который заканчивается строка $a_1 \dots a_{i+1}$, будет содержать строго меньше, чем $j+1$ символов. Чтобы найти длину этого префикса, нужно рассмотреть следующий по длине префикс x на который заканчивается $a_1 \dots a_i$, и пытаться продлевать уже его.

Чтобы иметь возможность находить длину такого префикса, алгоритм заранее строит по данной искомой строке x определённую структуру данных - префиксную сумму.

Определение 3. Префиксная сумма для строки $x = b_1 \dots b_m$ - это функция $\pi\{1, \dots, m\} \rightarrow \{0, \dots, m-1\}$, которая для всякого префикса $b_1 \dots b_i$ строки x выдаёт длину наибольшего суффикса подстроки $b_1 \dots b_i$, который также является префиксом x .

$$\pi(i) = \max\{j | j < i, b_1 \dots b_j = b_{i-j+1} \dots b_i\}$$

Примечание 11. Повторно применяя функцию π можно получить все префиксы x , являющиеся суффиксами $b_1 \dots b_i$, и алгоритм будет их перебирать, пока не найдёт такой, который можно продолжить следующим символом текста.

Утверждение 15. Имея готовую таблицу значений функции π , алгоритм КМР работает за время $\Theta(n)$.

Примечание 12. Осталось научиться быстро строить значение префиксной функции.

6.3 Поиск с помощью конечных автоматов.

Определение 4. Детермированный конечный автомат (DFA) - пятёрка $A = (\Sigma, Q, q_0, \delta, F)$, со следующим значением компонентов:

- Σ - алфавит;
- Q - конечное множество состояний;
- $q_0 \in Q$ - начальное состояние;
- $\delta : Q \times \Sigma \rightarrow Q$ - функция переходов. Если автомат находится в состоянии $q \in Q$ и читает символ $a \in \Sigma$, то его следующее состояние - $\delta(q, a)$;
- $F \subseteq Q$ - множество *принимаящих состояний*.

Для всякой входной строки $w = a_1 \dots a_l$, где $l \geq 0$ и $a_1, \dots, a_l \in \Sigma$, вычисление - последовательность состояний $p_0 \dots, p_{l-1}, p_l$, где $p_0 = q_0$, и всякое следующее состояние p_i , где $i \in \{1, \dots, l\}$, однозначно определено как $p_i = \delta(p_{i-1}, a_i)$.

$$p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_1} \dots \xrightarrow{a_l} p_l$$

Строка *принимается*, если последнее состояние p_l принадлежит множеству F - иначе *отвергается*. Множество строк, *распознаваемое* автоматом, обозначаемое через $L(A)$ - состоит из всех строк, которые он принимает.

Теорема 10. Для всякой строки $x = b_1 \dots b_m$ над алфавитом Σ существует DFA с $m+1$ состояниями, распознающий множество Σ^*x всех строк, заканчивающихся на x , и этот DFA можно построить за время $\Theta(|\Sigma| \cdot m)$.

Доказательство. Пока что записывать не буду. □

6.4 Представления множеств строк.

Алгоритм 41. Лексикографический порядок. Поэлементное сравнение слева направо.

Задача 17. Как хранить строки?

Алгоритм 42. В виде списка или двоичного дерева поиска, однако это всё хуйня.

Алгоритм 43. Префиксное дерево - структура для хранения множеств строк. Корневое дерево с дугами, помеченными символами алфавита. Дуги, исходящие из всякой вершины, должны быть помечены различными символами. Каждая вершина соответствует некоторой строке, которая хранится неявно в виде последовательности дуг на пути к ней. Также всякая вершина хранит один бит информации, принадлежит ли эта строка множеству; вместе с битом можно хранить любые дополнительные значения, сопоставленные в этой строке. Корень соответствует пустой строке. Если вершина соответствует строке w , и исходящая из неё дуга помечена символом $a \in \Sigma$, то эта дуга идёт в вершину, соответствующую строке wa .

Все операции с любой данной строкой длины l (поиск, вставка, удаление) выполняются за время $O(l)$, не зависящее от числа элементов в хранимом множестве.

Компактное префиксное дерево - объединяем несколько дуг в одну, если они идут одинаково параллельно.

Утверждение 16. В компактном префиксном дереве для множества из n элементов не более, чем n внутренних вершин.

6.5 Поиск нескольких строк одновременно: алгоритм Ахо-Корасик.

Задача 18. Пусть для данного текста нужно найти все вхождения не одной строки, а всех строк из конечного множества $K = \{x_1, \dots, x_k\}$.

Алгоритм 44. Алгоритм Ахо-Корасик. Данный алгоритм будет помнить самый длинный только что прочитанный префикс одной из строк x_1, \dots, x_k , а точнее сказать, самую длинную такую строку u , что u - префикс какой-то строки из K , и алгоритм только что прочитал u .

Примечание 13. Для решения обобщается префиксная сумма и построение функции π (работающее за линейное время), но пока что упущу подробности.

6.6 Суффиксные деревья.

Алгоритм 45. Суффиксное дерево - структура данных, строящаяся по данной строке и обеспечивающая эффективный поиск подстрок в этой строке. Понятно, что это такое.

7 Лекция 9.

7.1 Сжатие данных, основанное на повторении строк.

Задача 19. Нужно научиться использовать часто повторяющиеся подстроки, чтобы хранить их отдельно в памяти и эффективно использовать при хранении целого текста.

Алгоритм 46. Метод Лемпеля-Зива LZ77. Алгоритм сжатия данных, использующий повторяемость подстрок. Основная идея в том, что каждый раз, когда ранее прочитанная подстрока встречается повторно, алгоритм пишет вместо неё ссылку на предыдущее вхождение. На каждом шаге выводится тройка (d, l, a) , что означает: сперва повторяется подстрока длины l , ранее встречавшаяся d символов назад, а потом идёт символ a . Реализуется самым простым образом через суффиксное дерево, в котором на каждом шаге читаются следующие символы входной строки, пока не находится самая длинная раньше встречавшаяся подстрока. Обычно используется небольшое улучшение - "скользящее окно", то есть, подстроки ищутся (тем же суффиксным деревом, скорее всего) в окне из последних m элементов, а остальные забываются.

Теорема 11. Жадный LZ77 оптимален.

Алгоритм 47. Метод Лемпеля-Зива LZ78. Тут уже по мере чтения строки строится словарь из часто встречающихся строк в виде префиксного дерева. При декодировании строится в точности этот же словарь. В сжатом представлении строки словарь не хранится.

В начале в словаре содержится только один элемент под номером нуль: $T_0 = \varepsilon$, иными словами, префиксное дерево состоит из корня, помеченного номером 0. На каждом шаге читается самая длинная строка $T_j = v$, и уже имеющаяся в словаре, и выводится её код j ; также читается и выводится следующий символ a . При этом в словарь добавляется новая строка va - конкатенация только что прочитанной со следующим входным символом.

На префиксном дереве это выглядит так: после вывода очередной пары (j, a) алгоритм переходит в корень префиксного дерева, и дальше читает столько входных символов, сколько возможно. Когда очередной символ прочитать нельзя, создаётся новый лист, при этом выводится номер предыдущей вершины и прочитанный символ.

При декодировании строится то же самое префиксное дерево. Алгоритму при этом потребуется находить строку по номеру в таблице: можно, например, завести для этого отдельный массив строк, или же читать в префиксном дереве путь из вершины в корень. Прочитав очередную пару (j, a) алгоритм выводит строку T_j и символ a , после чего сразу перескакивает в вершину j и присоединяет к ней новый лист, с переходом по символу a .

Алгоритм 48. Преобразование Берроуза-Вилера. Перутся все циклические сдвиги строки $w = a_1 \dots a_n$. Они сортируются лексикографически, получается таблица T размера $n \times n$, где в каждой строчке - один из циклических сдвигов. наконец, берётся последовательность последних символов $b_1 \dots b_n$ в таблице T и номер t строки w в этой таблице.

Если в исходной строке какая-то подстрока ai часто повторялась, то есть много циклических сдвигов, начинающихся с i и заканчивающихся на a , и после сортировки они окажутся рядом. Поэтому в преобразованной строке будут длинные последовательности повторяющихся символов. Использование в качестве метода сжатия данных: преобразовать, а потом применять другие методы к преобразованной строке.

Утверждается, что по последовательности b_i , и t восстанавливается исходная w .

Утверждение 17. BWT можно вычислить за линейное время.

Утверждение 18. Обратное к BWT преобразование можно вычислить за линейное время.