



- Практикум
- Задача о рюкзаке
- Динамическое программирование. М...
- Динамическое программирование. А...
- Двумерная динамика
- Задача о рюкзаке. Динамическое реш...
- Наибольшая общая подпоследователь...
- Разбор задачи: переход от теории к п...
- Резюме
- Оцените пройденные уроки
- Резюме
- Оцените пройденные уроки
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 7
- Финальное задание спринта 7
- Резюме
- Оцените задание и сопровождение
- спринт 8
- Алгоритмы на строках
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы

# Сравнение строк

Алла: А какая самая простая задача при работе со строками?

Рита: Наверное, сравнение двух строк `first` и `second` на равенство.

Алла: Ну это же совсем легко. Строки равны, если у них одинаковая длина и все символы на соответствующих позициях — совпадают. То есть, если хотя бы один символ не совпадает, строки не равны.

```

функция compare(first, second):
    если длина(first) ≠ длина(second):
        вернуть False

    length = длина(first)
    для i от 0 до (length - 1):
        если first[i] ≠ second[i]:
            вернуть False

    вернуть True
    
```

Тимофей: Всё верно. И это сравнение работает за  $O(n)$ , где  $n$  — длина строки.

Алла: Но длина какой из строк имеется в виду, `first` или `second`?

Тимофей: Если эти строки имеют разную длину, то мы сразу, за  $O(1)$  можем сказать, что строки не равны. А если строки имеют одинаковую длину, то не важно, какую из этих длин брать. Поэтому  $n$  — длина любой из строк.

## Регистронезависимое сравнение

Гоша: Я вижу одну проблему. Люди часто пишут одни и те же слова то большими буквами, то маленькими. Но ведь получается, если у меня в телефонной книге сохранён контакт «Алла», а я напишу в поиске «АЛЛА», то ничего не найдётся.

Рита: Немного поправлю тебя: «Маленькие» буквы и «большие», а точнее строчные и прописные, называются в программировании буквами в нижнем и верхнем регистрах. Про слова, которые сочетают в себе оба этих типа букв, (как, например, «Алла»), говорят, что они записаны в смешанном регистре.

Чтобы сравнить две строки без учёта регистра, следует сначала привести их к единой форме. Например, заменить все буквы в верхнем регистре на аналогичные буквы в нижнем регистре: `compare(lowercase(first), lowercase(second))`.

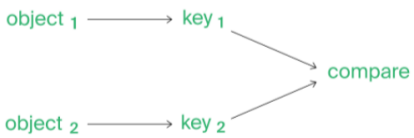
Алла → алла  
АЛЛА → алла

Приведение к «нормальной форме» — часто применяемый метод. Помните, как вы искали группы анаграмм в спринте про хеши? Для решения этой задачи нужно было составить хеш-таблицу, в которой значениями были оригинальные слова, а ключами — их нормализованные версии. Например, можно было отсортировать буквы слова в лексикографическом порядке, и тогда "LISTEN" и "SILENT" попадали в одну и ту же корзину "EILNST". Проверить, что два слова — анаграммы, можно следующим образом: `compare(sorted(first), sorted(second))`.

Обратите внимание, насколько разные сравнения можно проводить по одной и той же схеме:

- сначала мы по строке вычисляем некоторый ключ (например, нормальную форму, длину строки или любую другую её характеристику),
- а затем сравниваем значения ключей.

Мы подробно говорили об этом в спринте про сортировки.



Работа с ключами вместо объектов — распространённый приём. Он встречается и при работе с хеш-таблицами, и в сортировках, и при сравнении строк

Аналогичным образом можно сопоставлять слова, написанные кириллицей и транслитом. Или сравнивать слова по звучанию, преобразовывая их к специальному фонетическому алфавиту.

Гоша: Здорово. Про равенство двух строк теперь понятно. А как строки упорядочить? Какая строка больше, какая меньше?

Рита: Помните, мы говорили про лексикографическую сортировку: сначала сравниваются первые символы строки, если они равны, сравниваются вторые символы, и так далее, пока не найдётся различающийся символ. Сами символы сравниваются по их номеру в таблице кодировки.

Будьте аккуратны с кодировками, а также с регистром букв. Попробуем для примера сравнить две строки: "A\_a" и "Aa\_" , а затем привести их к нижнему и верхнему регистрам и снова сравнить уже нормализованные строки:

- "A\_a" < "Aa\_" даёт True — сравнение исходных строк;
- "a\_a" < "aa\_" даёт True — сравнение строк, приведённых к нижнему регистру;
- "A\_A" < "AA\_" даёт False — сравнение строк, приведённых к верхнему регистру.

Такой неожиданный результат сравнения получается из-за того, что символ подчёркивания идёт в таблице кодировки ASCII позднее буквы верхнего регистра, но раньше букв нижнего регистра: 'A' < '\_' < 'a' .

Решите задачи «Пограничный контроль» и «Сравнить две строки».

Вперёд









- Практикум
- Финальное задание спринта 7
- Финальное задание спринта 7
- Новое достижение
- Оцените задание и сопроводение
- спринт 8
- Алгоритмы на строках
- Строки. Простейшие операции
- Сравнение строк
- Строки. Простейшие операции
- Сравнение строк
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Строки. Простейшие операции
- Сравнение строк
- Подстроки, префиксы и суффиксы
- Поиск шаблона в строке. Наивный алг...
- Префикс-функция
- Вычисление префикс-функции
- Эффективный поиск шаблона в тексте
- Префиксное дерево
- Разбор задачи «Быстрое сравнение д...
- Резюме
- Оцените пройденные уроки
- Финальное задание спринта 8

# Эффективный поиск шаблона в тексте

Теперь рассмотрим алгоритм поиска подстроки в строке с использованием префикс-функции.

Пусть у нас есть шаблон  $p = \text{"pattern"}$  и текст  $t = \text{"text"}$ . Обе строки состоят из символов некоторого алфавита  $\Sigma$ . Составим из них одну общую строку, разделённую символом-сентинелом.

**Сентинел** — специальный символ, который не содержится в алфавите  $\Sigma$ . Его используют, чтобы при конкатенации двух строк обозначить место склейки. Мы в качестве сентинела возьмём символ  $\#$ .

Итак, у нас получилась строка  $\text{"pattern\#text"}$ . Именно от неё мы и будем вычислять префикс-функцию:  $\pi(\text{"pattern\#text"})$ .

В качестве примера возьмём шаблон  $\text{"sip"}$  и текст  $\text{"mississippi"}$ . Составим комбинированную строку  $\text{"sip\#mississippi"}$  и вычислим от неё префикс-функцию.

Как мы говорили, каждая подстрока является суффиксом некоторого префикса. Значит, если шаблон встречается в строке, то на него оканчивается какой-то префикс. Тогда соответствующая подстрока должна иметь вид  $\text{"pattern\#...pattern"}$ . В нашем случае такой подстрокой окажется  $s_{[0,13]}$ .

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
s[i]	s	i	p	#	m	i	s	s	i	s	i	p	p	i	
$\pi[i]$	0	0	0	0	0	0	1	2	1	1	3	0	0		

Вычислим значения префикс-функции для строки, полученной конкатенацией шаблона, сентинела и текста

Префикс-функция  $\pi_{i2}$  от такого префикса будет равна длине шаблона  $|p|$ . Это свойство любой позиции, на которой оканчивается вхождение шаблона. Значит, чтобы найти в тексте все вхождения шаблона, нам достаточно посчитать префикс-функцию комбинированной строки  $\text{"pattern\#text"}$  и найти позиции, где она равна  $|p|$ . Эти позиции будут соответствовать последней букве шаблона.

Стоит сказать, что  $|p|$  — это максимально возможное значение префикс-функции, ведь префикс большей длины обязательно захватывал бы сентинел, который встречается в строке лишь единожды. Чтобы суффикс начинался на  $\text{"pattern\#"}$ , он должен совпадать со всей строкой, но несобственные суффиксы в вычислении префикс-функции не участвуют.

В таком методе для поиска шаблона в тексте не требуется ничего, кроме вычисления префикс-функции.

## Эффективность алгоритма

Этот алгоритм состоит из двух шагов:

- построение префикс-функции,
- сканирование полученного массива в поисках нужного значения.

Оба этапа занимают линейное время  $O(|p| + |t|)$ . Это намного быстрее, чем  $O(|p| \cdot |t|)$  — скорость работы наивного алгоритма.

Описанный метод требует  $O(|p| + |t|)$  дополнительной памяти для хранения значений префикс-функции. Впрочем, его возможно модифицировать так, чтобы он тратил только  $O(|p|)$  памяти. Для этого следует хранить префикс-функцию только от первых  $|p|$  символов комбинированной строки (то есть от той части строки, которая соответствует шаблону), а также префикс-функцию для последнего обрабатываемого элемента.

В таком случае построение префикс-функции и поиск её максимального значения будут совмещены. Промежуточные значения префикс-функции можно не запоминать: динамическое программирование всегда обращается только к последнему вычисленному значению префикс-функции и первым  $|p|$  значениям.

```

функция поиск(p, text):
    # функция возвращает все позиции вхождения шаблона в тексте.
    result = []
    s = p + '#' + text
    pi = [0, None, None, ...] # Массив длины |p|.
    pi_prev = 0
    для i из [1 .. |s|]:
        k = pi_prev
        пока (k > 0) и (s[k] != s[i]):
            k = pi[k - 1]
        если s[k] == s[i], то:
            k += 1
        # Запоминаем только первые |p| значений pi-функции.
        если i < |p|, то:
            pi[i] = k
        # Запоминаем последнее значение pi-функции.
        pi_prev = k
    # Если значение pi-функции равно длине шаблона, то вхождение найдено.
    если k == |p|, то:
        # i - это позиция конца вхождения шаблона.
        # Дважды отнимаем от него длину шаблона, чтобы получить позицию начала:
        # - чтобы «переместиться» на начало найденного шаблона,
        # - чтобы не учитывать добавленное "pattern".
        result.добавить(i - 2 * |p|)
    вернуть result
    
```

Решите задачу «Глобальная замена».



Практикум	<
▼ Финальное задание спринта 7	✓
★ Финальное задание спринта 7	✓
Новое достижение	
Оцените задание и сопровождение	✓
спринт 8	
▼ Алгоритмы на строках	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Строки. Простейшие операции	✓
Сравнение строк	✓
Подстроки, префиксы и суффиксы	✓
Поиск шаблона в строке. Наивный алг...	✓
Префикс-функция	✓
Вычисление префикс-функции	✓
Эффективный поиск шаблона в тексте	✓
Префиксное дерево	✓
Разбор задачи «Быстрое сравнение д...	✓
Резюме	✓
Оцените пройденные уроки	✓
► Финальное задание спринта 8	✓

## Разбор задачи «Быстрое сравнение двух строк»

Задача сравнения двух строк решается достаточно просто за линейное время: если надо сравнить строки  $s$  и  $t$ , то достаточно найти их наибольший общий префикс и посмотреть на первый различающийся символ, если такой есть, или заключить, что одна строка является префиксом другой. Результат операции сравнения можно представить в виде числа. Если  $s < t$ , то результат равен  $-1$ , если  $s = t$ , то  $0$ , если  $s > t$ , то  $1$ .

Теперь расширим эту задачу. Пусть нам дан набор, состоящий из  $n$  строк. Наша задача — научиться быстро отвечать на запросы вида «сравнить строки  $s$  и  $t$  из этого набора». Заодно разберёмся, что в данном случае значит «быстро».

### Наивное решение

Наиболее простой подход состоит в том, чтобы применять обычный алгоритм сравнения двух строк. Тогда отвечать на запрос мы будем за  $O(L)$ , где  $L$  — длина самой длинной строки.

Если запросов будет поступать много, то имеет смысл воспользоваться кешированием ответов.

Кеширование — это техника запоминания уже запрошенных данных и полученных ответов на запросы. Например, в ленивом динамическом программировании мы сохраняли ответы для однажды посчитанных состояний, чтобы не пересчитывать их повторно.

В рассматриваемой задаче, сохраняя ответы для каждой запрашиваемой пары строк, мы сэконоим время на повторяющихся запросах, но увеличим пространственную сложность. На хранение всех результатов уйдёт  $O(n^2)$  дополнительной памяти, так как всего пар  $\frac{n(n-1)}{2}$ , и для каждой надо сохранить одно число — результат сравнения ( $-1$ ,  $0$  или  $1$ ). Время работы на подсчёт ответов для всех пар составит  $O(n^2L)$ .

### Улучшаем решение

Теперь научимся быстрее обрабатывать каждый запрос по отдельности. Рассмотрим утверждение для некоторого числа  $i$  и двух строк  $s$  и  $t$ : «Префиксы длины  $i$  строк  $s$  и  $t$  совпадают».

Заметим, что это утверждение до некоторого значения  $i$  истинно, но, начиная с некоторого  $i$ , ложно. Например, для строк  $s = \text{"abacaba"}$  и  $t = \text{"abadaba"}$  при  $i = 1, 2, 3$  утверждение верно, а при  $i \geq 4$  ложно.

i	1	2	3	4	5	6	7
s	a	b	a	c	a	b	a
t	a	b	a	d	a	b	a
результат	1	1	1	0	0	0	0

Такое свойство некоторого условия называется «монотонностью». На самом деле мы уже имели дело с этим свойством, когда говорили о бинарном поиске.

Бинарный поиск может применяться только на монотонных данных. Помните, мы разобрали поиск по словарю? Записи в нём упорядочены лексикографически — именно это и есть монотонность, благодаря которой мы успешно применили бинарный поиск и нашли нужное нам слово.

Обозначим  $T(k)$  — время на сравнение двух префиксов длины  $k$  на точное равенство. Тогда первую позицию, в которой строки  $s$  и  $t$  различаются, можно искать бинарным поиском за  $O(T(L) \cdot \log L)$ . Теперь подумаем, как можно сравнивать префиксы эффективно. Если сравнивать их наивно, за линейное время,  $T(k) = O(k)$ , то толку от использования бинарного поиска не будет.

Допустим, для строк  $s$  и  $t$  заранее вычислены массивы [префиксных хешей](#). Тогда сравнить два префикса длины  $k$  можно за константное время  $O(1)$ , обратившись к нужному префиксному хешу для каждой строки. Значит, сравнить строки  $s$  и  $t$  мы можем за  $O(\log L)$  при условии, что префиксные хеши уже посчитаны.

Мы готовимся к тому, что запросов будет приходиться много, поэтому выгодно посчитать префиксные хеши заранее, сделав предподсчёт. Префиксный хеш одной строки  $s$  вычисляется за  $O(|s|)$ , поэтому суммарно это займёт  $O(nL)$  времени, что сопоставимо с затратами на считывание строк. На хранение хешей уйдёт  $O(nL)$  дополнительной памяти. Зато время ответа на каждый запрос улучшилось с  $O(L)$  до  $O(\log L)$ !

```
PYTHON
all_strings = ... # считываем n строк.
n = len(all_strings)

# Вычисляем префиксные хеши.
# Как это делать -- изучали в спринте 4. :)
hashes = [get_hashes(s) for s in all_strings]
# В hashes[i] лежит массив префиксных хешей для i-ой строки.

queries = ... # считываем все запросы.
# Каждый запрос состоит из пары индексов строк, которые нужно сравнить.

for i, j in queries:
    # Итерация цикла сравнивает i-ую и j-ую строки.
    # В цикле перебираются пары индексов.
    s = all_strings[i]
    t = all_strings[j]

    # -- Бинарный поиск позиции первого различия --
    # Будем думать про left, mid и right как про длины общих префиксов.
    # Изначально мы можем гарантировать совпадение только нулевой длины.
    left = 0
    # Заведомо недостижимый ответ. Это, например, длина строки + 1.
    right = min(len(s), len(t)) + 1
    while right - left > 1:
        # Получаем середину.
        mid = (left + right) // 2
        # Так как mid - это длина, то [mid-1] - индекс последнего символа.
        if hashes[i][mid - 1] == hashes[j][mid - 1]:
            # Так как при длине префиксов mid они совпадают,
            # то переводим left.
            left = mid
```

Итого, используя алгоритм с предподсчётом префиксных хешей, мы получаем асимптотику  $O(nL + Q \log L)$ , если нам поступает  $Q$  запросов на сравнение пары строк. Иногда для алгоритма отдельно записывают время на предподсчёт и время ответа на запрос. В данном случае предподсчёт занял бы  $O(nL)$ , а ответ на каждый запрос  $O(\log L)$ .

Если бы мы захотели ещё улучшить рассмотренный алгоритм, то можно было бы добавить кеширование ответов. Проанализируем время работы алгоритма с добавленным кешированием. Предподсчёт хешей займёт  $O(nL)$ , затем на вычисление ответов для всех пар потребуется  $O(n^2 \log L)$  операций, далее ответ на любой запрос выполняется за  $O(1)$ . Затраты на память составили бы  $O(nL)$  на префиксные хеши плюс  $O(n^2)$  на кеш ответов.

**Вопрос.** Возможно ли подобрать две неравных строки  $s$  и  $t$  таких, чтобы их полиномиальные хеши совпадали? Стоит ли учитывать это при решении задачи?



▼ Финальное задание спринта 7 ✓

★ Финальное задание спринта 7 ✓

Новое достижение

Оцените задание и сопровождение ✓

СПРИНТ 8

▼ Алгоритмы на строках ✓

Строки. Простейшие операции ✓

Строки. Простейшие операции ✓

Сравнение строк ✓

Подстроки, префиксы и суффиксы ✓

Поиск шаблона в строке. Наивный алг... ✓

Префикс-функция ✓

Вычисление префикс-функции ✓

Эффективный поиск шаблона в тексте ✓

Префиксное дерево ✓

Разбор задачи «Быстрое сравнение д... ✓

Резюме ✓

Оцените пройденные уроки ✓

► Финальное задание спринта 8 ✓

## Резюме

При вставке текста в середину строки следует выделить больший блок памяти и сдвинуть участок строки, расположенный правее места вставки. Причём символы нужно смещать в конец выделенного блока памяти, начиная с последнего.

Мы разобрали, как сравнивать строки на равенство несколькими способами: от простого сравнения и сравнения без учёта регистра до более сложного — по произвольному признаку.

Мы рассказали про наивный алгоритм поиска подстроки в строке, работающий за квадратичное время и объяснили, на каких входных данных он будет работать максимально неэффективно.

Алгоритм Кнута-Морриса-Пратта основан на понятии префикс-функции  $\pi(s)$  и позволяет искать подстроку в строке за линейное время, эффективно сдвигая шаблон вдоль строки.

Напомним, что префикс-функция — это массив чисел, сформированный для строки так, что  $i$ -ый элемент в нём равен длине наибольшего собственного суффикса  $s_{[0,i)}$ , который также является префиксом  $s_{[0,i)}$ . Префикс-функция шаблона может быть вычислена за линейное время.

Также мы рассказали про структуру данных бор, или префиксное дерево. Оно используется для эффективного хранения наборов строк, когда одиночной строки в качестве шаблона оказывается недостаточно.

[Вперёд](#)

Новое достижение

Оцените задание и сопровождение ✓

СПРИНТ 8

▼ Алгоритмы на строках ✓

Строки. Простейшие операции ✓

Сравнение строк ✓

Подстроки, префиксы и суффиксы ✓

Поиск шаблона в строке. Наивный алг... ✓

Префикс-функция ✓

Вычисление префикс-функции ✓

Эффективный поиск шаблона в тексте ✓

Префиксное дерево ✓

Разбор задачи «Быстрое сравнение д... ✓

Резюме ✓

Оцените пройденные уроки ✓

▼ Финальное задание спринта 8 ✓

★ Финальное задание спринта 8 ✓

Оцените задание и сопровождение ✓

Решите оставшиеся задачи по этой теме в Яндекс.Контесте:  
<https://contest.yandex.ru/contest/26131/problems/>

Когда закончите, переходите к финальному заданию ниже.

## Финальное задание спринта 8

Вам предстоит решить две задачи и пройти код-ревью по ним.

Сдавайте решения в Яндекс.Контесте, а когда получите ОК по обоим задачам, — запакуйте исходные файлы в один архив и загрузите на ревью. В начале каждого решения в комментарии:

- напишите объяснение так же, как в предыдущем спринте;
- укажите ID успешной отправки, чтобы ревьюер мог удостовериться, что решение рабочее.

Важно: выберите один язык программирования и сдавайте финальные задачи на ревью только на нём. Задачи, которые не требуют ревью, вы можете сдавать на любом из доступных языков.

Успехов!

Вперед

