

Алгоритм Рабина-Карпа

Николай Пономарев

Проблема

- Поиск шаблона в строке достаточно частая проблема в информатике
- Самым простым решением является последовательное «прикладывание» шаблона к строке
- Сложность такого решения $O(nm)$
- Можно оптимизировать выходом из цикла при первом несовпадении, но гарантий на значительное ускорение нет

В информатике существуют так называемые хеш-функции, которые позволяют преобразовать произвольные данные в конечное множество значений

- У одинаковых входных данных одинаковый хеш
- Однако обратное неверно, т.к. происходят коллизии

Коллизия – совпадение результата хеш-функции для разных входных данных

- У «хороших» хеш-функций вероятность коллизий достаточно мала

Задумка

Алгоритм Рабина-Карпа был придуман Ричардом Карпом и Майклом Рабином в 1987 году.

Идея алгоритма заключается в подсчете хеша для подстроки и сравнении его с хешем образца, и только в случае равенства посимвольно сравнивать подстроку и образец

Реализация

```
1 def rabin_karp(s: str, pattern: str) → int:
2     n = len(s)
3     m = len(pattern)
4     pattern_hash = hash(pattern)
5     for i in range(n - m):
6         substring = s[i : i + m]
7         substring_hash = hash(substring)
8         if pattern_hash == substring_hash:
9             if pattern == substring:
10                 return i
11     return -1
```

Сложность

- Хеширование и посимвольное сравнение имеют сложность $O(m)$
- Хеширование производится на каждом шаге цикла, итоговая сложность $O(nm)$

Но такая же сложность и у «наивного» алгоритма!

Для решения этой проблемы необходимо использовать скользящий (rolling) хеш.

Скользящий хеш

Скользящий хеш — это такая хеш-функция, для пересчета которой достаточно знать ее предыдущее значение и что изменилось в данных.

Пример

Самый простой хеш, который складывает все коды символов в строке $s = a_1 a_2 \dots a_n$:

$$h(s) = \sum_{i=0}^n a_i,$$

является скользящим. После сдвига на 1 вправо его можно пересчитать за константное время:

$$h(s[2..n + 1]) = h(s[1..n]) - a_1 + a_{n+1}$$

Полиномиальный хеш

Полиномиальный хеш является достаточно простым скользящим хешем, однако он обеспечивает достаточно низкую вероятность коллизий.

$$h(s) = \left(\sum_{i=1}^n s[i]x^{n-i} \right) \bmod q$$

На практике взятие остатка обычно выполняется после каждой операции:

```
1 def polynomial_hash(s: str, x: int, q: int) → int:
2     hash = 0
3     for a in s:
4         hash = (hash * x) % q
5         hash = (hash + a) % q
6     return hash
```


Полиномиальный хеш. Сдвиг

За основу возьмем строку $s = a_1 a_2 a_3 a_4 \dots a_n$

Посчитаем хеш для подстроки $s_1 = a_1 a_2 a_3$:

$$h(s_1) = [([([([a_1 x] \bmod q + a_2] \bmod q) x] \bmod q) + a_3] \bmod q$$

Теперь сдвинемся на 1 символ вправо и подсчитаем хеш для подстроки $s_2 = a_2 a_3 a_4$:

$$h(s_2) = [(h(s_1) + q - [(a_1 x^2) \bmod q]) x + a_4] \bmod q$$

Полиномиальный хеш. Выбор q и x

Вариант Рабина и Карпа:

- фиксированный $x = 2$
- случайное простое число $q \in [2, n^3]$, где n – длина строки

Вариант Дитзфелбингера и др.:

- случайное $x \in (0, q - 1]$
- фиксированное простое q (например простое число Мерсенна: $2^{31} - 1$ или $2^{61} - 1$)

Сложность с «хорошим» хешем

- Вычисление скользящего хеша занимает $O(n)$
- Сравнение строк при совпадении хеша – $O(m)$

Тогда сложность алгоритма $O(n + m)$

Применение

- Алгоритм легко модифицируется для поиска набора образцов в тексте сохраняя хорошее время работы
- Из-за этого часто применяется для проверки на плагиат

Материалы

- Русская Википедия
- Английская Википедия
- Rabin M. O., Karp R. M. Efficient randomized pattern-matching algorithms – статья Рабина и Карпа с описанием алгоритма
- Dietzfelbinger M., Gil J., Matias Y., Pippenger N. Polynomial hash functions are reliable – статья Дитцфелбингера и др. о полиномиальном хеше
- Rabin M. O. Fingerprinting by random polynomials – статья про хеш Рабина