

# Поиск подстроки

Обозначения:

$T[0:m]$  - текст

$T[i:j]$  - подстрока

$P[0:n]$  - образец

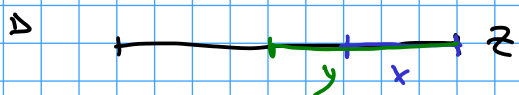
$x, y$  - строки

$x \sqsubseteq y$  -  $x$  - префикс  $y$

$x \supseteq y$  -  $x$  - суффикс  $y$

Утв:

$x \supseteq z, y \supseteq z, |x| \leq |y| \Rightarrow x \supseteq y$



## Задача о поиске подстроки

Вход:  $T[0:m], P[0:n]$

Найти  $i$  (наименьшее):

$$T[i, i+n] = P$$

Наивный алгоритм

for  $i = 0$  to  $m-n$

$$O((m-n) \cdot n) = O(m \cdot n)$$

for  $j = 0$  to  $n$

$$O(n)$$

if  $P[j] \neq T[i+j]$

continue

return  $i$

## Алгоритм Карпа - Рабина

$T[0:m], P[0:n]$

Для всех  $i$  от  $0$  до  $m-n$

вычисляем  $\text{hash}(T[i, i+n]) = h_i$

//  $m-n$  хешей

for  $i = 0$  to  $m - n$ :

if  $h_i = \text{hash}(P)$ :

if  $P = T[i, i+n]$

// цена  $O(n)$

return  $i$

Полиномиальная хеш

$S[0:n]$

$$\text{hash}(S) = q^n S[0] + q^{n-1} S[1] + \dots + q^1 S[n-1] + S[n] \pmod N$$

$$h_i = \text{hash}(T[i, i+n])$$

$$h_{i+1} = (h_i - q^n T[i]) \cdot q + T[i+n+1] \pmod N$$

$\Rightarrow$  Мы можем вычислить все  $h_i$

$$\text{for } C_1 n + (m - n) \cdot C_2 = O(m)$$

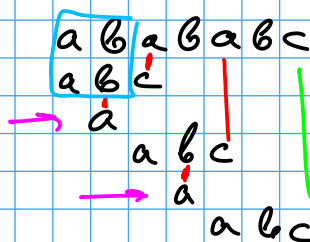
Оценим. В предположении гипотезы равномерного хеширования вероятность коллизии  $\frac{1}{N} \leq \frac{1}{m}$

$$O(m + m + \cancel{m/N} \cdot n) = O(m + n) \leq 1$$

Комбинаторный подход

$T = \text{abababc}$

$P = \text{abc}$



$z$  - функция

$S = \underline{\text{abababc}}$

$$z(S)[i] = j \Leftrightarrow$$

$S[i, i+j] \subseteq S, j - \max$

$z(S) = 0040200$

Поиск подстроки:

$T[0, m], P[0, n]$

$z(P \# T)$

$z(abc \# abcabcabc)$

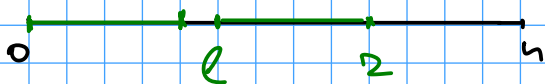
00002020300

Временные:

Наивно  $O(n^2)$

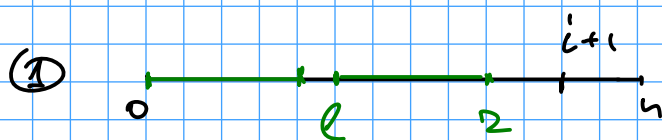
Худший случай: aaaaaaaaaab

Временные  $O(n)$ :

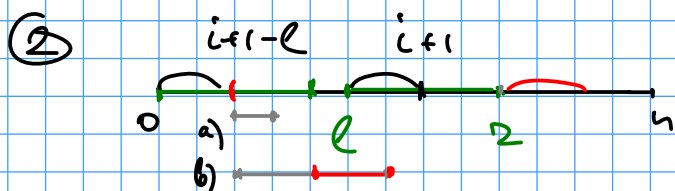


$z[i+1] - ?$  Все  $z[i]$  вычисляем

$l, z$ :  $S[l, z]$  - наибольший префикс который мы уже нашли



Вычисляем  $z[i+1]$  наилучшим образом



a)  $z[i+1] = \min(z[i+1-l], n-i)$

b)  $z[i+1] = \min(z[i+1-l], z-i)$

проверяем символы  $O(z)$

Пример:

$abc \# abcabcabc$   
 $00002020300$

Z-Function ( $S[0, n]$ ):

$z(0, n)$

$l, r = 0$

for  $i = 1$  to  $n$

if  $i \leq r$ :

$z[i] = \min(r - i + 1, z[i - l])$

while  $z[i] + i \leq n$  and  $S[z[i]] = S[z[l] + i]$

$++ z[i]$

if  $i + z[i] > r$

$l = i, r = i + z[i]$

$O(n)$

В итоге: поиск подстроки  $\downarrow$   $O(m+n)$   
(находим все вхождения)

Алгоритм Кнута - Мориса - Пратта

= Префиксы - функции  $\pi$

$abcabcabc$   
 $P = abcabcabd$

сверяем на  $i - \pi(i)$

$P_1 = P[0]$

$P_2 = P[0, 1]$

$\vdots$

$\uparrow$

все префиксы  $P$

$P \quad abcabcabd$

$\pi \quad 000123450$

$\pi(i) = \max_{k < i} k : P_k \text{ } P_i$

KMP( $T, P$ )

$\pi \leftarrow$  Prefix Function( $P$ )

$k = 0$  // *сбрасываем курсор*

for  $i = 0$  to  $m$ :

while  $k > 0$  and  $P[k+1] \neq T[i]$

$k = \pi(k)$

if  $P[k+1] = T[i]$

$k = k + 1$

if  $k = n$ :

return  $i - n + 1$

$O(m)$

Prefix Function( $P$ ):

$k = 0$

$\pi[i] = 0$

for  $i = 1$  to  $n$

while  $k > 0$  and  $P[k+1] \neq P[i]$

$k = \pi(k)$

if  $P[k+1] = P[i]$

$k = k + 1$

$\pi[i] = k + 1$