

Поиск подстроки

Обозначения:

$T[1:m]$ - текст

$P[1:n]$ - шаблон

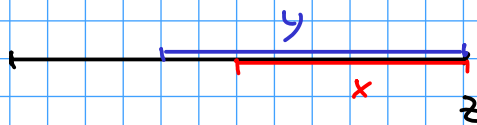
x, y - строки

$x \supset y$ - x - суффикс y

$x \sqsubset y$ - x - префикс y

Утв:

$x \supset z, y \supset z \quad |x| \leq |y| \Rightarrow x \supset y$



Задача о поиске шаблона

Вход: T, P

- Выход:
1. мин i : $P \supset T[i:i+n-1]$
 2. все такие i

Наивный подход

$O(n \cdot m)$

$T = a a a \dots a$

$P = a a a a b$

Алгоритм Карпа - Рабина

$O(n+m)$ в среднем

Z - функция

$Z[i] =$ длина макс префикса T начинающегося в i

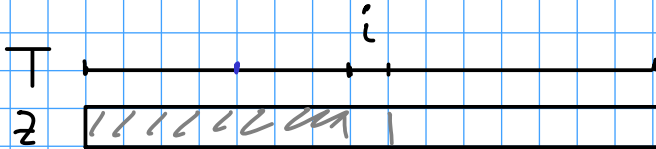
$$T = \underline{a} \underline{b} c a \underline{b} a c a \underline{b} a \underline{b} c$$

$$Z = 1200201020300$$

↖ не встречается в P и T
 Найти $Z(P \# T)$ мы решим обе задачи

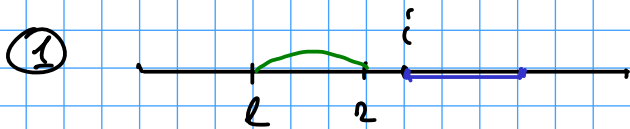
] $Z[1] \dots Z[i-1]$ - посчитаем

$$Z[i] = ?$$

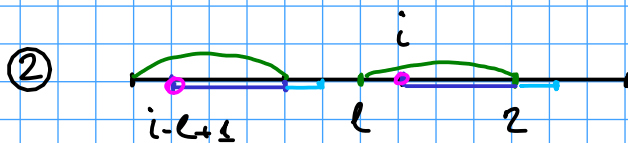


] $T[l, r]$ - самый правый префикс P, который мы до этого момента посчитали.

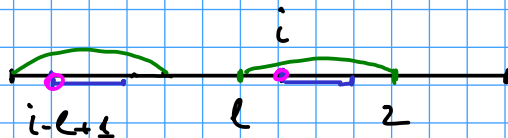
$$\text{т.е. } Z[l] = |T[l, r]|$$



мы вычислим $Z[i]$ напрямую



$$Z[i] = Z[i-l+1] + \text{проверить границы}$$



$$Z[i] = Z[i-l+1]$$

$z(T[1:m]):$

$$z[i] = 0$$

$O(m)$

$$z[1] = m$$

$$l = r = 0$$

for $i = 2$ to m

if $i \leq l$

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

$O(m)$ while $z[i] + i \leq m$ & $T[z[i]+1] = T[i+z[i]]$

$$z[i] = z[i] + 1$$

if $i + z[i] - 1 > r$

$$l = i$$

$$r = i + z[i] - 1$$

\Rightarrow Поиск маджона (т.е. $z(P \# T)$)
работает за $O(m+n)$

Алгоритм Кнехта - Мориса - Пратта

$\pi(P)$ - массив

$\equiv \pi[i] = \max j < i : P[1:j] \supseteq P[1:i]$

$P = \underline{a} \underline{b} \underline{c} \underline{a} \underline{b} \underline{a} \underline{c} \underline{a} \underline{b} \underline{a} \underline{b} \underline{c} \equiv$ префикс - суффиксы

$\pi = 0 0 1 2 1 0 1 2 1 2 3$

Аналогично: $\pi(P \# T)$ то мы решим
обе задачи поиска

$T = \begin{array}{|c|c|c|c|c|c|c|} \hline a & b & c & a & b & c & a & b & c \\ \hline a & b & c & a & b & & & & \\ \hline \end{array}$

$P = a b c a b d$

\rightarrow
 \rightarrow
 $\quad a$
 $\quad \quad a$
 $\quad \quad \quad a b c a b d$

Prefix Function ($P[1:n]$)

$k = 0$ // *собрано совпадений*

$O(n)$

for $i = 1$ to n

$O(n)$ | while $k > 0$ && $P[k+1] \neq P[i]$
 $k = \pi[k]$
if $P[k+1] = P[i]$ |
 $k = k + 1$
 $\pi[i] = k$

abcab | abcabc
 $k=5$
 $k=2$

Search ($P[1:n], T[1:m]$)

$k = 0$ // *собрано совпадений*

$O(n)$ $\pi \leftarrow$ Prefix Function (P) $O(n+m)$

for $i = 1$ to m

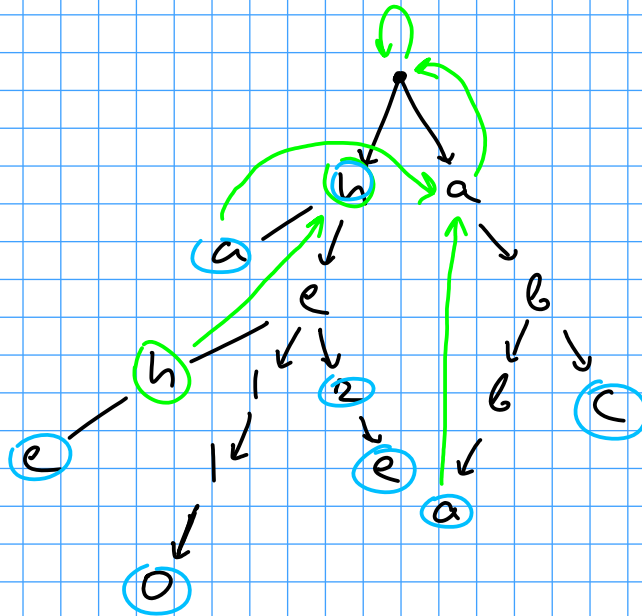
$O(m)$ | while $k > 0$ && $P[k+1] \neq T[i]$
 $k = \pi[k]$
if $P[k+1] = T[i]$
 $k = k + 1$
if $k = n$
 return i

$P =$ abcabacababc
 $\pi = 000121012123$

Алгоритм Ахо-Корасик

Задача: множество $\{P_i\}$ - набор слов
Найти все вхождения P_i в T

Бор $\{hello, her, here, abc, abba\}$
 h, ha, he



~ конечному автомату

1. Строим бор из наборов

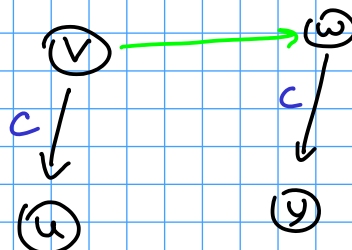
$$O(\sum |P_i|) \quad \text{Память: } O(\sum |P_i| \cdot |A|)$$

$= n$

\uparrow
алфавит

2. Восстановив набор суррисакции стрелки

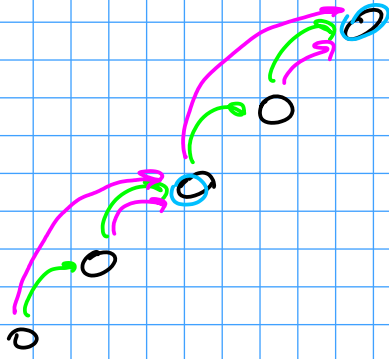
Поиск в ширину



$$O(n)$$

3. Для нахождения всех вхождений
добавляем "сжатые стрелки":
заменяем только в

поиск узла маджона



$O(m)$

Сложность: $O(n) + O(m) + O(\# \text{выход})$

↑
построили
дерево +
списки

↑
прошли
текут

$O(m + n + \# \text{выход})$