

# Vyhľadávanie podreťazcov v reťazci (String matching)

## 1. Naive method (Priama metóda)

Pre každý z možných posunov sa pokúsi vyhľadať reťazec  $P$ .

**Zložitosť:**  $O((n - m + 1)m)$

**Najhorší prípad:** vyhľadávanie podreťazca  $a^m$  v  $a^n$ ,  $m < n$

$T$  - reťazec v ktorom sa vyhľadáva

$P$  - hľadaný podreťazec

$|\cdot|$  - dĺžka reťazca

$P[a..b] = a..b$  prvok v reťazci  $P$

1.  $n := |T|$
2.  $m := |P|$
3. for  $s:=0$  to  $n-m$  do
4.   if  $P[1..m] = T[s+1..s+m]$  then
5.     print "Podreťazec sa našiel pri posunutí",  $s$

## 2. Rabin-Karp algorithm (Rabin-Karpov algoritmus)

Použitá abeceda:  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

Toto nám umožní brať podreťazce  $T$  a samotné  $P$  ako číslo v desiatkovej sústave.

Podmienka  $P[1..m] = T[s + 1..s + m]$  sa tak prevedie na porovnávanie dvoch čísel. Hodnota sa vypočítava z Hornerovej schémy. Pre každý posun v  $T$  stačí vynechať prvú cifru, celé vynásobiť 10 a pripočítať nasledujúcu cifru.

$$\begin{array}{c} h \\ \boxed{a_k \mid a_{k+1} \mid \dots \mid a_{k+m-1} \mid a_{k+m}} \end{array} \rightarrow \begin{array}{c} h' \\ \boxed{a_k \mid a_{k+1} \mid \dots \mid a_{k+m-1} \mid a_{k+m}} \end{array}$$
$$h = 10^{m-1}a_k + 10^{m-2}a_{k+1} + \dots + 10a_{k+m-2} + a_{k+m-1} \rightarrow$$
$$h' = 10 * (h - 10^{m-1}a_k) + a_{k+m}$$

Pre dlhé podreťazce  $P$  môžeme dostávať veľmi veľké čísla. (Efektívne by bolo dobré dostávať čísla, ktoré sú menšie ako dĺžka slova na procesore -  $0 - 2^{32} - 1$ , (32bit),  $0 - 2^{64} - 1$  (64bit).)

Preto sa to rieši pomocou výpočtu modulo  $q$ . Najprv sa porovnajú získané čísla modulo  $q$  a v prípade rovnosti sa ešte porovnajú člen, po člene. (Pretože, v prípade nízkej hodnoty modula sa kľudne môže stať, že dva rôzne podreťazce zodpovedajú rovnakému číslu modulo  $q$ .)

- Podobne ako pri základe 10, možno predošlé urobiť pre iný číselný základ  $d$ .

**Zložitosť:**

Preprocessing (príprava) -  $\Theta(m)$

Výpočet (Computation) -  $O((n - m + 1)m)$

**Najhorší prípad:** vyhľadávanie podreťazca  $a^m$  v  $a^n$ ,  $m < n$

- Oproti Naive dáva lepšie priemerné časy.

Rabin-Karp Matcher( $T, P, d, q$ )

1.  $n := |T|$
2.  $m := |P|$
3.  $h := d^{m-1} \pmod q$
4.  $p := 0$
5.  $t_0 := 0$
6. for  $i:=1$  to  $m$  do {preprocessing}
7.      $p := (d * p + P[i]) \pmod q$
8.      $t_0 := (d * t_0 + T[i]) \pmod q$
9. for  $s:=0$  to  $n - m$  do {computation}
10.     if  $p = t_s$  then
11.         if  $P[1..m] = T[s + 1..s + m]$  then
12.             print "Podreťacec sa našiel s posunom",  $s$
13.     if  $s < n - m$  then
14.          $t_{s+1} := (d * (t_s - T[s + 1]h) + T[s + m + 1]) \pmod q$

### 3. Konečný automat

Def.: *Konečný automat (finite automaton)* je usporiadaná 5-tica  $(Q, q_0, A, \Sigma, \delta)$ , kde

- $Q$  je konečná množina stavov
- $q_0$  počiatočný stav
- $A$  množina koncových stavov (akceptujúce)
- $\Sigma$  je použitá abeceda
- $\delta$  je tzv. *prechodová funkcia* z  $Q \times \Sigma$  do  $Q$ .

Rozšírenie  $\delta$  funkcie -  $\delta^* : Q \times \Sigma^* \rightarrow Q$  je definované indukzívne:

$$\begin{aligned} \delta^*(q, \epsilon) &= q \\ \delta^*(q, wa) &= \delta(\delta^*(q, w), a) \end{aligned}$$

Final-state function - vracia stav automatu po spracovaní nejakého slova

*Suffix funkcia* pre  $P$ ,  $|P| = m$  je  $\sigma : \Sigma^* \rightarrow \{0, 1, \dots, m\}$  definovaná ako

$$\sigma(x) = \max\{k : P_k \sqsupseteq x\},$$

kde  $u \sqsupseteq v$  znamená, že  $u$  je sufixom  $v$  a  $P_k = P[1..k]$ .

**Definícia automatu:** pre  $P$ ,  $|P| = m$

$$Q = \{0, 1, \dots, m\}, q_0 = 0, A = \{m\}, \delta(q, a) = \sigma(P_q a).$$

Vždy, keď sa počas simulácie vstupného slova  $T$  na automate dostaneme do stavu  $m$ , našiel sa podvýraz  $P$  a jeho posun je rovný o  $m$  menej ako je aktuálna pozícia v reťazci.

Platia nasledujúce vety:

V (suffix-function inequality): Pre každý reťazec  $x$  a znak  $a$  platí:  $\sigma(xa) \leq \sigma(x) + 1$ .

V (suffix-function recursion lemma): Pre každý reťazec  $x$  a znak  $a$ , ak  $q = \sigma(x)$ , tak  $\sigma(xa) = \sigma(P_q a)$ .

VÝPOČET PRECHODOVEJ FUNKCIE ( $P, \Sigma$ )

1.  $m := |P|$
2. for  $q := 0$  to  $m$  do
3.     for each symbol  $a \in \Sigma$  do
4.          $k := \min(m + 1, q + 2)$
5.         repeat  $k := k - 1$
6.         until  $P_k \sqsupset P_q a$
7.          $\delta(q, a) := k$
8. return  $\delta$

Zložitosť tejto funkcie je  $O(m^3|\Sigma|)$ . Dá sa zlepšiť na  $O(m|\Sigma|)$ .  
Zložitosť samotného výpočtu je  $\Theta(n)$ .

#### 4. KMP (Knuth, Morris, Pratt)

Prefixová funkcia (Prefix function)  $\pi$  pre  $P$ ,  $|P| = m$ :

$\pi : \{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, m - 1\}$   $\pi(q) = \max\{k : k < q, P_k \sqsupset P_q\}$ .

KMP-MATCHER ( $T, P$ )

1.  $n := |T|$
2.  $m := |P|$
3.  $\pi := \text{COMPUTE-PREFIX-FUNCTION}(P)$
4.  $q := 0$
5. for  $i := 1$  to  $n$  do
6.     while  $q > 0$  and  $P[q + 1] \neq T[i]$  do
7.          $q := \pi(q)$
8.     if  $P[q + 1] = T[i]$  then
9.          $q := q + 1$
10.     if  $q = m$  then
11.         print "Podvýraz sa vyskytol s posunom",  $i - m$
12.      $q := \pi(q)$

COMPUTE-PREFIX-FUNCTION ( $P$ )

1.  $m := |P|$
2.  $\pi(1) := 0$
3.  $k := 0$
4. for  $q := 2$  to  $m$  do
5.     while  $k > 0$  and  $P[k + 1] \neq P[q]$  do
6.          $k := \pi(k)$
7.     if  $P[k + 1] = P[q]$  then
8.          $k := k + 1$
9.      $\pi(q) := k$
10. return  $\pi$

**Typy úloh:**

1. Vypočítajte  $\pi$  pre  $P$  a  $\Sigma$
2. Zostavte automat na vyhľadávanie podreťazca  $P$  a znázornite ho.

## Riešené príklady

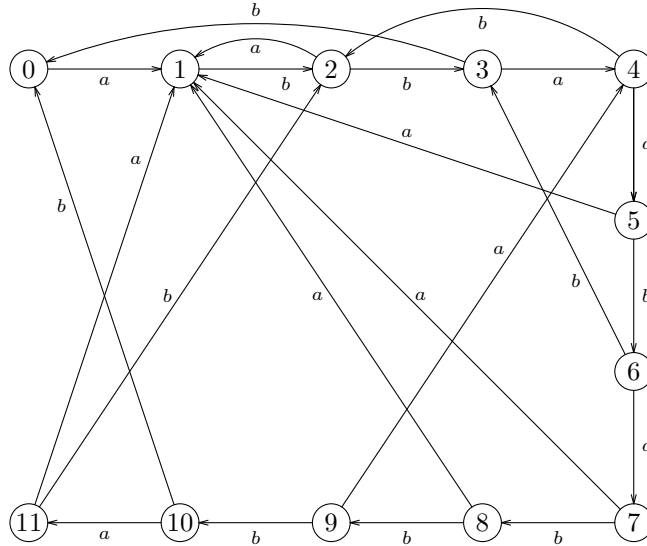
- Zostavte automat na vyhľadávanie reťazca  $P = abbaababba$  a znázornite ho.

$$|P| = 11$$

$$\text{Stavy: } 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, q_0 = 0 \text{ a } A = \{11\}$$

- $\delta(0, a) = \sigma(P_0a) = \sigma(\mathbf{a}) = 1$
- $\delta(0, b) = \sigma(P_0b) = \sigma(b) = 0$
- $\delta(1, a) = \sigma(P_1a) = \sigma(\mathbf{aa}) = 1$
- $\delta(1, b) = \sigma(P_1b) = \sigma(\mathbf{ab}) = 2$
- $\delta(2, a) = \sigma(P_2a) = \sigma(\mathbf{aba}) = 1$
- $\delta(2, b) = \sigma(P_2b) = \sigma(\mathbf{abb}) = 3$
- $\delta(3, a) = \sigma(P_3a) = \sigma(\mathbf{abba}) = 4$
- $\delta(3, b) = \sigma(P_3b) = \sigma(\mathbf{abbb}) = 0$
- $\delta(4, a) = \sigma(P_4a) = \sigma(\mathbf{abbaa}) = 5$
- $\delta(4, b) = \sigma(P_4b) = \sigma(\mathbf{abbab}) = 2$
- $\delta(5, a) = \sigma(P_5a) = \sigma(\mathbf{abbaaa}) = 1$
- $\delta(5, b) = \sigma(P_5b) = \sigma(\mathbf{abbaab}) = 6$
- $\delta(6, a) = \sigma(P_6a) = \sigma(\mathbf{abbaaba}) = 7$
- $\delta(6, b) = \sigma(P_6b) = \sigma(\mathbf{abbaabb}) = 3$
- $\delta(7, a) = \sigma(P_7a) = \sigma(\mathbf{abbaabaa}) = 1$
- $\delta(7, b) = \sigma(P_7b) = \sigma(\mathbf{abbaabab}) = 8$
- $\delta(8, a) = \sigma(P_8a) = \sigma(\mathbf{abbaababa}) = 1$
- $\delta(8, b) = \sigma(P_8b) = \sigma(\mathbf{abbaababb}) = 9$
- $\delta(9, a) = \sigma(P_9a) = \sigma(\mathbf{abbaababba}) = 4$
- $\delta(9, b) = \sigma(P_9b) = \sigma(\mathbf{abbaababbb}) = 10$
- $\delta(10, a) = \sigma(P_{10}a) = \sigma(\mathbf{abbaababbbba}) = 11$
- $\delta(10, b) = \sigma(P_{10}b) = \sigma(\mathbf{abbaababbbb}) = 0$
- $\delta(11, a) = \sigma(P_{11}a) = \sigma(\mathbf{abbaababbbbaa}) = 1$
- $\delta(11, b) = \sigma(P_{11}b) = \sigma(\mathbf{abbaababbbbab}) = 2$

	0	1	2	3	4	5	6	7	8	9	10	11
a	1	1	1	4	5	1	7	1	1	4	11	1
b	0	2	3	0	2	6	3	8	9	10	0	2



- Určte prefixovú funkciu (z KMP algoritmu) pre reťazec  $P$  z predošlej úlohy.

	1	2	3	4	5	6	7	8	9	10	11
$\pi$	0	0	0	1	1	2	1	2	3	0	1

Napr.  $\pi(9) = 3$ :

$P[1..9]$  - *abbaababb*

$P[1..8]$  - *abbaabab* - nie je suffix v  $P[1..9]$

$P[1..7]$  - *abbaaba* - nie je suffix v  $P[1..9]$

$P[1..6]$  - *abbaab* - nie je suffix v  $P[1..9]$

$P[1..5]$  - *abbaa* - nie je suffix v  $P[1..9]$

$P[1..4]$  - *abba* - nie je suffix v  $P[1..9]$

$P[1..3]$  - *abb* - je suffix v  $P[1..9]$