

Úvod

Porovnávanie reťazcov

- Porovnávanie reťazcov, presnejšie hľadanie výskytov reťazca v reťazci (string-matching) je pomerne dôležitou súčasťou širokej domény zaoberajúcej sa spracovaním textu. Algoritmy na porovnávanie textov sa využívajú pri implementácii softvérových systémov, ktoré sú reálne nasadené v praxi. Takisto však hrajú dôležitú rolu v teoretickej informatike, kde môžu byť výzvou pre navrhovanie efektívnejších algoritmov.

1

2

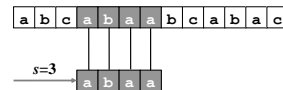
základné pojmy

- $S = \text{"AGCTTGA"}$
- $|S| = 7$, dĺžka reťazca S
- **podreťazec**: $S_{i,j} = S_i S_{i+1} \dots S_j$
 - príklad: $S_{2,4} = \text{"GCT"}$
- **podpostupnosť reťazca** S : vymazaním niekoľkých (vrátane žiadneho) znakov z S
 - "ACT" and "GCTT" sú podpostupnosti.
- **predpona** S : $S_{1,k}$
 - "AGCT" je predpona S .
- **prípona** S : $S_{h,|S|}$
 - "CTTGA" je prípona S .

3

základné pojmy

- Uvažujme 2 reťazce:
 - Vzor $P[1 \dots m]$, ktorý má dĺžku m
 - Text $T[1 \dots n]$, ktorý má dĺžku n
- Vzor P sa nachádza v texte T s posunutím s ak platí:
 - $T[s+1 \dots s+m] = P[1 \dots m]$
- Príklad: $T = \text{abcabaabcbac}$, $P = \text{abaa}$
 - $m=4$, $n=13$, $s=3$



4

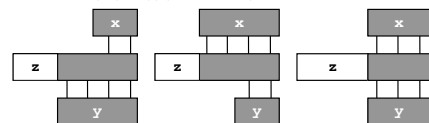
základné pojmy

- **Predpona (prefix)**: reťazec w je predponou reťazca x , ak (existuje reťazec y taký, že) $x = wy$, kde y je akýkoľvek reťazec z použitej abecedy Σ , t.j. prvok z množiny Σ^*
 - Napr: $\text{pre}(\text{ab}, \text{abcca})$
- **Prípona (suffix)**: reťazec w je príponou reťazca x , ak (existuje reťazec y taký, že) $x = yw$, kde y je akýkoľvek reťazec z použitej abecedy Σ , t.j. prvok z množiny Σ^*
 - Napr: $\text{suf}(\text{cca}, \text{abcca})$

5

Lema

- Predpokladajme, že " x ", " y " a " z " sú reťazce, pre ktoré platí $\text{suf}(x, z)$ a $\text{suf}(y, z)$, potom:
 - ak $|x| \leq |y|$, tak $\text{suf}(x, y)$
 - ak $|x| \geq |y|$, tak $\text{suf}(y, x)$
 - ak $|x| = |y|$, tak $x = y$



6

najznámejšie algoritmy

Algoritmus	fáza predspracovania	Vyhľadávacia fáza
Naivný	0	$O((n-m+1)m)$
Rabin-Karp	$\Theta(m)$	$O((n-m+1)m)$
Konečný automat	$O(m \Sigma)$	$\Theta(m)$
Knuth-Morris-Pratt	$\Theta(m)$	$\Theta(m)$

- Ďalšie algoritmy, ich opisy, vizualizácie a zdrojové kódy môžete nájsť napr. na :
<http://www.lgm.univ-mlv.fr/~lecroq/string/>

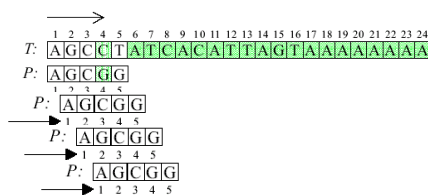
7

Naivné hľadanie výskytu reťazca v reťazci

- Nemá fázu predspracovania
- Vždy sa posúva len o 1 pozíciu doprava
- Porovnávanie môže prebiehať v akomkoľvek poradí
- Veľká časová zložitosť
- Vykoná sa $2n$ porovnávaní textu

8

naivné = hrubou silou



čas: $O(mn)$ kde $m=|P|$ a $n=|T|$.

9

Naivné hľadanie výskytu reťazca v reťazci

- Samotný algoritmus pozostáva z porovnávanie znakov na všetkých miestach medzi 0 a $n-m$. Pri každom kroku sa posúva iba o 1 miesto doprava.

Naivne porovnavanie (T, P)

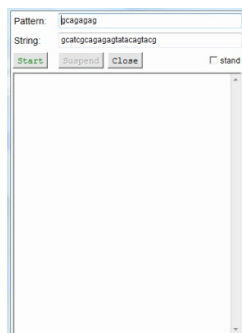
```

n ← length[T]
m ← length[P]
for s ← 0 to n - m
    do if P[1 . . m] = T[s+1 . . s+m]
        then
            print "vyskytuje sa s posunutim" s

```

10

Naivné hľadanie výskytu reťazca v reťazci



11

Naivné hľadanie výskytu reťazca v reťazci

C++:

```

void BF(char *x, int m, char *y, int n)
{
    int i, j;
    /* Searching */
    for (j = 0; j <= n - m; ++j)
    {
        for (i = 0; i < m && x[i] == y[i + j]; ++i); if (i >= m)
            OUTPUT(j);
    }
}

```

12

2 druhy algoritmov hľadania výskytu reťazca

- predspracovanie vzoru P (P sa nemení, T sa mení)
 - napr: dopyt P do databázy T
 - algoritmy:
 - Knuth – Morris – Pratt
 - Boyer – Moore
- predspracovanie textu T (T sa nemení, P sa mení)
 - napr: hľadá sa vzor P v slovníku T
 - algoritmus:
 - príponový strom

13

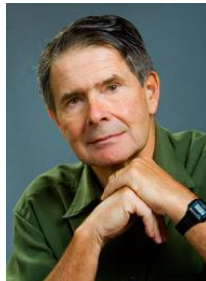
Predspracovanie vzoru

- dvojfázové algoritmy
 - fáza 1 : vygeneruj pole, ktoré bude indikovať smer pohybu.
 - fáza 2 : použi to pole na pohyb a hľadanie výskytu
- príklady
 - [algoritmus KMP](#):
 - navrhli Knuth, Morris a Pratt v 1977.
 - [algoritmus BM](#):
 - navrhli Boyer a Moore v 1977.

14

James H. Morris

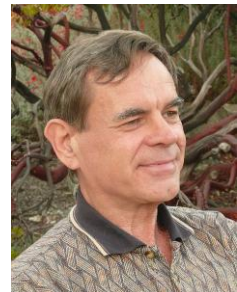
- 1941, Pittsburgh –
- Bc CMU
- Master manažment, MIT
- PhD informatika, MIT
- profesor, CMU
- práce
 - Donald Knuth, James H. Morris, Jr., and Vaughan Pratt. Fast pattern matching in strings. SIAM Journal on Computing, 6(2):323–350. 1977.



15

Vaughan Pratt

- 1944 –
- emeritný profesor
- Stanfordova univerzita
- práce
 - Donald Knuth, James H. Morris, Jr., and Vaughan Pratt. Fast pattern matching in strings. SIAM Journal on Computing, 6(2):323–350. 1977.



16

Knuthov-Morrisov-Prattov algoritmus

- KMP algoritmus vychádza z analýzy algoritmu naivného vyhľadávania. V určitých situáciách vie využiť informáciu získanú čiastočným porovnávaním vybraného podreťazca a vzoru a posunúť podreťazec o viac než jeden znak.

prvý prípad v KMP algoritme

- prvý symbol vzoru P sa viac vo vzore T nenachádza.
- možno sa posunúť až ku T_k , pretože $T_k \neq P_k$ (a $T_i = P_i$ pre prvé tri pozície $i=1,2,3$) v (a).

→
 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
 T: AGCCTATTCACATTTAGTAAAAA
 P: AGCCTG
 1 2 3 4 5

(a)

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
 T: AGCCTATTCACATTTAGTAAAAA
 P: AGCCTG
 → 1 2 3 4 5

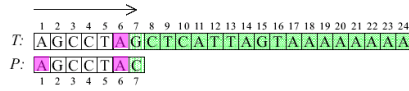
(b)

17

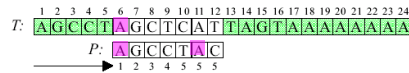
18

druhý případ v KMP algoritme

- prvý symbol vo vzore P sa v ňom ešte nachádza na niektorom ďalšom mieste.
- $T_7 \neq P_7$ v (a). treba sa posunúť ku T_6 , pretože $P_6 = P_1 = T_6$.



(a)

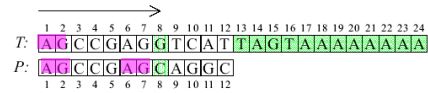


(b)

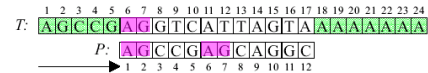
19

tretí prípad v KMP algoritme

- predpona v zoru P sa vyskytuje v P ešte raz.
- $T_8 \neq P_8$ v (a). Treba sa posunúť k T_6 , pretože $P_{6,7} = P_{1,2} = T_{6,7}$.



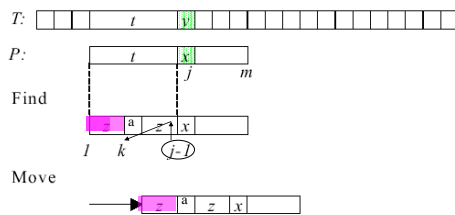
(a)



(b)

20

základná myšlienka KMP algoritmu



21

Knuth-Morris-Pratt príklad

$t[0]$	$t[1]$	$t[2]$	$t[3]$	$t[4]$	$t[5]$	$t[6]$	$t[7]$	$t[8]$	$t[9]$	$t[10]$	$t[11]$	$t[12]$	$t[13]$
A	B	C		A	B	C	D	A	B		A	B	C

$p[0]$	$p[1]$	$p[2]$	$p[3]$	$p[4]$	$p[5]$	$p[6]$
A	B	C	D	A	B	D

Y Y Y N

$$m = 0$$

22

Knuth-Morris-Pratt príklad

t[0]	t[1]	t[2]	t[3]	t[4]	t[5]	t[6]	t[7]	t[8]	t[9]	t[10]	t[11]	t[12]	t[13]
A	B	C		A	B	C	D	A	B		A	B	C

$p[0]$	$p[1]$	$p[2]$	$p[3]$	$p[4]$	$p[5]$	$p[6]$
A	B	C	D	A	B	D

Y Y Y Y Y Y N

$$m = 4$$

23

Knuth-Morris-Pratt príklad

t[0]	t[1]	t[2]	t[3]	t[4]	t[5]	t[6]	t[7]	t[8]	t[9]	t[10]	t[11]	t[12]	t[13]
A	B	C		A	B	C	D	A	B		A	B	C

$p[0]$	$p[1]$	$p[2]$	$p[3]$	$p[4]$	$p[5]$	$p[6]$
A	B	C	D	A	B	D

N

$m = 10$

24

Knuth-Morris-Pratt príklad

$t[0]$	$t[1]$	$t[2]$	$t[3]$	$t[4]$	$t[5]$	$t[6]$	$t[7]$	$t[8]$	$t[9]$	$t[10]$	$t[11]$	$t[12]$	$t[13]$
A	B	C		A	B	C	D	A	B		A	B	C

$p[0]$	$p[1]$	$p[2]$..
A	B	C	..
Y	Y	Y	

$m = 11$

25

Knuthov-Morrisov-Prattov algoritmus

algorithm *kmp_search*:

Input pole znakov T (text, v ktorom sa hľadá výskyt vzoru)
pole znakov P (text, vzor, ktorého výskyt sa hľadá)

Output celé číslo (miesto v poli T, začínajúce ktorým sa našiel výskyt vzoru P)

define variables:

integer $m \leftarrow 0$ (začiatočné miesto súčasného výskytu v T)

integer $i \leftarrow 0$ (miesto súčasného znaku v P)

pole celých čísel π (tabuľka, ktorá sa vypočíta inde)

while $m + i < \text{dĺžka } \pi$ **do**:

if $P[i] = T[m + i]$ **then**

let $i \leftarrow i + 1$

if $i = \text{dĺžka } P$ **then** return m

else

let $m \leftarrow m + i - \pi[i]$,

if $i > 0$ **then** let $i \leftarrow \pi[i]$

26

KMP algoritmus – pomocná tabuľka

algorithm *kmp_table*:

input: pole znakov P (vzor, ktorého výskyt sa hľadá)

output: pole celých čísel π (tabuľka, ktorú treba vyplniť)

define variables:

integer $i \leftarrow 2$ (súčasná hodnota v π , pre ktorú sa počíta hodnota π)

integer $j \leftarrow 0$ (pozícia znaku v P, počítajúca od nuly, ktorý je nasledujúcim znakom v súčasnom podreťazci, ktorý je kandidátom na zhodu)

let $\pi[0] \leftarrow -1$, $\pi[1] \leftarrow 0$

while $i < \text{dĺžka } P$ **do**:

(prvý prípad: podreťazec pokračuje)

if $P[i - 1] = P[j]$

then let $\pi[i] \leftarrow j + 1$, $i \leftarrow i + 1$, $j \leftarrow j + 1$

(druhý prípad: nepokračuje, ale možno sa vrátiť)

else if $j > 0$ then let $j \leftarrow \pi[j]$

(tretí prípad: nie sú ďalší kandidáti. Všimnime si, že $j = 0$)

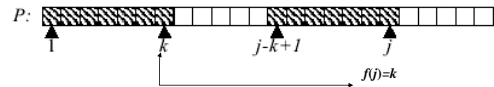
else let $\pi[i] \leftarrow 0$, $i \leftarrow i + 1$

27

definícia funkcie počítajúcej predponu

$f(j)$ = najväčšie $k < j$ také, že $P_{1,k} = P_{j-k+1,j}$

$f(j) = 0$ ak také k neexistuje

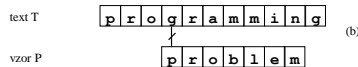
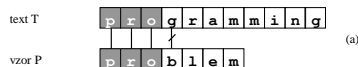


j	1	2	3	4	5	6	7	8	9	10	11	12
P :	A	T	C	A	C	A	T	C	A	T	C	A
f :	0	0	0	1	0	1	2	3	4	2	3	4

28

Knuthov-Morrisov-Prattov algoritmus

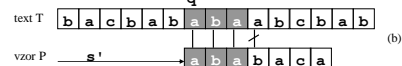
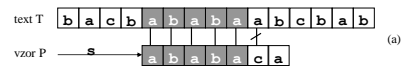
Porovnávanie vzoru s reťazcom začína na prvom znaku zľava (vzor je zarovnaný s reťazcom). Algoritmus postupuje, kým nenarazí na nezhodu na štvrtej pozícii medzi znakmi **b** a **g** (obr. a). Z predchádzajúcich znakov okamžite vieme, že posun vzoru o jeden alebo dva znaky nemá význam. Preto nastane posun o tri znaky. Tým sa vzor zarovná s textom nad znakom, kde nastala nezhoda. Od tohto miesta môže ďalej pokračovať porovnávanie.



29

Knuthov-Morrisov-Prattov algoritmus

- V tomto príklade vidíme, že vzor je posunutý o **s** a nastáva nové porovnávanie. Pri ňom sa zistilo, že nezhoda nastala na 6. pozícii reťazca, čo indikuje posun o 5 znakov (**g**). V tomto prípade však takýto posun nie je možný. Posunúť sa môžeme iba o 2 znaky, pretože na 3. znaku sa nachádza zhoda medzi týmto znakom a prvým znakom vzoru (na tomto mieste môže začínať vzor)



30

Knuthov-Morrisov-Prattov algoritmus

- Posun pri prehľadávaní je nezávislý od prehľadávaného reťazca. Jeho veľkosť určuje tzv. predponová funkcia.
- predponová funkcia (Prefix function)
 π pre P , $|P| = m$:
 $\pi : \{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, m-1\}$
 $\pi(q) = \max\{k : k < q, P_{1,k} \supseteq P_{q-k+1,q}\}.$
- Jednotlivé posuny sa vypočítavajú vo fáze predspracovania.

31

Knuthov-Morrisov-Prattov algoritmus

```

KMP POROVNANIE(T, P)
n ← dĺžka(T)
m ← dĺžka(P)
π ← PREDPONOVÁ FUNKCIA(P)
q ← 0
for i ← 1 to n //prehľadávaj ďalší zľava doprava
do while q > 0 and P[q+1] ≠ T[i]
do q ← π[q] //nehodzuje sa, bude „ústup“
if P[q+1] = T[i]
then q ← q + 1 //zhoduje sa
if q = m
then print "Vzor sa v reťazci vyskytuje s posunom " i-m
q ← π[q]

```

32

Knuthov-Morrisov-Prattov algoritmus

```

PREDPONOVÁ FUNKCIA (P)
// funkcia sa reprezentuje tabuľkou π, jej
// hodnoty sa tu vypočítavajú a zapisujú do π
m ← dĺžka(P)
π[1] ← 0
k ← 0
for q ← 2 to m
do while k > 0 and P[k+1] ≠ P[q]
do k ← π[k]
if P[k+1] = P[q]
then k ← k + 1
π[q] ← k
return π

```

33

príklad: vypočítať π pre vzor 'P':

p a b a b a c a

na začiatku: m = dĺžka[P] = 7
 $\pi[1] = 0$
 $k = 0$

krok 1: q = 2, k = 0
 $\pi[2] = 0$

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0					

krok 2: q = 3, k = 0,
 $\pi[3] = 1$

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1				

krok 3: q = 4, k = 1
 $\pi[4] = 2$

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1	2			

krok 4: q = 5, k = 2
 $\pi[5] = 3$

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1	2	3		

krok 5: q = 6, k = 3
 $\pi[6] = 1$

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1	2	3	1	

krok 6: q = 7, k = 1
 $\pi[7] = 1$

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1	2	3	1	1

po 6 opakovaníach je výpočet funkcie
 π úplný: →

q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1	2	3	1	1

príklad: nech sú dané reťazec 'S' a vzor 'P':

S b a c b a b a b a b a c a c a

P a b a b a c a

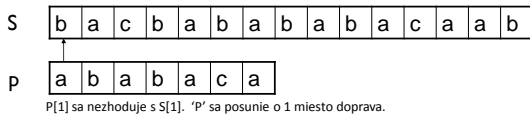
Na nájdenie výskytu vzoru p v texte S použijeme algoritmus KMP.

nojprv sa pre vzor 'P' vypočíta predponová funkcia π daná tabuľkou (pozri prechádzajúci príklad)

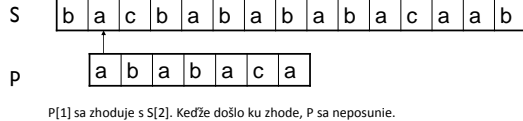
q	1	2	3	4	5	6	7
p	a	b	a	b	a	c	a
π	0	0	1	2	3	1	1

Na začiatku: $n = \text{dĺžka}(S) = 15$;
 $m = \text{dĺžka}(P) = 7$

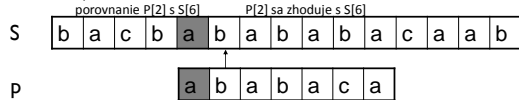
krok 1: $i = 1, q = 0$
 porovnanie $P[1]$ s $S[1]$



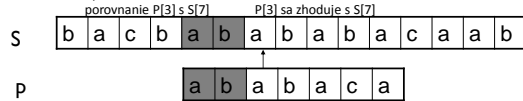
krok 2: $i = 2, q = 0$
 porovnanie $P[1]$ s $S[2]$



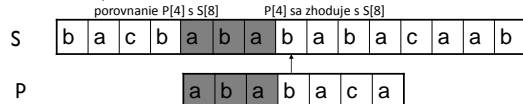
krok 6: $i = 6, q = 1$



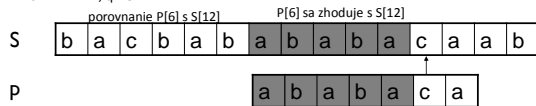
krok 7: $i = 7, q = 2$



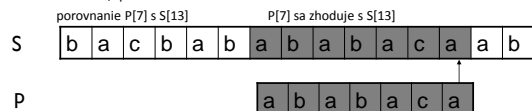
krok 8: $i = 8, q = 3$



krok 12: $i = 12, q = 5$

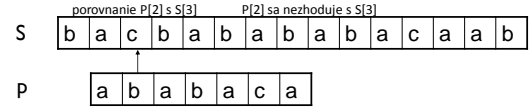


krok 13: $i = 13, q = 6$

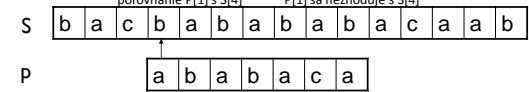


Vzor 'P' sme našli (jeho úplný výskyt) v texte 'S'. Celkový počet posunov, ktoré sa vykonali, aby sa našiel výskyt je: $i - m = 13 - 7 = 6$.

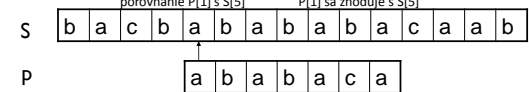
krok 3: $i = 3, q = 1$



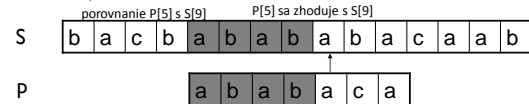
krok 4: $i = 4, q = 0$



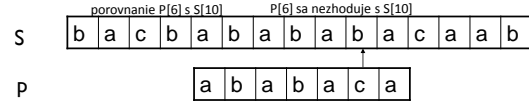
krok 5: $i = 5, q = 0$



krok 9: $i = 9, q = 4$

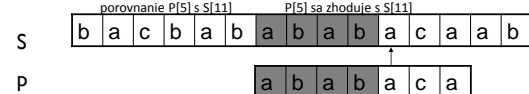


krok 10: $i = 10, q = 5$



ústup v P, porovnanie $P[4]$ s $S[10]$, preto lebo po nezhode $q = \lceil m/5 \rceil + 3$

krok 11: $i = 11, q = 4$



odhad času výpočtu

• PREDPONOVÁ FUNKCIA (P)

```

1 m ← dĺžka[p] // p' vzor
2 m[1] ← 0
3 k ← 0
4 for q ← 2 to m
5   do while k > 0 and p[k+1] != p[q]
6     do k ← m[k]
7     If p[k+1] = p[q]
8       then k ← k + 1
9   m[q] ← k
10 return m

```

Cyklus for v riadkoch 4 až 10 sa vykoná 'm' ráz. Kroky 1 až 3 sa vykonávajú v konštantnom čase. Preto odhad času výpočtu predponovej funkcie je $O(m)$.

• KMP POROVNANIE(T, P)

```

1 n ← dĺžka[T]
2 m ← dĺžka[p]
3 m ← predponová funkcia(p)
4 q ← 0
5 for i ← 1 to n
6   do while q > 0 and p[q+1] != T[i]
7     do q ← m[q]
8   if p[q+1] = T[i]
9     then q ← q + 1
10  if q = m
11    then print "Vzor sa v retazci vyskytuje s posunom" i - m
12    q ← m[q]

```

Cyklus začínajúci v riadku 5 sa vykoná 'n' ráz, t.j. toľko ráz, ako je dlhý text T. Keďže kroky 1 až 4 sa vykonávajú v konštantnom čase, o celkovom čase výpočtu rozhoduje cyklus. Preto odhad času výpočtu KMP porovnania je $O(n)$.

KMP algoritmus – časová zložitosť

- časová zložitosť: $O(m+n)$
 - $O(m)$ výpočet funkcie f
 - $O(n)$ hľadanie vzoru P

43

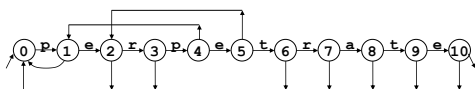
KMP algoritmus a konečný automat

- KMP algoritmus do určitej miery súvisí s konečnými automatmi. Predpokladajme, že máme vzor P dĺžky m . Definujeme si konečný automat, ktorý bude mať $m+1$ stavov. Prechody medzi jednotlivými stavmi budú postupne určené jednotlivými písmenami vzoru. Teda napr. prechod medzi nulovým a prvým stavom bude podľa písmena p_1 , prechod medzi prvým a druhým stavom podľa p_2 atď. Ostatné prechody (teda akési chybové) bude určovať práve predponová funkcia. Vstupným stavom bude stav 0 a výstupným stav m . Samotné vyhľadávanie bude realizované ako práca takéhoto automatu so vstupom, ktorý odpovedá zadanému reťazcu. Rozdiel je iba v tom, že pokiaľ sa pomocou predponovej funkcie vrátíme do niektorého predchádzajúceho stavu, okamžite skúsime cez ten istý znak (ktorý spôsobil nezhodu) prejsť do nasledujúceho stavu.

44

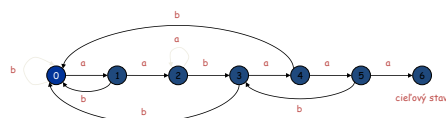
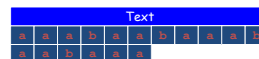
KMP algoritmus a konečný automat

- Príklad konečného automatu pre vzor `perpetrate`



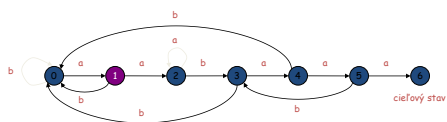
45

Knuth-Morris-Pratt

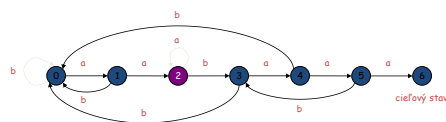
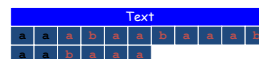


Knuth-Morris-Pratt

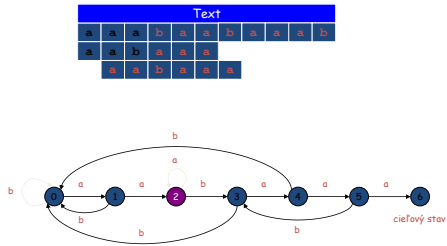
• .



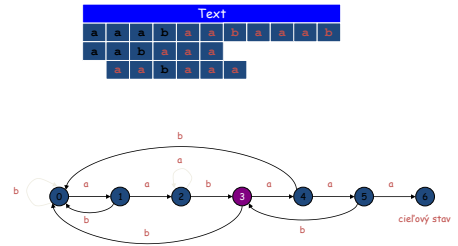
Knuth-Morris-Pratt



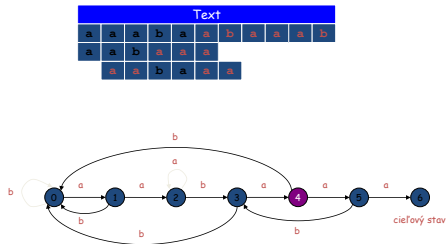
Knuth-Morris-Pratt



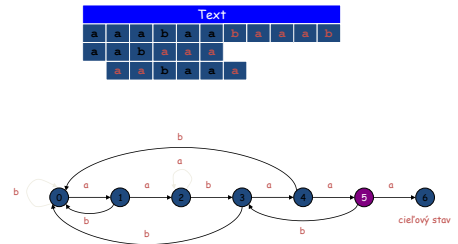
Knuth-Morris-Pratt



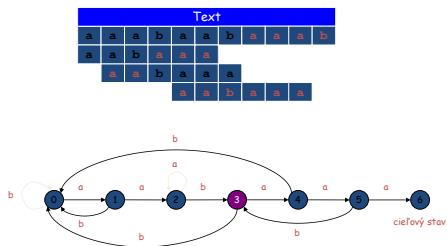
Knuth-Morris-Pratt



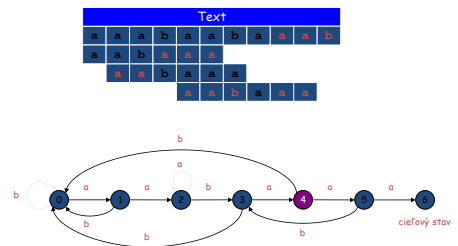
Knuth-Morris-Pratt



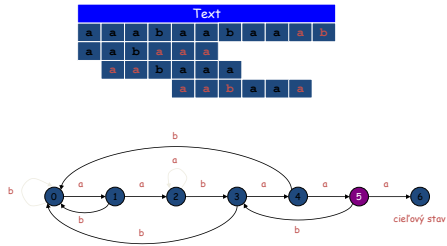
Knuth-Morris-Pratt



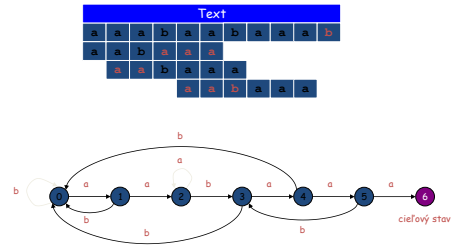
Knuth-Morris-Pratt



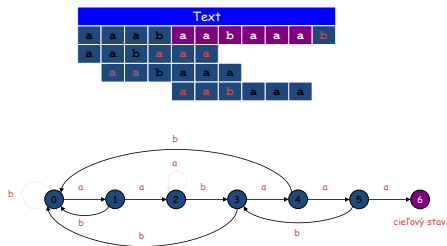
Knuth-Morris-Pratt



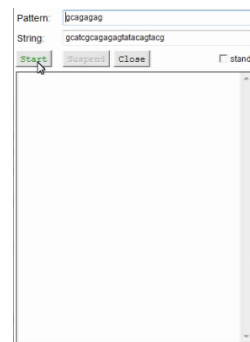
Knuth-Morris-Pratt



Knuth-Morris-Pratt



Knuthov-Morrisov-Prattov algoritmus



58

Knuthov-Morrisov-Prattov algoritmus -
príklad implementácie

```
void preKmp(char *x, int m, int kmpNext[])
{
    int i, j; i = 0; j = kmpNext[0] = -1;
    while (i < m)
    {
        while (j > -1 && x[i] != x[j])
            j = kmpNext[j];
        i++;
        j++;
        if (x[i] == x[j])
            kmpNext[i] = kmpNext[j];
        else
            kmpNext[i] = j;
    }
}
```

59

Knuthov-Morrisov-Prattov algoritmus -
príklad implementácie

```
void KMP(char *x, int m, char *y, int n)
{
    int i, j, kmpNext[XSIZE];
    /* Preprocessing */
    preKmp(x, m, kmpNext);
    /* Searching */
    i = 0;
    while (i < n)
    {
        while (i > -1 && x[i] != y[i])
            i = kmpNext[i];
        i++;
        j++;
        if (i >= m)
        {
            OUTPUT(j - i);
            i = kmpNext[i];
        }
    }
}
```

60

Robert S. Boyer

- emeritný profesor
- Computer Science Department
- The University of Texas at Austin
- práce:
 - BOYER R.S., MOORE J.S., 1977, A fast string searching algorithm. Communications of the ACM. 20:762-772.



61

J Strother Moore

- profesor
- Computer Science Department
- The University of Texas at Austin
- práce:
 - BOYER R.S., MOORE J.S., 1977, A fast string searching algorithm. Communications of the ACM. 20:762-772.
 - dokazovanie teorém



62

Boyerov-Moorov algoritmus - príklad

t[0]	t[1]	t[2]	t[3]	t[4]	t[5]	t[6]	t[7]	t[8]	t[9]	t[10]
A	B	C	E	F	G	A	B	C	D	E

p[0]	p[1]	p[2]	p[3]
A	B	C	D

N

Vo vzore nie je žiadny znak E: z toho plynie, že vzor nemožno nájsť nikde pred znakom t[3]. Preto sa možno a teda aj treba posunúť o 4 miesta doprava.

63

Boyerov-Moorov algoritmus - príklad

t[0]	t[1]	t[2]	t[3]	t[4]	t[5]	t[6]	t[7]	t[8]	t[9]	t[10]
A	B	C	E	F	G	A	B	C	D	E

p[0]	p[1]	p[2]	p[3]
A	B	C	D

N

Nevyskytuje sa. Ale vo vzore sa B vyskytuje. Treba sa posunúť o dve miesta doprava.

64

Boyerov-Moorov algoritmus - príklad

t[0]	t[1]	t[2]	t[3]	t[4]	t[5]	t[6]	t[7]	t[8]	t[9]	t[10]
A	B	C	E	F	G	A	B	C	D	E

p[0]	p[1]	p[2]	p[3]
A	B	C	D

Y Y Y Y

65

Boyerov-Moorov algoritmus

- porovnáva sprava doľava
- 2 predvypočítané funkcie
 - posun o dobrú príponu
 - posun o zlý znak

66

Boyerov-Moorov algoritmus

myšlienka č. 1: porovnávať sprava doľava

```

12345678901234567
T:  xpbctbxbabpqxctbpq
P:  tpabxab

```

67

Boyerov-Moorov algoritmus

```

12345678901234567
T:  spbctbsabpqscctbpq
P:  tpabsab
P:  tpabxab

```

myšlienka č. 2: pravidlo zlého znaku

$R(x)$: najpravejší výskyt znaku x v P . $R(x)=0$ ak sa x nevyskytuje. $R(t)=1$, $R(s)=5$.

i : miesto, kde došlo k nezhode v P . $i=3$

k : zodpovedajúce miesto v T . $k=5$. $T[k]=t$

pravidlo zlého znaku: Vzor P sa posunie doprava o $\max\{1, i-R(T[k])\}$. čiže ak je najpravejší výskyt znaku $T[k]$ v P na mieste j ($j < i$), tak $P[j]$ sa očitne pod $T[k]$ po posunutí.

68

Boyerov-Moorov algoritmus

- myšlienkou pravidla zlého znaku je posunúť P o viac než o jeden znak, ak je to možné.
- pravidlo neúčinné, ak $j > i$
- bohužiaľ, prípad $j > i$ je častý

```

12345678901234567
T:  spbctbsatpqscctbpq
P:  tpabsat
P:  tpabsat

```

69

Boyerov-Moorov algoritmus

označme $x=T[k]$ znak, ktorý sa nezhoduje so vzorom v T .

myšlienka č. 3: **rozšírené pravidlo zlého znaku**: P sa posunie doprava tak, že najbližšie x vľavo od miesta i v P bude pod $T[k]$.

```

12345678901234567
T:  spbctbsatpqscctbpq
P:  tpabsat
P:  tpabsat

```

70

Boyerov-Moorov algoritmus

aby sme mohli používať rozšírené pravidlo zlého znaku, potrebujeme poznať:
pre každé miesto i v P , pre každý znak x v abecede, miesto najbližšieho výskytu x vľavo od miesta i .

možný prístup: dvojrozmerným poľom: $n * |\Sigma|$

pamäťovo a časovo príliš náročné

71

Boyerov-Moorov algoritmus

iný prístup: prezrieť P sprava doľava a pre každý znak x si držať zoznam miest, kde sa x vyskytuje (v klesajúcom poradí).

```

P:  tpabsat
t → 7, 1
a → 6, 3 ...

```

ak dôjde k nezhode medzi $P[i]$ a $T[k]$, (označme $x=T[k]$), prezri zoznam výskytov znaku x , nájdi prvé číslo (označme ho j), ktoré je menšie než i a posuň P doprava tak, aby $P[j]$ sa dostalo pod $T[k]$.

```

12345678901234567
T:  spbctbsatpqscctbpq
P:  tpabsat
P:  tpabsat

```

ak také číslo j neexistuje, tak posuň P za $T[k]$
čas a pamäť: lineárna

72

Boyerov-Moorov algoritmus

myšlienka č. 4: pravidlo dobrej prípony



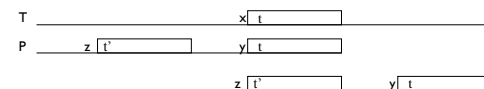
t je prípona vzoru P taká, že sa zhoduje s podreťazcom t textu T a $x \neq y$
 t' je najpravejší výskyt t v P taký, že t' nie je príponou P a $z \neq y$

73

Boyerov-Moorov algoritmus

pravidlo dobrej prípony:

(1) ak existuje t' tak posuň P doprava tak, že t' v P je pod t v T



123456789012345678
T: prstabst**ub**abvqxr**st**
P: qcab**d**ab**d**ab
P: qcab**d**ab**d**ab
P: qcab**d**ab**d**ab

74

Boyerov-Moorov algoritmus

- rozšírené pravidlo zlého znaku sa sústreďuje na znaky.
- pravidlo dobrej prípony sa sústreďuje na podreťazce.
- ako teraz získať informáciu, ktorú treba pre pravidlo dobrej prípony? ako pre nejaké t nájsť t' ?

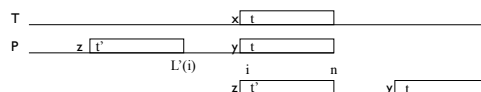
75

Boyerov-Moorov algoritmus

$L'(i)$: pre každé i, $L'(i)$ je najväčšie miesto menšie než n také, že podreťazec $P[i, \dots, n]$ sa zhoduje s príponou reťazca $P[1, \dots, L'(i)]$ a navyše znak predchádzajúci tejto prípony je rôzny od $P[i-1]$.

ak také miesto neexistuje tak $L'(i) = 0$.

nech $t = P[i, \dots, n]$, potom $L'(i)$ je miesto pravého konca t' .



T: prstabst**ub**abvqxr**st**
P: qcab**d**ab**d**ab
 1234567890
 $L'(9) = 4, L'(10) = 0, L'(8) = ?, L'(7) = ? L'(6) = ?$

76

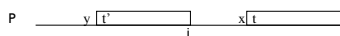
Boyerov-Moorov algoritmus

nech $t = P[i, \dots, n]$, potom $L'(i)$ je miesto pravého konca t' .

aby sa dalo používať pravidlo dobrej prípony, treba poznať $L'(i)$ pre všetky $i = 1, \dots, n$.

pre vzor P:

N_i je dĺžka najdlhšieho podreťazca, ktorý končí na mieste j a ktorý je príponou P.

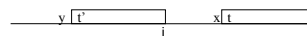


$t = t'$
 $j = |t'| = |t|$
 $x \neq y$

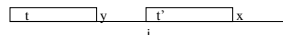
77

Boyerov-Moorov algoritmus

N_i : dĺžka najdlhšieho podreťazca, ktorý končí na mieste j a ktorý je príponou P.



Z_i : dĺžka najdlhšieho podreťazca, ktorý začína na mieste i a ktorý je predponou P.



78

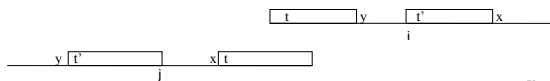
Boyerov-Moorov algoritmus

N je obrátené Z

P' je obrátený reťazec P

všimnime si, že $N_i(P) = Z_{n-i+1}(P')$

1 2 3 4 5 6 7 8 9 0	1 2 3 4 5 6 7 8 9 0
P: q c a b d a b d a b	P': b a d b a d b a c q
N _i : 0 0 0 2 0 0 5 0 0 0	Z _i : 0 0 0 5 0 0 2 0 0 0



79

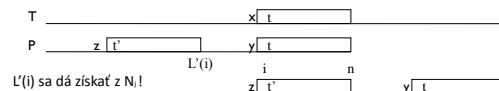
Boyerov-Moorov algoritmus

pre daný vzor P ,

N_i (pre $j=1, \dots, n$) sa dá vypočítať v lineárnom čase $O(n)$ algoritmom Z .

prečo treba určiť N_i ?

aby sa dalo použiť pravidlo dobrej prípony, treba poznať $L'(i)$ pre všetky $i=1, \dots, n$.

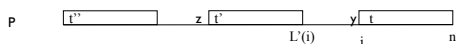


80

Boyerov-Moorov algoritmus

Pre miesto i , nech $t = P[i, \dots, n]$.

$L'(i)$ je najväčšie miesto j menšie než n také, že $N_i = |t|$



1 2 3 4 5 6 7 8 9 0	1 2 3 4 5 6 7 8 9 0
P: q c a b d a b d a b	P': b a d b a d b a c q
N _i : 0 0 0 2 0 0 5 0 0 0	Z _i : 0 0 0 5 0 0 2 0 0 0
L'(i): 0 0 0 0 0 7 0 0 4 0	

81

Boyerov-Moorov algoritmus

Ako získať $L'(i)$ z N_i v lineárnom čase?

vstup: vzor P

výstup: $L'(i)$ pre $i=1, \dots, n$

Algoritmus

vypočítaj N_i pre $j=1, \dots, n$ na základe algoritmu Z

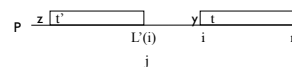
for $i=1; i \leq n; i++$

$L'(i)=0;$

for $j=1; j \leq n; j++$

$i = n - N_j + 1$

$L'(i)=j;$

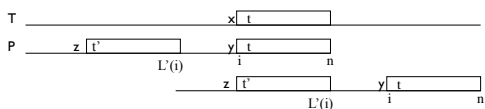


82

Boyerov-Moorov algoritmus

Pravidlo dobrej prípony:

(1) ak existuje t' , tak posuň P doprava takým spôsobom, že t' v P je pod t v T



123456789012345678	
T: prstabstubbavqxrst	
P: qcabdabdb	i=9; L'(9)=4
P: qcabdabdb	

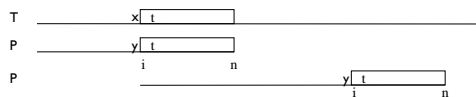
83

Boyerov-Moorov algoritmus

Pravidlo dobrej prípony:

(1) Ak nastáva nezhoda na mieste $i-1$ v P a $L'(i) > 0$ (t.j. existuje t'), tak podľa pravidla dobrej prípony možno posunúť P o $n - L'(i)$ miest doprava.

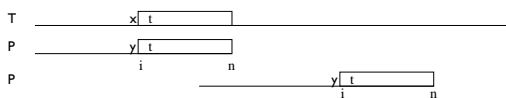
(2) Čo ak nastane nezhoda na mieste $i-1$ v P a $L'(i) = 0$ (t.j. t' neexistuje)? Možno P posunúť takto



84

Boyerov-Moorov algoritmus

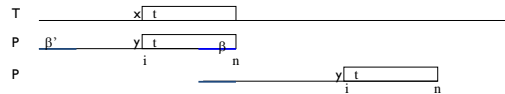
Dá sa však spraviť viac!



85

Boyerov-Moorov algoritmus

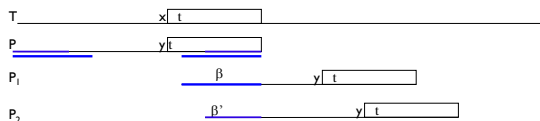
Pozorovanie 1: ak β je predponou P a je tiež príponou P, tak...



86

Boyerov-Moorov algoritmus

pozorovanie 2: ak je viac kandidátov β , tak posuň P o čo najmenšiu dĺžku

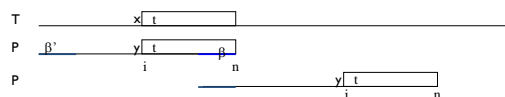


87

Boyerov-Moorov algoritmus

Pravidlo dobrej prípony: keď nastane nezhoda na mieste $i-1$ vzoru P

- (1) ak $L'(i) > 0$ (t.j. t' existuje), tak podľa pravidla dobrej prípony možno posunúť P o $n-L'(i)$ miest doprava.
- (2) inak ak $L'(i) = 0$ (t.j. t' neexistuje)? Možno posunúť P poza ľavý koniec t o najmenší počet miest taký, že predpona posunutého vzoru sa zhoduje s príponou t.

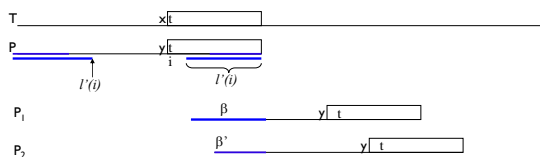


88

Boyerov-Moorov algoritmus

$L'(i)$: dĺžka najväčšej prípony $P[i..n]$ takej, že je aj predponou vzoru P. Ak taká neexistuje, tak $L'(i) = 0$.

$L'(i)$ je dĺžka prekrytia medzi neposunutým a posunutým vzorom.

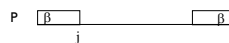


89

Boyerov-Moorov algoritmus

$L'(i)$ sa rovná najväčšiemu $j \leq |P[i..n]|$ takému, že $N=j$

1. $N=j$ tak β je predponou P a tiež príponou P



2. a chceme najväčšie j



90

Boyerov-Moorov algoritmus

$L'(i)$ sa rovná najväčšiemu $j \leq |P[i, \dots n]|$ takému, že $N_j = i$

	1	2	3	4	5	6	7	8	9	0
P:	a	b	d	a	b	a	b	d	a	b
N _b :	0	2	0	0	5	0	2	0	0	0
L'(i):	5	5	5	5	5	5	2	2	2	0

91

Boyerov-Moorov algoritmus

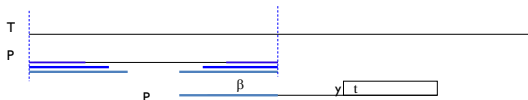
Ako vypočítať $L'(i)$ z N_j v lineárnom čase?

92

Boyerov-Moorov algoritmus

Čo ak sa nájde zhoda?

posuň P o 1 miesto...ale...



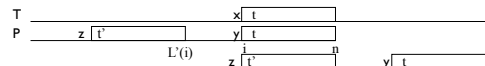
Posuň P o najmenší počet miest taký, že predpona posúvaného vzoru sa zhoduje s t, t.j. posuň P doprava o $n - L'(i)$

93

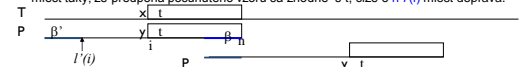
Boyerov-Moorov algoritmus

Pravidlo dobrej prípony: keď nastane nezhoda na mieste $i-1$ vzoru P

(1) ak $L'(i) > 0$ (t.j. t' existuje), tak podľa pravidla dobrej prípony možno posunúť P o $n - L'(i)$ miest doprava.



(2) Inak ak $L'(i) = 0$ (t.j. t' neexistuje)? Možno posunúť P po ľavý koniec t o najmenší počet miest taký, že predpona posunutého vzoru sa zhoduje s t, čiže o $n - f(i)$ miest doprava.



(3) Ak sa nájde nezhoda, tak posuň P doprava o $n - f(2)$

94

Boyerov-Moorov algoritmus

rozšírené pravidlo zlého znaku vs. pravidlo dobrej prípony

123456789012345678	T: prstabstuvabvqxrst	123456789012345678	T: prstabstuvabvqxrst
P: qcabd dab	P: qcabd dab	P: qcabd dab	P: qcabd dab
P: qcabdabdad	P: qcabdabdad	P: qcabdabdad	P: qcabdabdad
P: qcabdabdad	P: qcabdabdad	P: qcabdabdad	P: qcabdabdad

95

Boyerov-Moorov algoritmus

Posuň P o najväčší počet miest určený niektorým z oboch pravidiel. To je podstata Boyerovho-Moorovho algoritmu!

vstup: text T, vzor P; výstup: nájdi výskyt vzoru P v T

Algoritmus **Boyer-Moore**

vypočítaj $L'(i)$, $L'(i)$ a $R(x)$

$k = n$;

while ($k \leq m$) do

$i = n$

$h = k$

 while $i > 0$ and $P[i] = T[h]$ do

$i--$;

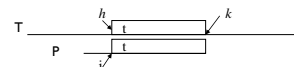
$h--$;

 if $i = 0$

 oznám výskyt vzoru P v T končiaci na mieste k;

$k = k + n - f(2)$

 else posuň P (zvýš k) o väčší počet miest z počtu určeného rozšíreným pravidlom zlého znaku a počtu určeného pravidlom dobrej prípony.



96

Boyerov-Moorov algoritmus

- výkonnosť závisí od dĺžky vzoru
- $O(n/m)$
- Dlhšie vzory = lepšia výkonnosť
- Najmenší vzor = $m = 1$
 $O(n)$ – lineárne hľadanie

97

porovnávanie pomocou automatu

- Na základe vzoru sa vytvorí minimálny deterministický konečný automat, pomocou ktorého sa rozpoznáva vzor v zadanom reťazci.

98

porovnávanie pomocou automatu

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)
- Σ je vstupná abeceda
- δ je tzv. prechodová funkcia z $Q \times \Sigma$ do Q .

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

$\delta^*(q, \epsilon) = q$

$\delta^*(q, wa) = \delta(\delta^*(q, w), a)$

99

porovnávanie pomocou automatu

funkcia koncového stavu - vracia stav automatu po spracovaní nejakého slova

príponová 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 príponou v a $P_k = P[1..k]$.

100

porovnávanie pomocou automatu

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 dostane automat do stavu m , našiel sa podvýraz P a jeho posun je daný aktuálnym miestom v reťazci zmenšeným o m .

Platí:

- Pre každý reťazec x a znak a platí: $\sigma(xa) \leq \sigma(x) + 1$.
- Pre každý reťazec x a znak a , ak $q = \sigma(x)$, tak $\sigma(xa) = \sigma(P_q a)$

101

porovnávanie pomocou automatu

POROVNANIE KONEČNÝM AUTOMATOM(T, δ, m)

$n \leftarrow \text{dĺžka}(T)$

$q \leftarrow 0$

for $i \leftarrow 1$ **to** n

do $q \leftarrow \delta(q, T[i])$

if $q = m$

then

 print "Vzor sa v reťazci vyskytuje s posunom " $i-m$

102

porovnávanie pomocou automatu

- VÝPOČET PRECHODOVEJ FUNKCIE (P, Σ)


```

m := |P|
for q := 0 to m do
  for každý symbol a ∈ Σ do
    k := min(m+1, q+2)
    repeat k := k - 1
    until  $P_k \supset P_q a$ 
     $\delta(q, a) := k$ 
      
```

 return δ
- Zložitosť algoritmu:
 - Preprocesná fáza: $O(m|\Sigma|)$
 - Fáza vyhľadávania: $O(n)$

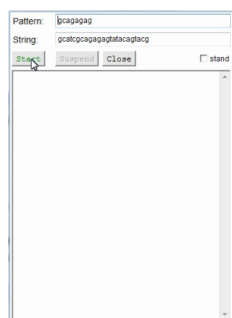
103

porovnávanie pomocou automatu

- Zložitosť algoritmu:
 - fáza predspracovania: $O(m|\Sigma|)$
 - fáza vyhľadávania: $O(n)$

104

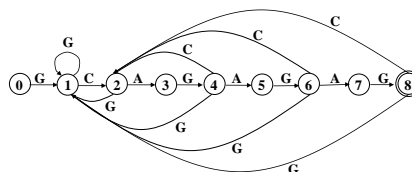
porovnávanie pomocou automatu



105

porovnávanie pomocou automatu

- DFA z príkladu:



106

Michael O. Rabin

- 1931 Breslau (Nemecko)
dnes Wrocław (Poľsko)
- 1953, M.Sc. Hebrew
University of Jerusalem
- 1956, Ph.D. Princeton
University
- práce:
 - výpočtová zložitosť,
 - matematická logika
 - Karp, Richard M.; Rabin,
Michael O. (March 1987).
Efficient randomized pattern-
matching algorithms



107

Richard M. Karp

- 3.1.1935, Boston –
- Bc 1955, Harvard
- Master 1956, Harvard
- PhD 1959 aplikovaná
matematika, Harvard
- International Computer
Science Institute in Berkeley
- práce:
 - Karp, Richard M.; Rabin,
Michael O. (March 1987).
Efficient randomized pattern-
matching algorithms



108

Rabinov-Karpov algoritmus

- Namiesto priameho porovnávania reťazca so vzorom sa porovnávajú výstupy hešovacích funkcií. Porovnáva sa heš vzoru s hešom vybraného podreťazca (vyberá sa podreťazec taký dlhý ako je dĺžka vzoru). Pokiaľ sa heše zhodujú, uskutočňuje sa porovnávanie jednotlivých znakov.

Rabinov-Karpov algoritmus – príklad 1/1

$t[0]$	$t[1]$	$t[2]$	$t[3]$	$t[4]$	$t[5]$	$t[6]$	$t[7]$	$t[8]$	$t[9]$	$t[10]$
A	C	D	E	A	C	A	C	C	D	E

← Hash(ACDE) = 10 →

$p[0]$	$p[1]$	$p[2]$	$p[3]$
A	C	D	B

← Hash(ACDB) = 5 →

109

110

Rabinov-Karpov algoritmus – príklad 1/2

$t[0]$	$t[1]$	$t[2]$	$t[3]$	$t[4]$	$t[5]$	$t[6]$	$t[7]$	$t[8]$	$t[9]$	$t[10]$
A	C	D	E	A	C	A	C	C	D	E

← Hash(CDEA) = 6 →

$p[0]$	$p[1]$	$p[2]$	$p[3]$
A	C	D	B

← Hash(ACDB) = 5 →

111

Rabinov-Karpov algoritmus

```

algorithm RabinKarp:
Input   pole znakov T, dĺžka n
         pole znakov P, dĺžka m

    hP := hash(P[1..m])
    hT := hash(T[1..m])
    for i from 1 to n-m+1
        if hT = hP
            if T[i..i+m-1] = P
                return i
            hT := hash(T[i+1..i+m])
    return nenašiel sa výskyt
  
```

112

Rabinov-Karpov algoritmus – príklad 2/1

$t[0]$	$t[1]$	$t[2]$	$t[3]$	$t[4]$	$t[5]$	$t[6]$	$t[7]$	$t[8]$	$t[9]$	$t[10]$
A	C	D	E	A	C	A	C	C	D	E

← Hash(ACDE) = 10 →

$p[0]$	$p[1]$	$p[2]$	$p[3]$
A	C	D	B

← Hash(ACDB) = 5 →

$j[0]$	$j[1]$	$j[2]$	$j[3]$
F	M	D	T

← Hash(FMDT) = 62 →

113

Rabinov-Karpov algoritmus – príklad 2/2

$t[0]$	$t[1]$	$t[2]$	$t[3]$	$t[4]$	$t[5]$	$t[6]$	$t[7]$	$t[8]$	$t[9]$	$t[10]$
A	C	D	E	A	C	A	C	C	D	E

← Hash(CDEA) = 6 →

$p[0]$	$p[1]$	$p[2]$	$p[3]$
A	C	D	B

← Hash(ACDB) = 5 →

$j[0]$	$j[1]$	$j[2]$	$j[3]$
F	M	D	T

← Hash(FMDT) = 62 →

114

Rabinov-Karpov algoritmus

- Hešovacia funkcia by mala mať tieto vlastnosti:
 - jednoducho vypočítateľná
 - s malou pravdepodobnosťou kolízií
 - Heš posunutého podreťazca by mal byť jednoducho odvodený z predchádzajúceho hešu (táto vlastnosť výrazne uľahčí výpočet a algoritmus sa tým stáva omnoho efektívnejší ako naivný)

115

Rabinov-Karpov algoritmus

- Hešovacia funkcia:

Predpokladáme, že nahradíme reťazec M znakov určitým celým číslom. Ak použijeme konštantu b - maximálny počet možných znakov v abecede, tak definujeme:

 - $x = t[i]b^M + t[i+1]b^{M-1} + \dots + t[i+M]$

Pokročíme v texte o jeden znak dopredu a hodnota x' bude:

 - $x' = t[i+1]b^M + t[i+2]b^{M-1} + \dots + t[i+M+1]$

ak preskúmame x a x' , tak zistíme, že:

 - $x' = (x - t[i]b^M)b + t[i+M+1]$

116

Rabinov-Karpov algoritmus

Hešovacia funkcia:

- Tretej požiadavke vyhovuje napríklad hešovacia funkcia definovaná ako polynóm $(m-1)$. stupňa, kde hodnoty znakov vystupujú ako koeficienty. Aby sme sa vyhli problémom s príliš veľkými číslami pri výpočtoch, používa sa modulo aritmetika:
 - $pathash = (f^{m-1} \text{ord}(pat_0) + f^{m-2} \text{ord}(pat_1) + \dots + f \text{ord}(pat_{m-2}) + \text{ord}(pat_{m-1})) \bmod p$
 - $texthash_i = (f^{m-1} \text{ord}(text_i) + f^{m-2} \text{ord}(text_{i+1}) + \dots + f \text{ord}(text_{i+m-2}) + \text{ord}(text_{i+m-1})) \bmod p$
 - $texthash_{i+1} = (f^{m-1} \text{ord}(text_{i+1}) + f^{m-2} \text{ord}(text_{i+2}) + \dots + f \text{ord}(text_{i+m-1}) + \text{ord}(text_{i+m})) \bmod p$
- $= (f (texthash_i - f^{m-1} \text{ord}(text_i)) + \text{ord}(text_{i+m})) \bmod p$
- Hešovaciú funkciu ovplyvňujú parametre f a p .

117

Rabinov-Karpov algoritmus

RABIN-KARP POROVNANIE(T, P, d, q)

```

 $n \leftarrow \text{dĺžka}(T)$ 
 $m \leftarrow \text{dĺžka}(P)$ 
 $h \leftarrow d^{m-1} \bmod q$ 
 $p \leftarrow 0$ 
 $t_0 \leftarrow 0$ 
for  $i \leftarrow 1$  to  $m$  //spracovanie
do  $p \leftarrow (dp + P[i]) \bmod q$ 
 $t_0 \leftarrow (dt_0 + T[i]) \bmod q$ 
for  $s \leftarrow 0$  to  $n - m$  //párovanie
do if  $p = t_s$ 
then if  $P[1..m] = T[s+1..s+m]$ 
then print "Vzor sa v reťazci vyskytuje s posunom" s
if  $s < n - m$ 
then  $t_{s+1} \leftarrow (d(t_s - T[s+1]h) + T[s+m+1]) \bmod q$ 

```

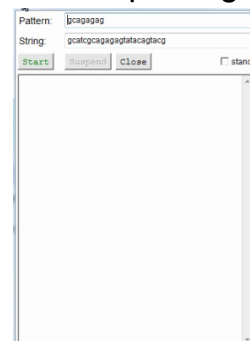
118

Rabinov-Karpov algoritmus

- Zložitosť algoritmu:
 - fáza predspracovania má časovú zložitosť $O(m)$
 - vyhľadavacia fáza má časovú zložitosť $O(mn)$
 - Očakávaná doba behu algoritmu je $O(n+m)$

119

Rabinov-Karpov algoritmus



120

Rabinov-Karpov algoritmus

```

C++:
#define REHASH(a,b,h)((h)-(a)*d)<1)&&(b))
void KR(char *x,int m, char *y,int n)
{
    int d,h=1,i,j;
    /* preprocessing */
    /* computes d = 2^M(m-1) with the left shift operator */
    for (d = 1; i = 1; m; ++i)
        d = (d<<i);
    for (h = h * i + 1; i < m; ++i)
    {
        h = ((h<<i) + x[i]);
        h = ((h<<i) + y[i]);
    }
    /* Searching */
    j = 0;
    while (j < n-m)
    {
        if (h == h*d && memcmp(x,y+j,m) == 0)
            OUTPUT(j);
        h = REHASH(h,d,y+j,m,h);
        ++j;
    }
}

```

121

časová efektívnosť

- Jeden vzor
 - BM $O(n/m)$
 - KMP $O(n)$
 - RK $O(mn)$
- Viac vzorov
 - BM, KMP $O(n.k)$
 - RK $O(n + k)$

122

aplikácie

- BM - textové editory – search/replace
- Karp-Rabin – vyhľadavanie plagiátov

123

prípony

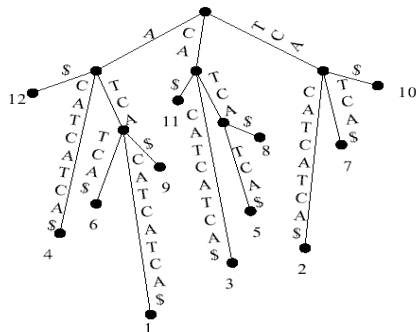
- Prípony reťazca $T = \text{"ATCACATCATCA"}$

ATCACATCATCA	$T_{(1)}$
TCACATCATCA	$T_{(2)}$
CACATCATCA	$T_{(3)}$
ACATCATCA	$T_{(4)}$
CATCATCA	$T_{(5)}$
ATCATCA	$T_{(6)}$
TCATCA	$T_{(7)}$
CATCA	$T_{(8)}$
ATCA	$T_{(9)}$
TCA	$T_{(10)}$
CA	$T_{(11)}$
A	$T_{(12)}$

124

príponový strom

- príponový strom reťazca $T = \text{"ATCACATCATCA"}$



vlastnosti príponového stromu

- každá hrana je ohodnotená nejakým podreťazcom reťazca T .
- každý vnútorný vrchol na najmenej dvoch potomkov.
- každá prípona $T_{(i)}$ má svoju ohodnotenú cestu z koreňa do listu pre $1 \leq i \leq n$.
- príponový strom má n listov.
- hrany vychádzajúce z toho istého vrchola majú ohodnotenia, ktoré sa určite nezačínajú tým istým znakom.

126

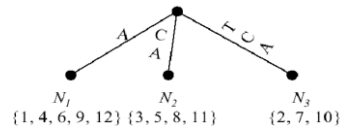
algoritmus vytvorenia príponového stromu

krok 1: rozdeľ všetky prípony do skupín podľa prvého znaku a vytvor vrchol.

krok 2: pre každú skupinu: ak obsahuje len jednu príponu, tak vytvor listový vrchol a hranu ohodnotenú touto príponou, inak nájdi najdlhšiu spoločnú predponu medzi príponami v tejto skupine a vytvor hranu vychádzajúcu z vrchola ohodnotenú najdlhšou spoločnou predponou. Vymaž túto predponu zo všetkých prípon v tejto skupine.

krok 3: opakuj predchádzajúci postup pre každý vrchol, ktorý nie je ukončený.

príklad vytvorenia príponového stromu

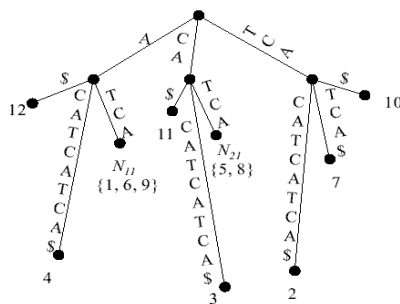


- $T = \text{"ATCACATCATCA"}$.
- začiatkové znaky: "A", "C", "T"
- v N_3 ,
 $T(2) = \text{"TCACATCATCA"}$
 $T(7) = \text{"TCATCA"}$
 $T(10) = \text{"TCA"}$
- najdlhšia spoločná predpona N_3 je "TCA"

127

128

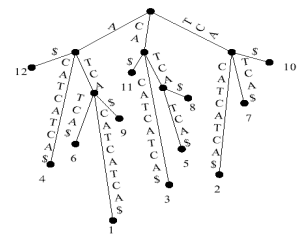
- $T = \text{"ATCACATCATCA"}$.
- druhá rekurzia:



129

nájdenie podreťazca pomocou príponového stromu

- $T = \text{"ATCACATCATCA"}$
- $P = \text{"TCAT"}$
 - P je na mieste 7 v T .
- $P = \text{"TCA"}$
 - P je na mieste 2, 7 a 10 v T .
- $P = \text{"TCATT"}$
 - P sa nenachádza v T .



130

príponový strom – časová zložitosť

- Príponový strom pre textový reťazec T dĺžky n sa dá zostrojiť v čase $O(n)$ time (pomocou zložitého algoritmu).
- Vyhľadanie vzoru P dĺžky m v príponovom strome vyžaduje $O(m)$ porovnaní.
- Hľadanie presného výskytu reťazca: $O(n+m)$ time

131

Príponové pole

- V príponovom poli sú všetky prípony reťazca T v neklesajúcom lexikálnom poradí.
- napr. $T = \text{"ATCACATCATCA"}$

i	1	2	3	4	5	6	7	8	9	10	11	12
A	12	4	9	1	6	11	3	8	5	10	2	7

4	ATCACATCATCA	$T_{(1)}$
11	TCACATCATCA	$T_{(2)}$
7	CACATCATCA	$T_{(3)}$
2	ACATCATCA	$T_{(4)}$
9	CATCATCA	$T_{(5)}$
5	ATCATCA	$T_{(6)}$
12	TCATCA	$T_{(7)}$
8	CATCA	$T_{(8)}$
3	ATCA	$T_{(9)}$
10	TCA	$T_{(10)}$
6	CA	$T_{(11)}$
1	A	$T_{(12)}$

1	A	$T_{(12)}$
2	ACATCATCA	$T_{(4)}$
3	ATCA	$T_{(9)}$
4	ATCACATCATCA	$T_{(1)}$
5	ATCATCA	$T_{(6)}$
6	CA	$T_{(11)}$
7	CACATCATCA	$T_{(3)}$
8	CATCA	$T_{(8)}$
9	CATCATCA	$T_{(5)}$
10	TCA	$T_{(10)}$
11	TCACATCATCA	$T_{(2)}$
12	TCATCA	$T_{(7)}$

132

Hľadanie v príponovom poli

- ak T sa reprezentuje príponovým poľom, vzor P sa dá nájsť v T v čase $O(m \cdot \log n)$ [binárnym vyhľadávaním](#).
- príponové pole sa dá určiť v čase $O(n)$ [lexikálnym hľadaním do hĺbky](#) v príponovom poli.
- celkový čas: $O(n + m \cdot \log n)$

133

hľadanie približného výskytu reťazca

- Textový reťazec T , $|T|=n$
vzor (reťazec) P , $|P|=m$
 k chýb, kde chybami môžu byť náhrada, vymazanie alebo vloženie znaku.
- napr:
 $T = \text{"pttapa"} , P = \text{"patt"} , k = 2$,
 $T_{1,2} , T_{1,3} , T_{1,4}$ a $T_{5,6}$ sú všetko reťazce vzdialené nie viac než 2 chyby od vzoru P .

-134

vzdialenosť editovania prípony

- Nech sú dané 2 reťazce S_1 a S_2 , vzdialenosť editovania prípony je najmenší počet náhrad, vložení a výmazov, ktoré prepíšu nejakú príponu S_1 do S_2 .
- napr:
 - $S_1 = \text{"ptt"} , S_2 = \text{"pt"} , \text{vzdialenosť editovania prípony medzi } S_1 \text{ a } S_2 \text{ je } 1$.
 - $S_1 = \text{"ptt"} , S_2 = \text{"p"} , \text{vzdialenosť editovania prípony medzi } S_1 \text{ a } S_2 \text{ je } 2$.
 - $S_1 = \text{"pt"} , S_2 = \text{"patt"} , \text{vzdialenosť editovania prípony medzi } S_1 \text{ a } S_2 \text{ je } 2$.

135

vzdialenosť editovania prípony

- Nech T a P sú reťazce, ak aspoň jedna zo vzdialeností editovania prípony medzi $T_{1,1} , T_{1,2} , \dots , T_{1,n}$ a P nie je väčšia než k , tak P sa približne vyskytuje v T s chybou nie väčšou než k .
- napr: $T = \text{"pttapa"} , P = \text{"patt"} , k = 2$
 - pre $T_{1,1} = \text{"p"} , P = \text{"patt"} , \text{vzdialenosť editovania prípony je } 3$.
 - pre $T_{1,2} = \text{"pt"} , P = \text{"patt"} , \text{vzdialenosť editovania prípony je } 2$.
 - pre $T_{1,5} = \text{"pttap"} , P = \text{"patt"} , \text{vzdialenosť editovania prípony je } 3$.
 - pre $T_{1,6} = \text{"pttapa"} , P = \text{"patt"} , \text{vzdialenosť editovania prípony je } 4$.

136

hľadanie približného výskytu reťazca

- riešenie metódami dynamického programovania
 - nech $E(i,j)$ označuje najmenšiu vzdialenosť editovania prípony medzi $T_{1,j}$ a $P_{1,i}$
 - nech $m = |P|$ označuje dĺžku vzoru,
 - potom $E(m,j)$ nám pomôže vyriešiť:
- úloha nájsť „najpribližnejší“ (t.j. s čo najmenším počtom chýb k) výskyt vzoru P v texte T :
 - vezmi podreťazec textu T , pre ktorý je $E(m,j)$ najmenšie.
- úloha nájsť približný výskyt vzoru P v texte T s najviac k chybami:
 - vezmi ľubovoľný podreťazec textu T , pre ktorý je $E(m,j) \leq k$.

137

najmenšia vzdialenosť editovania prípony

- nech $E(i,j)$ označuje najmenšiu vzdialenosť editovania prípony medzi $T_{1,j}$ a $P_{1,i}$
- pre všetky miesta i vo vzore P a všetky miesta j v texte T spočítame najmenšiu vzdialenosť editovania medzi prvými i znakmi vzoru P a ľubovoľným podreťazcom $T_{h,j}$ reťazca T , ktorý sa končí na mieste j (t.j. nejakou príponou $T_{1,j}$)
- to znamená prejsť všetky podreťazce T končiace v mieste j a určiť, ktorý z nich má najmenšiu vzdialenosť editovania prípony ku vzoru P . Jeho vzdialenosť je hodnota $E(i,j)$
- efektívne sa to dá počítať pomocou formuly (porovnaj s počítaním Levensteinovej vzdialenosti medzi reťazcami):

$$E(i,j) = E(i-1,j-1) \quad \text{if } P_i = T_j$$

$$E(i,j) = \min\{E(i,j-1), E(i-1,j), E(i-1,j-1)\} + 1 \quad \text{if } P_i \neq T_j$$

138

najmenšia vzdialenosť editovania prípony

$$E(i, j) = E(i-1, j-1) \quad \text{if } P_i = T_j$$

$$E(i, j) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1 \quad \text{if } P_i \neq T_j$$

- výsledok bude matica $n \times m$, $n = |T|$, $m = |P|$. Ako ju zrátať? Postupne, podľa formuly. Pomôže pridať 1 riadok s nulami.
- keďže vieme, že výsledok chceme použiť nielen na zistenie najmenšej vzdialenosti, ale aj na určenie približného podreťazca vyskytujúceho sa v texte T (ktorý je toľko vzdialený od vzoru P), je užitočné si zapamätať, ktorú z troch alternatív $E(i, j-1)$ alebo $E(i-1, j)$ alebo $E(i-1, j-1)$ sme v tom-ktorom kroku použili.

139

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0					
	2 a	2						
	3 t	3						
	4 t	4						

- $E(1, 1) = E(i-1, j-1)$ if $P_i = T_j$

140

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1				
	2 a	2						
	3 t	3						
	4 t	4						

- $E(1, 2) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

↑

141

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1				
	2 a	2						
	3 t	3						
	4 t	4						

- alebo

- $E(1, 2) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

↖

142

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1			
	2 a	2						
	3 t	3						
	4 t	4						

- $E(1, 3) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

143

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1		
	2 a	2						
	3 t	3						
	4 t	4						

- $E(1, 4) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

144

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	
	2 a	2						
	3 t	3						
	4 t	4						

- $E(1,5) = E(i-1, j-1)$ if $P_i = T_j$

145

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2						
	3 t	3						
	4 t	4						

- $E(1,6) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

146

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1					
	3 t	3						
	4 t	4						

- $E(2,1) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

147

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1				
	3 t	3						
	4 t	4						

- $E(2,2) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

148

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2			
	3 t	3						
	4 t	4						

- $E(2,3) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

149

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2			
	3 t	3						
	4 t	4						

- alebo
- $E(2,3) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

150

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2			
	3 t	3						
	4 t	4						

• alebo

- $E(2,3) = \min\{E(i,j-1), E(i-1,j), E(i-1,j-1)\} + 1$ if $P_i \neq T_j$

151

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1		
	3 t	3						
	4 t	4						

- $E(2,4) = E(i-1,j-1)$

if $P_i = T_j$

152

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	
	3 t	3						
	4 t	4						

- $E(2,5) = \min\{E(i,j-1), E(i-1,j), E(i-1,j-1)\} + 1$ if $P_i \neq T_j$

153

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3						
	4 t	4						

- $E(2,6) = E(i-1,j-1)$

if $P_i = T_j$

154

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2					
	4 t	4						

- $E(3,1) = \min\{E(i,j-1), E(i-1,j), E(i-1,j-1)\} + 1$ if $P_i \neq T_j$

155

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1				
	4 t	4						

- $E(3,2) = E(i-1,j-1)$

if $P_i = T_j$

156

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1			
	4 t	4						

- $E(3,3) = E(i-1, j-1)$ if $P_i = T_j$

157

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2		
	4 t	4						

- $E(3,4) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

158

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	
	4 t	4						

- $E(3,5) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

159

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4						

- $E(3,6) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

160

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3					

- $E(4,1) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

161

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		<i>T</i>						
		0	1	2	3	4	5	6
		p t t a p a						
<i>P</i>	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2				

- $E(4,2) = E(i-1, j-1)$ if $P_i = T_j$

162

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2	1			

- $E(4,3) = E(i-1, j-1)$ if $P_i = T_j$

163

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2	1	2		

- $E(4,4) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

164

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2	1	2	3	

- $E(4,5) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

165

počítanie matice E

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2	1	2	3	2

- $E(4,6) = \min\{E(i, j-1), E(i-1, j), E(i-1, j-1)\} + 1$ if $P_i \neq T_j$

166

hľadanie najbližšieho výskytu reťazca

- napr: $T = \text{"pttapa"}, P = \text{"patt"}$;
- pozrieme sa na posledný riadok, najmenšia vzdialenosť je v 3. stĺpci: $E(4,3)=1$. Z miesta (4,3) ideme naspäť podľa šípiek až do prvého riadku. Je to miesto (1,1). Stĺpec určuje začiatok podreťazca v T, ktorý je najbližším výskytom vzoru P – $T_{1,3} = \text{ptt}$. Jediná chyba sa „napraviť“ jedným vložením znaku a.

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2	1	2	3	2

167

hľadanie približného výskytu reťazca s chybou najviac 2

- napr: $T = \text{"pttapa"}, P = \text{"patt"}, k=2$;
- pozrieme sa na posledný riadok. Hodnoty $E(4,j) \leq k$ sú v stĺpcoch 2, 3, 4 a 6. Zodpovedajú približným výskytom $T_{1,2} = \text{pt}$, $T_{1,3} = \text{ptt}$, $T_{1,4} = \text{ptta}$, $T_{5,6} = \text{pa}$.

		T						
		0	1	2	3	4	5	6
		p t t a p a						
P	0	0	0	0	0	0	0	0
	1 p	1	0	1	1	1	0	1
	2 a	2	1	1	2	1	1	0
	3 t	3	2	1	1	2	2	1
	4 t	4	3	2	1	2	3	2

168

Hľadanie najdlhšej spoločnej podpostupnosti

- Longest common subsequence (LCS) problém
 - Dané sú dve postupnosti $x[1..m]$ a $y[1..n]$; máme nájsť najdlhšiu podpostupnosť, ktorá sa vyskytuje v oboch postupnostiach
 - Podpostupnosť: prvky v pôvodnej postupnosti nemusia byť nevyhnutne vedľa seba, ale ich poradie ostáva nezmenené
- $x = \{A \text{ B C B D A B}\}$, $y = \{B \text{ D C A B A}\}$
 - $\{B \text{ B A}\}$ je podpostupnosť oboch postupností x a y
- Algoritmus hrubej sily
 - Pre každú podpostupnosť v x , zisti či nie je podpostupnosťou y . Vráť najdlhšiu.
 - Koľko podpostupností je v x ?
 - 2^m
 - Aká by bola časová náročnosť?
 - 2^m podpostupností x porovnať s n prvkami postupnosti y
 - $O(n \cdot 2^m)$

169

Hľadanie najdlhšej spoločnej podpostupnosti

- Úloha: Porovnanie dvoch DNA reťazcov
- $X = \{A \text{ B C B D A B}\}$, $Y = \{B \text{ D C A B A}\}$
- Algoritmom hrubej sily porovnáme každú podpostupnosť X so znakmi v Y
 - $X = A \text{ B C B D A B}$
 - $Y = B \text{ D C A B A}$
- LCS problém má optimálnu subštruktúru: riešenie čiastkových problémov je časť konečného riešenia
- Čiastkový problém
 - Nájsť najdlhšiu spoločnú podpostupnosť párov prefixov X a Y
- Na vyriešenie tohto problému môžeme použiť dynamické programovanie!

170

Definícia dĺžky najdlhšej spoločnej podpostupnosti

- V prvom rade nájdeme dĺžku LCS. Neskôr zmodifikujeme algoritmus pre nájdenie LCS samotnej.
- Definujeme X_i a Y_j ako predpony X a Y dĺžky i respektíve j
- Definujeme $c[i, j]$ ako dĺžku LCS X_i a Y_j
- Potom dĺžka LCS X a Y bude $c[m, n]$
- Rekurzívna definícia $c[i, j]$

$$c[i, j] = \begin{cases} c[i-1, j-1] + 1 & \text{ak } x[i] = y[j], \\ \max(c[i, j-1], c[i-1, j]) & \text{inak} \end{cases}$$

171

počítanie dĺžky najdlhšej spoločnej podpostupnosti

- Začneme s $i=j=0$ (prázdna podpostupnosť X a Y)
 - Pretože X_0 a Y_0 sú prázdne reťazce, ich LCS je vždy prázdna ($c[0, 0] = 0$)
- LCS prázdneho reťazca a hocikákeho reťazca je pre každé i a j : $c[0, j] = c[i, 0] = 0$
- Pre výpočet $c[i, j]$ sa rozhodujeme medzi dvoma prípadmi:
 - $x[i] = y[j]$
 - Pri zhode symbolu v postupnostiach X a Y je dĺžka LCS X_i a Y_j rovnaká ako dĺžka LCS menšej postupnosti X_{i-1} a Y_{j-1} , plus 1
 - $x[i] \neq y[j]$
 - Ak sa symboly nezhodujú dĺžka ostáva nezmenená ($\max(c[i-1, j], c[i, j-1])$)

172

Algoritmus na nájdenie dĺžky najdlhšej spoločnej podpostupnosti

```

LCS DĹŽKA(X, Y)
m = length(X)
n = length(Y)
for i = 1 to m c[i, 0] = 0 //X0
for i = 1 to m c[i, 0] = 0 //Y0
for i = 1 to m //pre všetky Xi
    for j = 1 to n //pre všetky Yj
        if (Xi == Yj)
            c[i, j] = c[i-1, j-1] + 1
        else
            c[i, j] = max(c[i-1, j], c[i, j-1])
return c

```

- Príklad: $X = \text{ABCB}$; $Y = \text{BDCAB}$
 - $\text{LCS}(X, Y) = \text{BCB}$
 - $X = A \text{ B C B}$
 - $Y = B \text{ D C A B}$

173

Najdlhšia spoločná podpostupnosť – príklad (inicializácia)

		j	0	1	2	3	4	5
			Y _j	B	D	C	A	B
i		0	X _i					
1	A							
2	B							
3	C							
4	B							

ABCB
BDCAB

$X = \text{ABCB}$; $m = |X| = 4$
 $Y = \text{BDCAB}$; $n = |Y| = 5$
 alokácia 2-rozmerného poľa $c[0..4, 0..5]$

174

Najdlhšia spoločná podpostupnosť – príklad (1)

	j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B	
0	X _i	0	0	0	0	0	0
1	A	0					
2	B	0					
3	C	0					
4	B	0					

ABCB
BDCAB

for $i = 1$ to m $c[i,0] = 0$
for $j = 1$ to n $c[0,j] = 0$

175

Najdlhšia spoločná podpostupnosť – príklad (2)

	j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B	
0	X _i	0	0	0	0	0	0
1	A	0	0				
2	B	0					
3	C	0					
4	B	0					

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

176

Najdlhšia spoločná podpostupnosť – príklad (3)

	j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B	
0	X _i	0	0	0	0	0	0
1	A	0	0	0	0		
2	B	0					
3	C	0					
4	B	0					

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

177

Najdlhšia spoločná podpostupnosť – príklad (4)

	j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B	
0	X _i	0	0	0	0	0	0
1	A	0	0	0	0	1	
2	B	0					
3	C	0					
4	B	0					

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

178

Najdlhšia spoločná podpostupnosť – príklad (5)

	j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B	
0	X _i	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	B	0					
3	C	0					
4	B	0					

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

179

Najdlhšia spoločná podpostupnosť – príklad (6)

	j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B	
0	X _i	0	0	0	0	0	0
1	A	0	0	0	0	1	1
2	B	0	1				
3	C	0					
4	B	0					

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

180

Najdlhšia spoločná podpostupnosť – príklad (7)

j	0	1	2	3	4	5
i	Y_i	B	D	C	A	B
0	X_i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	
3	C	0				
4	B	0				

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

181

Najdlhšia spoločná podpostupnosť – príklad (8)

j	0	1	2	3	4	5
i	Y_i	B	D	C	A	B
0	X_i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0				
4	B	0				

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

182

Najdlhšia spoločná podpostupnosť – príklad (9)

j	0	1	2	3	4	5
i	Y_i	B	D	C	A	B
0	X_i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1		
4	B	0				

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

183

Najdlhšia spoločná podpostupnosť – príklad (10)

j	0	1	2	3	4	5
i	Y_i	B	D	C	A	B
0	X_i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	
4	B	0				

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

184

Najdlhšia spoločná podpostupnosť – príklad (11)

j	0	1	2	3	4	5
i	Y_i	B	D	C	A	B
0	X_i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	2
4	B	0				

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

185

Najdlhšia spoločná podpostupnosť – príklad (12)

j	0	1	2	3	4	5
i	Y_i	B	D	C	A	B
0	X_i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	2
4	B	0	1			

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

186

Najdlhšia spoločná podpostupnosť – príklad (13)

j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B
0	X _i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	2
4	B	0	1	1	2	2

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

187

Najdlhšia spoločná podpostupnosť – príklad (13)

j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B
0	X _i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	2
4	B	0	1	1	2	2

ABCB
BDCAB

if ($X_i == Y_j$)
 $c[i,j] = c[i-1,j-1] + 1$
 else $c[i,j] = \max(c[i-1,j], c[i,j-1])$

$c[4,5]$ obsahuje dĺžku najdlhšej spoločnej podpostupnosti

188

Analýza algoritmu pre nájdenie najdlhšej spoločnej podpostupnosti

- LCS algoritmus vypočíta hodnoty každého vstupu podľa $c[m,n]$. Aký je teda výpočtový čas?
- $O(m \cdot n)$
 - Každá hodnota $c[i,j]$ je spočítaná v konštantnom čase, a v poli máme $m \cdot n$ prvkov
- Zatiaľ sme našli len dĺžku najdlhšej spoločnej podpostupnosti.
- Ďalej je potrebné nájsť najdlhšiu spoločnú podpostupnosť.
- Musíme modifikovať algoritmus aby nám dával výstup najdlhšej spoločnej podpostupnosti postupnosti X a Y
 - V poli $c[i,j]$ máme všetko zaznamenané
 - Každá hodnota $c[i,j]$ závisí na $c[i-1,j]$ alebo $c[i,j-1]$
 - Pre každú hodnotu $c[i,j]$ vieme určiť ako sme ju dosiahli

189

Hľadanie najdlhšej spoločnej podpostupnosti

$$c[i,j] = \begin{cases} c[i-1,j-1] + 1 & \text{ak } x[i] = y[j], \\ \max(c[i,j-1], c[i-1,j]) & \text{inak} \end{cases}$$

2	2
2	3

V tomto prípade

$$c[i,j] = c[i-1,j-1] + 1 = 2 + 1 = 3$$

1	1
2	2

V tomto prípade

$$c[i,j] = \max(c[i,j-1], c[i-1,j]) = \max(2, 1) = 2$$

- Môžeme začať z $c[m,n]$ a ísť späť
- Ak $c[i,j] = c[i-1,j-1] + 1$, zapamätáme si $x[i]$
 - $x[i]$ je časť z najdlhšej spoločnej podpostupnosti
- Ak $i = 0$ alebo $j = 0$ (dosiahneme začiatok), výstupom sú písmená odpamätané v X, usporiadané v opačnom poradí

190

Hľadanie najdlhšej spoločnej podpostupnosti

j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B
0	X _i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	2
4	B	0	1	1	2	2

191

Hľadanie najdlhšej spoločnej podpostupnosti

j	0	1	2	3	4	5
i	Y _i	B	D	C	A	B
0	X _i	0	0	0	0	0
1	A	0	0	0	0	1
2	B	0	1	1	1	2
3	C	0	1	1	2	2
4	B	0	1	1	2	2

Najdlhšia spoločná postupnosť (od zadu): **B C B**
 Najdlhšia spoločná postupnosť: **B C B**

192