

Algoritmusok és adatszerkezetek II

Mintaillesztés témakör jegyzete

Készült Ásványi Tibor előadásai és gyakorlatai alapján

Sárközi Gergő, 2021-22-1. félév

Nincsen lektorálva!

Tartalomjegyzék

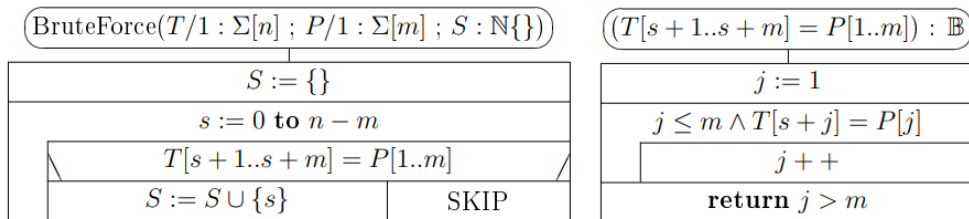
1. Mintaillesztés	2
2. Egyszerű (brute force) algoritmus	2
3. Quicksearch	3
3.1. Quicksearch példa	4
4. Knuth-Morris-Pratt (lineáris) algoritmus RÖVIDEN	5
4.1. Példa	6
5. Knuth-Morris-Pratt (lineáris) algoritmus	7
5.1. Jelölések	7

1. Mintaillesztés

- Abécé: $\Sigma = \{\sigma_1, \sigma_2, \dots, \sigma_d\}$ ($1 \leq d < \infty$ konstans)
- Szöveg, amiben keresünk: $T/1 : \Sigma[n]$ ($1 \leq n$)
- Minta, amit keresünk: $P/1 : \Sigma[m]$ ($1 \leq m \leq n$)
- $s \in 0..(n - m)$ P érvényes eltolása T -n $\Leftrightarrow T[s + 1..s + m] = P[1..m]$
- A cél az érvényes eltolások S halmazának megállapítása

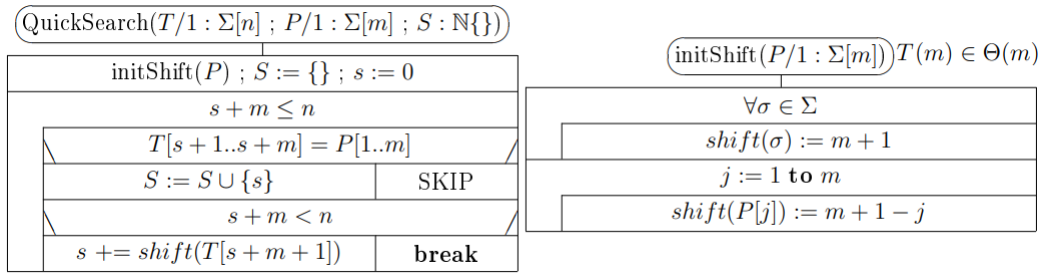
2. Egyszerű (brute force) algoritmus

- Minden lehetséges s értékre, egymástól függetlenül, próbáljuk a mintát
- Időkomplexitás: $MT(n, m) \in \Theta(n * m)$ és $mT(n, m) \in \Theta(n)$
 - Alapból $MT \in \Theta((n - m + 1) * m)$ és $mT \in \Theta(n - m + 1)$
 - $m \leq n \implies (n - m + 1) \in \Theta(n)$
 - Tehát $MT \in \Theta(n * m)$ és $mT \in \Theta(n)$ (mint legfelül)
 - Ha m nem elhanyagolható n -hez képest ($m \geq \epsilon * n$ ahol $0 < \epsilon < 1$) akkor $(n * m) \in \Theta(n^2) \implies MT \in \Theta(n^2)$



3. Quicksearch

- Egynél nagyobb lépésekben növeli az s eltolását (de nem ugrik át egy érvényes eltolást sem)
- Előkészítő fázis: Ábécé minden σ eleméhez $shift(\sigma) \in 1..m+1$ címke
 - Csak a mintától függ, a szövegtől nem
- $shift(\sigma)$ működése:
 - σ mindig a minta utáni első karakter a szövegben: $\sigma = T[s+m+1]$
 - Megmondja $T[s+1..s+m]$ megnézése után mennyivel nőjön s
 - Ha $\sigma \in P$: s mennyivel nőjön, hogy a minta illeszkedhessen a $T[s+m+1]$ karakterre (pl. ha $P[m] = \sigma$ akkor $shift(\sigma) = 1$)
 - Ha $\sigma \notin P$: minta átugorja $T[s+m+1]$ karaktert ($shift(\sigma) = m+1$)
- Időkomplexitás:
 - $mT \in \Theta(\frac{n}{m+1} + m)$ (pl. T és P diszjunktak)
 - * Jobb, mint a brute force megoldás
 - $MT \in \Theta((n-m+2) * m)$ (pl. T és P mind azonos σ sokszor)
 - * Azonos brute force-szal, de gyakorlatban lassabb
 - Átlagosan gyorsabb, mint a brute force, de azért nem optimális



3.1. Quicksearch példa

- Bal fenti ábra:
 - $xxxx$ jelöli a mintával az eltolás előtt összehasonlított szövegrészt
 - Az eltolás mértékét mutatja be: az eltolás utána állapot látható

Szöveg: ...xxxxA.....xxxxB.....xxxxC.....xxxxD...
 Minta: CADA CADA CADA CADA

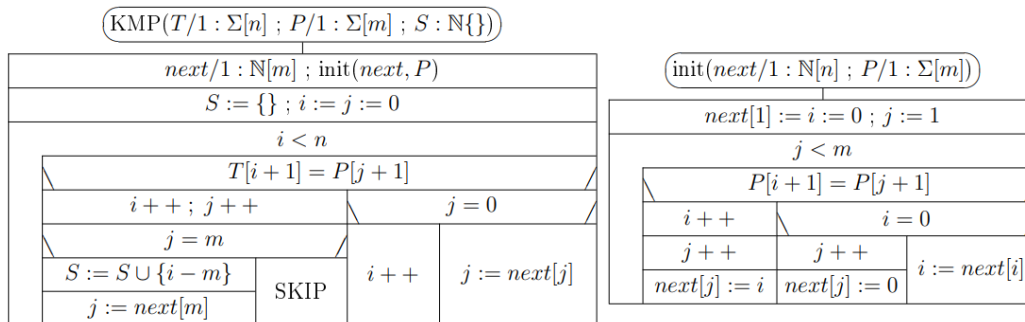
σ	A	B	C	D
$shift(\sigma)$	1	5	4	2

σ	A	B	C	D	
initial $shift(\sigma)$	5	5	5	5	5
C			4		4
A	3				3
D				2	2
A	1				1
final $shift(\sigma)$	1	5	4	2	

$i =$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
$T[i] =$	A	D	A	B	A	B	C	A	D	A	B	C	A	B	A	D	A	C	A	D	A	D	A
	\emptyset	A	D	A																			
		\emptyset	A	D	A																		
$s = 6$							<u>C</u>	<u>A</u>	<u>D</u>	<u>A</u>													
												<u>C</u>	<u>A</u>	B	A								
														\emptyset	A	D	A						
$s = 17$																		<u>C</u>	<u>A</u>	<u>D</u>	<u>A</u>		
																				\emptyset	A	D	A

4. Knuth-Morris-Pratt (lineáris) algoritmus RÖVIDEN

- Lineáris időben végzi el a feladatot
- Nem kell minden esetben a minta elejétől kezdeni az illesztést:
a prefixet nem kell újra vizsgálni, ha az egyezik a szuffixszel
- Előfeldolgozás: $(\Theta(m))$ idő alatt végbemegy
 - megadunk egy *next* függvényt, ami megadja a leghosszabb megegyező prefix-suffix párok hosszát minden minta kezdőszeletre (hossza)
 - *next*(*j*) a leghosszabb olyan *P* prefix hossza, amely *P* első *j* karakterének szuffixe (de nem egyezik meg vele), azaz $next(j) \in 0..(j-1)$
- A szövegben nem kell visszaugrani, azaz buffer nélkül is használható.
(Minden karaktert csak egyszer olvasunk ki, és csak "előrefelé" haladunk.)
- A mintát sikeres/sikertelen illesztés esetén annyival toljuk előrebb, amekkora a sikeresen illesztet részminta hossza MÍNUSZ a sikeresen illesztet részminta legnagyobb szuffixe, ami egyben prefix.
Azaz ez a legnagyobb szuffix lesz a minta kezdete.
- Időkomplexitás: $MT = mT \in \Theta(n)$
 - $\Omega(n)$, mert *i* egyesével nő és *n*-ig megy
 - $O(n)$, $2i - j$ értéke mindig szig. mon. nő, tehát max $2n$ iteráció



- $next[1] = 0$
- $next[i+1] \leq next[i] + 1$
- $next(j) \in 0..(j-1)$ ($j \in 1..m$)

- *init* ciklusának invariánsa:
 - $i \leq j \leq m$
 - P első i karaktere szuffixe P első j karakterének
 - és $\forall l \in (i+2)..j : P$ első l karaktere nem szuffixe P első $j+1$ karakterének, de egyenlőek lehetnek
 - és $next[1..j] = next(1..j)$ (azaz a tömb a fv alapján van töltve)
- *KMP* ciklusának invariánsa:
 - $i \in 0..n$ és $j \in 0..(m-1)$ és $j \leq i$
 - és $S = \{s \in 0..(i-m) \mid T[(s+1)..(s+m)] = P\}$
 - és P első j karaktere szuffixe T első i karakterének (vagy egyenlők)
 - és $\forall l \in (j+2)..m : P$ első l karaktere nem szuffixe T első $i+1$ karakterének (és nem is egyenlők)

4.1. Példa

i	j	$next[j]$	1 \underline{A}	2 \underline{B}	3 \underline{A}	4 \underline{B}	5 \underline{B}	6 \underline{A}	7 \underline{B}	8 \underline{A}
0	1	0		A						
0	2	0			<u>A</u>					
1	3	1			<u>A</u>	<u>B</u>				
2	4	2			<u>A</u>	<u>B</u>	A			
0	4	2					A			
0	5	0						<u>A</u>		
1	6	1						<u>A</u>	<u>B</u>	
2	7	2						<u>A</u>	<u>B</u>	<u>A</u>
3	8	3								

Minta:

P = ABABBABA

A végeredmény:

$P[j] =$	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>
$j =$	1	2	3	4	5	6	7	8
$next[j] =$	0	0	1	2	0	1	2	3

$i =$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
$T[i] =$	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>
	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	B												
$s=2$			<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>							
$s=7$								<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>		
													<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	B
															<u>A</u>	<u>B</u>	<u>A</u>

$$S = \{2; 7\}$$

5. Knuth-Morris-Pratt (lineáris) algoritmus

5.1. Jelölések

- Akár teljes prefix: $x \sqsubseteq y \Leftrightarrow \exists z : x + z = y$
- Igazi prefix: $x \sqsubset y \Leftrightarrow x \sqsubseteq y \wedge x \neq y$
- Akár teljes szuffix: $x \sqsupseteq y \Leftrightarrow \exists z : z + x = y$
- Igazi szuffix: $x \sqsupset y \Leftrightarrow x \sqsupseteq y \wedge x \neq y$
- Az üres sztring mindennek a prefixe és a szuffixe is.
- Kezdőszelet: $A_j = A[1..j]$ (ezt a jelölést ritkán használjuk)
 - A_0 az üres sztring
- Prefix-szuffix: $x \sqbox y \Leftrightarrow x \sqsubset y \wedge x \sqsupset y$
- i . legnagyobb elem: $\max_i H \quad (i \in 1..|H|)$
 - $\max_1 H = \max H$ és $\max_{|H|} H = \min H$
- $H(j) = \{h \in 0..j-1 \mid P_h \sqsupset P_j\} = \{|x| \mid x \sqbox P_j\} \quad (j \in 1..m)$

Azaz azon sztring hosszak, amelyek prefixek és szuffixek is P -nek egyben.
- $next(j) = \max H(j) \quad (j \in 1..m)$

Leghosszab P -beli prefix hossza, ami egyben valódi szuffixe P_j -nek.

NINCS BEFEJEZVE, NAGYON HIÁNYOS