

Haladó informatikai algoritmusok

Mintaillesztés

Knuth-Morris-Pratt algoritmus

Boyer-Moore algoritmus

Knuth-Morris-Pratt (KMP) algoritmus

Lineáris idejű mintaillesztő algoritmus.

Kiküszöböli a δ átmeneti függvény kiszámítását és az illesztési ideje $\Theta(n)$.

Egy $\pi[1..m]$ **segédfüggvényt** használunk, amelyet előzetesen számítunk ki a minta alapján $\Theta(m)$ idő alatt.

A π tömb segítségével a δ átmeneti függvény értékei "menet közben" hatékonyan számolhatók.

Ez azt jelenti, hogy bármely $q = 0, 1, \dots, m$ állapotra és $a \in \Sigma$ jelre, $\pi(q)$ érték adja meg $\delta(q, a)$ számítási módját és ez független a -tól.

A π tömbnek csak m eleme van, míg δ -nak $\sigma(m|\Sigma|)$, így az előfeldolgozás során megspóroljuk a $|\Sigma|$ tényezőt, ha π -t határozzuk meg δ helyett.

Alapötlet

A léptetés mértéke előre kiszámítható a mintára, továbbá a léptetés után az összehasonlítást nem kell a minta elejéről kezdeni.

Legyen a minta S és $m = |S|$.

A minta prefix függvénye: Egy minta prefix függvénye tartalmazza azokat az ismereteket, amelyek megadják, hogyan illeszkedik a minta önmaga eltoltjaira:

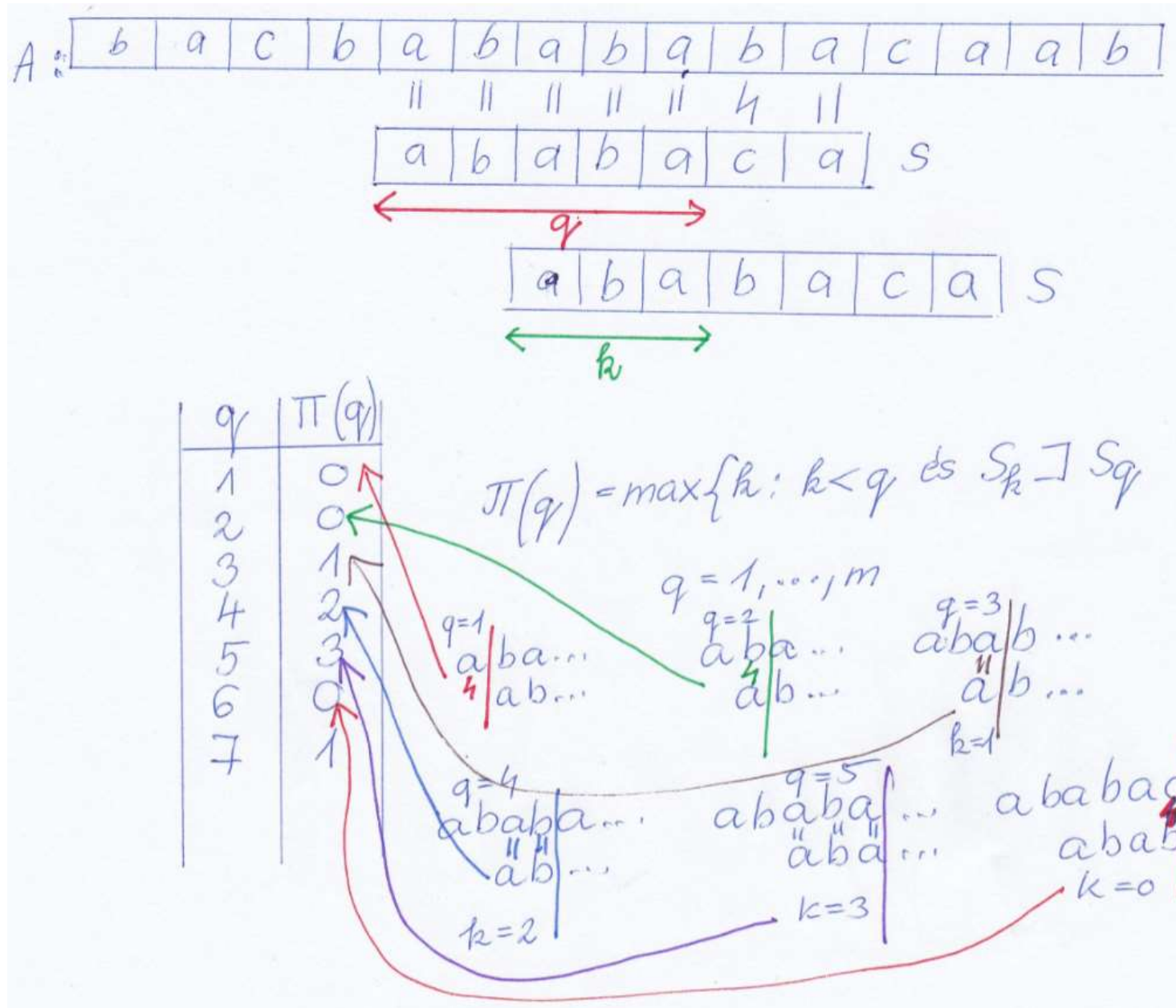
$$\pi : \{1, \dots, m\} \rightarrow \{0, \dots, m - 1\}$$

$$\pi(q) = \max\{k : k < q \text{ és } S_k \text{ szuffixe } S_q\} \quad q = 1, \dots, m$$

(a leghosszabb olyan S -beli prefix hossza, amely valódi szuffixe S_q -nak)

Ha $S[1..q] = A[i - q .. i - 1]$, de $S[q + 1] \neq A[i]$, akkor legalább annyit kell léptetni a mintát, hogy S_k szuffixe S_q teljesüljön.

Példa



Eljárás

KMP-ILLESZTO(T,P)

1. $n := \text{hossz}[A]$
2. $m := \text{hossz}[S]$
3. $\pi := \text{PREFIX} - \text{FUGGVENY} - \text{SZAMITAS}(S)$
4. $q := 0$
5. **for** $i := 1$ **to** n
6. **do while** $q > 0$ **and** $S[q + 1] \neq A[i]$
7. **do** $q := \pi[q]$
8. **if** $S[q + 1] = A[i]$
9. **then** $q := q + 1$
10. **if** $q = m$
11. **then** print "A minta illeszkedik az ("i-m+1")-edik pozícióra"
12. $q := \pi[q]$

Eljárás

PREFIX-FUGGVENY-SZAMITAS(S)

1. $m := \text{hossz}[S]$
2. $\pi[1] := 0$
3. $k := 0$
4. **for** $q := 2$ **to** m
5. **do while** $k > 0$ **and** $S[k + 1] \neq S[q]$
6. **do** $k := \pi[k]$
7. **if** $S[k + 1] = S[q]$
8. **then** $k := k + 1$
9. $\pi[q] := k$
10. **return** π

Algoritmus végigkövetése

$A: \quad 1. \ 2. \ 3. \ 4. \ 5. \ 6. \ 7. \ 8. \ 9. \ 10.$
 $ab \ c \ b \ ab \ ab \ a \ b$

$S: \quad ab \ ab$
 $1. \ 2. \ 3. \ 4.$

$n=10$
 $m=4$

Π : prefix-fő-szalmitás (s)
 $m=4 \quad \Pi[1] = \emptyset \quad k = \emptyset$

$q=2$
 $1. \ k > 0 \ X$
 $2. \ S[1] \stackrel{?}{=} S[2]$
 $a \stackrel{?}{=} b \ X$
 $3. \ \Pi[2] := \emptyset$
 $\Pi[1..4] = (0, 0, 1, 2)$

$q=3$
 $1. \ k > 0 \ X$
 $2. \ S[1] \stackrel{?}{=} S[3]$
 $a \stackrel{?}{=} a \ \checkmark$
 $3. \ \Pi[3] = 1$

$q=4$
 $1. \ k=1 > 0 \ \checkmark$
 $S[2] \stackrel{?}{=} S[4]$
 $b = b \ \checkmark$
 $2. \ \Pi[4] = 2$

$n=10 \quad m=4$
 $q=0 \quad i=1 \quad A: a$
 $q > 0 \ X$
 $S[1] \stackrel{?}{=} A[1]$
 $a = a \ \checkmark$
 $q = m$
 $1 = 4 \ X$
nem érünk a minta végére

$q=0$ marad

$i=2 \quad q=1 \quad A: ab$
 $S: ab$
 $q > 0 \ \checkmark \quad S[2] \neq A[2]$
 $b \neq b \ X$
 $S[2] \stackrel{?}{=} A[2]$
 $b = b \ \checkmark$
 $q = m$
 $2 = 4 \ X$
 $q=2$

$i=3 \quad q=2$
 $q > 0 \ \checkmark \quad S[3] \neq A[3]$
 $a \neq c \ X$
 $S[3] \stackrel{?}{=} A[3]$
 $a = c \ X$
 $q = m$
 $3 = 4 \ X$
 $q=1$

$i=4 \quad A: abcb$
 $S: a$
 $q > 0 \ X$
 $0 > 0 \ X$
 $S[4] \stackrel{?}{=} A[4]$
 $a = b \ X$
 $0 = 4 \ X$
 $q=0$

$i=5 \quad A: abcba$
 $S: a$
 $0 > 0 \ X$
 $S[1] \stackrel{?}{=} A[5]$
 $a = a \ \checkmark$
 $1 = 4 \ X$
 $q=1$

$i=6 \quad A: abcba$
 $S: ab$
 $1 > 0 \ \checkmark \quad S[2] \neq A[6]$
 $b \neq b \ X$
 $S[2] = A[6]$
 $a = a \ \checkmark$
 $2 = 4 \ X$
 $q=2$

$i=7 \quad A: abcba$
 $S: aba$
 $2 > 0 \ \checkmark \quad S[3] \neq A[7]$
 $a \neq a \ X$
 $S[3] \stackrel{?}{=} A[7]$
 $a = a \ \checkmark$
 $3 = 4 \ X$
 $q=3$

$i=8 \quad A: abcba$
 $S: abab$
 $3 > 0 \ \checkmark \quad S[4] \neq A[8]$
 $b \neq b \ X$
 $S[4] \stackrel{?}{=} A[8]$
 $b = b \ \checkmark$
 $4 = 4 \ \checkmark$
 $q=4$
A minta illeszkedik az 5. pozícióra.
 $q=2$ (A következő illeszkedik következőre)

$i=9$
 $2 > 0$ és $S[3] \neq A[9]$
 $a \neq a \ X$
 $S[3] \stackrel{?}{=} A[9]$
 $a = a$
 $q=3$

$3 \neq 4$

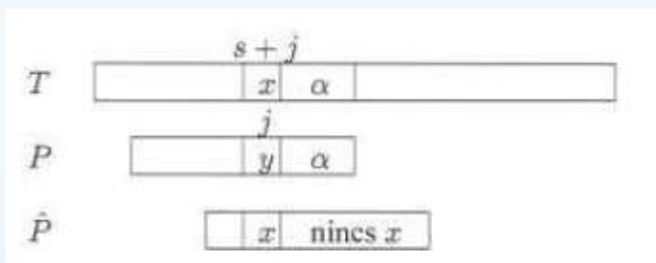
$i=10$
 $10 > 0$ és $S[4] \neq A[10]$
 $b \neq b \ X$
 $S[4] \stackrel{?}{=} A[10]$
 $b = b \ \checkmark$
 $q=4$
 $4 = 4$
A minta illeszkedik a 10-4+1 = 7. pozícióra
 $q = \Pi[4] = 2$
 $q=2$

Boyer-Moore algoritmus

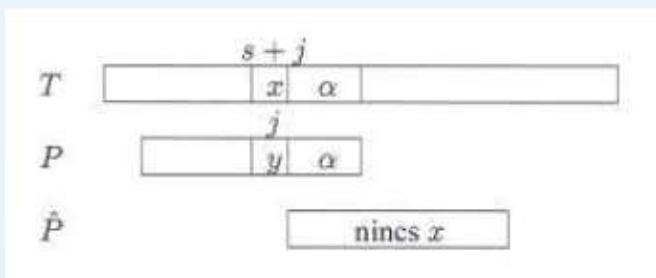
A mintát balról jobbra mozgatja, de **jobbról balra ellenőrzi az illeszkedést**.

Heurisztikák:

1. Az *"utolsó karakter"* heurisztika minden egyes betűhöz (karakterhez) megadja, hogy az hol fordult elő utoljára a mintában. Ha a szöveg épp vizsgált karaktere nem fordul elő a mintában, akkor a mintát eltolhatjuk eme karakter után. Ha viszont szerepel benne a karakter, akkor annak az utolsó előfordulását illesztjük a szöveg megfelelő karakteréhez. Utolsó pozíció heurisztika, amikor a nem egyező karakter szerepel a mintában:



Utolsó pozíció heurisztika, amikor a nem egyező karakter nem szerepel a mintában:



Példák

①

T: abcacbc**caab**

P: caab

a.) caab

a.) caab

a.) caab

a.) **caab**
megtaláltuk

1/a.)

T		^{stj}	X	α	
P			y		
P			X	nincs X	

ha a szövegbeni karakter szerepel a mintában, akkor annak az utolsó előfordulását illesztjük a szöveg megfelelő karakteréhez

②

T: acbacbc**caab**

P: caa

b.) caa

b.) **caa**

③

T: abcabc**babab**

b.) baba

a.) baba

b.) **baba**

1/b.)

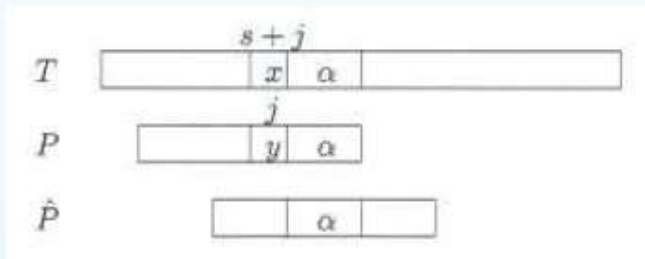
T		^{stj}	X	α	
P			y	α	
P			nincs X		

ha a szöveg épp visszáért karaktere nem fordul elő a mintában, akkor a mintát eltoljuk eme karakter után

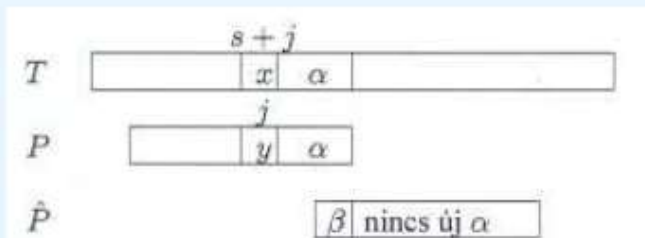
Boyer-Moore algoritmus

2. A "jó suffix" heurisztikánál úgy mozgatjuk a mintát, hogy a szövegnek arra a részére, amelyre már illeszkedett a minta vizsgált része, újra (legalább részben) illeszkedjék. Ha az illeszkedő részminta újra előfordul a szövegben, akkor ezt a részt kell illeszteni, egyébként eme minta egy suffixét kell illeszteni a szöveghez.

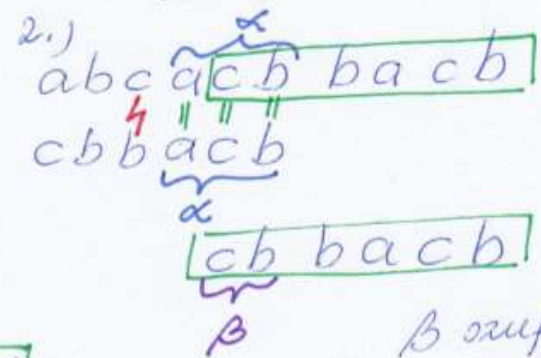
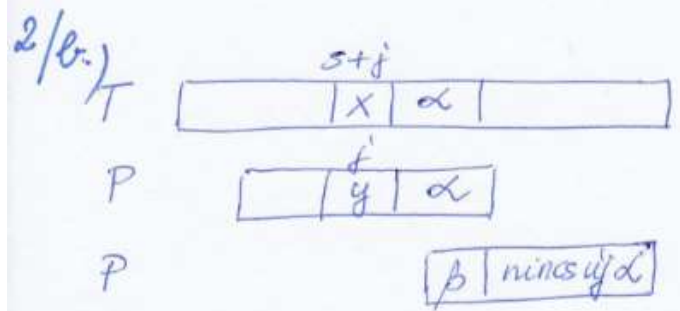
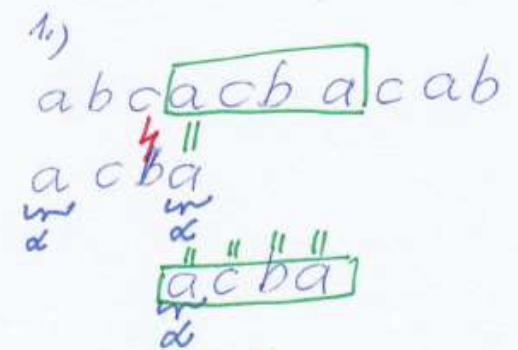
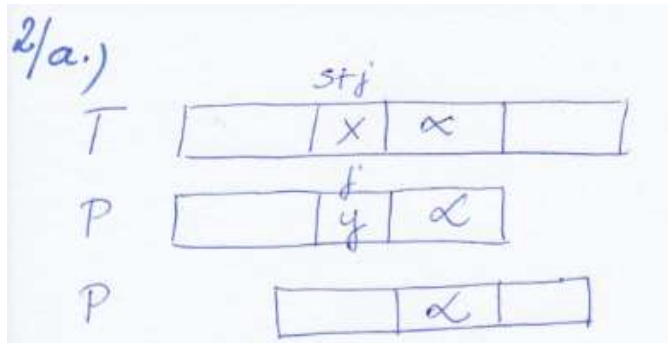
A jó suffix heurisztika, amikor az α újra előfordul:



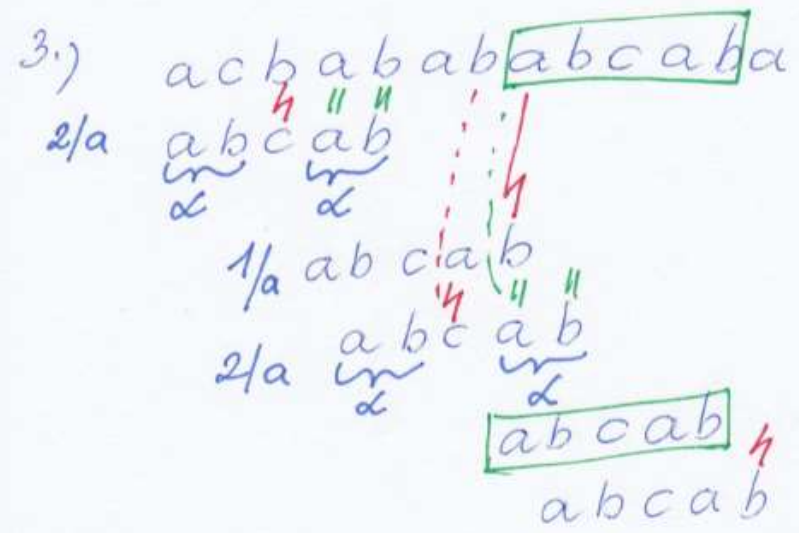
A jó suffix heurisztika, amikor az α nem fordul elő újra:



Példák



β szuffixe a mintának



1 egyezés van

Függvény1

Function UTOLSÓ-POZÍCIÓ(P, m, S)

Input: P minta, m minta hossza, S ábécé

Output: u utolsó pozíció heurisztika értékei

1. **forall the** $a \in S$ **do** $u[a] \leftarrow 0$
2. **for** $j \leftarrow 1$ **to** m **do** $u[P[j]] \leftarrow j$
3. **return** u

Függvény2

Function JÓ-SZUFFIX(P, m)

Input: P minta, m minta hossza

Output: g jó szuffix heurisztika értékei

1. $p_1 \leftarrow \text{PREFIX-FÜGGVÉNY-SZÁMÍTÁS}(P)$
2. $P' \leftarrow \text{FORDÍT}(P)$
3. $p_2 \leftarrow \text{PREFIX-FÜGGVÉNY-SZÁMÍTÁS}(P')$
4. **for** $j \leftarrow 0$ **to** m **do** $g[j] \leftarrow m - p_1[m]$
5. **for** $l \leftarrow 1$ **to** m **do**
6. $j \leftarrow m - p_2[l]$
7. **if** $g[j] > l - p_2[l]$ **then** $g[j] \leftarrow l - p_2[l]$
8. **endfor**
9. **return** g

Függvény3

BOYER-MOORE-ILLESZTŐ(T, P, S)

Input: T szöveg, P minta, S ábécé

Output: illeszkedés esetén üzenet

1. $n \leftarrow T.hossz$
2. $m \leftarrow P.hossz$
3. $u \leftarrow \text{UTOLSÓ-POZÍCIÓ}(P, m, S)$
4. $g \leftarrow \text{JÓ-SZUFFIX}(P, m)$
5. $s \leftarrow 0$
6. **while** $s \leq n - m$ **do**
7. $j \leftarrow m$
8. **while** $j > 0$ **and** $P[j] = T[s + j]$ **do** $j \leftarrow j - 1$
9. **if** $j = 0$ **then** *kiír* illeszkedés az $s + 1$ pozíción
10. **else** $s \leftarrow s + \max(g[j], j - u(T[s + j]))$
11. **endif**
12. **endwhile**

T = a c b a c b c a c a

S = {a, b, c}

P = c a c

n = 10

m = 3

u[a] = 0 u[b] = 0 u[c] = 0

j=1 u[P[1]] = u[c] = 1 j=2 u[P[2]] = u[a] = 2

j=3 u[P[3]] = u[c] = 3

u[a] = 2 u[b] = 0 u[c] = 3

UTOLSO-POZICIO(P, m, S)

JO-SZUFFIX(P, m)

P ₁	cac
	c
	c
	c
	c
	c

q	k
1	0
2	0
3	1

P₁ = [c, 0, 1]

α κ
m₁ m₂
cac
w cac

↓
2-vel tudom
are'le' osszetartni

P' = FORDIT(P) = cac

P₂ cac

q	k
1	0
2	0
3	1

P₂ = [0, 0, 1]

j=0 g[0] = 3 - P₁[3] = 3 - 1 = 2

j=1 g[1] = 3 - P₁[3] = 2

j=2 g[2] = 2

j=3 g[3] = 2

g = [2, 2, 2, 1] cac
c

l=1 j = m - P₂[1] = 3 - 0 = 3

g[3] > 1 - P₂[1] ✓ g[3] = 1 - P₂[1] = 1

cac
cac

l=2 j = m - P₂[2] = 3 - 0 = 3

g[3] > 2 - P₂[2] -

1 > 2 ✓

l=3 j = 3 - P₂[3] = 2

g[3] > 3 - P₂[3] ✓

1 > 2 ✓

g = [2, 2, 2, 1]

neu az
kacsom
mar
(c)
(a)c
(c)ac

$S = 0$
 while 1: $S \leq n - m$
 $0 \leq 10 - 3 = 7 \checkmark$

T: 1. 2. 3. 4. 5. 6. 7. 8. 9. 10.
 a c b a c b c a c a
 P: c a c 4 || || ||
 c a c c a c

$j = 3$
 while 2: $j > 0 \checkmark$ és $P[3] \stackrel{?}{=} T[0+3]$
 $c \stackrel{?}{=} b \checkmark$

$j \stackrel{?}{=} 0 \checkmark$
 $S = 0 + \max(g[3], 3 - u(T[3]))$
 $u[b] = 0$

$S = \max(2, 3) = 3$
 while 1: $3 \leq 7 \checkmark$

$j = 3$
 while 2: $3 > 0 \checkmark$ és $P[3] \stackrel{?}{=} T[3+3]$
 $c \stackrel{?}{=} b \checkmark$

$j \stackrel{?}{=} 0 \checkmark$

$S = 3 + \max(g[3], 3 - u(T[6]))$
 $u[b] = 0$

$S = 3 + \max(2, 3) = 6$
 while 1: $6 \leq 7 \checkmark$

$j = 3$
 while 2: $3 > 0 \checkmark$ és $P[3] \stackrel{?}{=} T[9]$
 $c \stackrel{?}{=} c \checkmark$

$j = 2$
 while 2: $2 > 0 \checkmark$ és $P[2] \stackrel{?}{=} T[8]$
 $a \stackrel{?}{=} a \checkmark$

$j = 1$
 while 1: $1 > 0 \checkmark$ és $P[1] \stackrel{?}{=} T[7]$
 $c \stackrel{?}{=} c \checkmark$

$j = 0$
 $j \stackrel{?}{=} 0 \checkmark$
 elérkedés az s pozícióra
 $S = 6 + 1 = 7$

"utolsó karakter" huondzika
 b nem fordul elő a
 mintában ezért 3-mal
 arébb tolhatom

itt is ugyan ez

a minta megegyezik
 egy rögz karakter
 sorozattal T-ben