

Haladó informatikai algoritmusok
Mintaillesztés
Véges determinisztikus
automata használata

Fogalmak és jelölések

- Legyen Σ véges jelkészlet, ábécé \Rightarrow **betűk**
- Jelölje Σ^* a Σ elemeiből képezhető összes véges jelsorozat halmazát, beleértve az üres sortozatot is, amit λ -val jelölünk.
- $\Sigma^* = \{x_1 \dots x_n : x_i \in \Sigma, i = 1, \dots, n\} \cup \{\lambda\} \Rightarrow$ **szavak**
- $X \in \Sigma^*$ szó hossza az $X = x_1 \dots x_n$ sorozat elemeinek a száma, jele $|X| = n$.
- A továbbiakban feltételezzük, hogy a szavakat tömb ábrázolja, tehát az X szó i -edik elemére az $X[i]$ tömbelem-kiválasztással hivatkozunk.

Fogalmak és jelölések

- $x = x_1 \dots x_m$, $y = y_1 \dots y_n \in \Sigma^*$ szavak **konkatenációja**:
 $xy = x_1 \dots x_m \cdot y_1 \dots y_n$
- Az u szó kezdő szelete, vagy **prefixe** v -nek, jele $u \sqsubset v$, ha $\exists w \in \Sigma^*$, hogy $uw = v$
- Az u szó végződése, vagy **szuffixe** v -nek, jele $u \sqsupset v$, ha $\exists w \in \Sigma^*$, hogy $wu = v$
- Az $X \in \Sigma^*$ szó első i betűjéből álló prefixére az $X_i = X[1..i]$ jelölést használjuk.
- Általában, egy $X = x_1 \dots x_n \in \Sigma^*$ szó és $1 \leq i \leq j \leq |X|$ esetén $X_{i,j} = X[i..j] = x_i \dots x_j$ jelölést használjuk.
- Azt mondjuk, hogy az S szó **előfordul i eltolással** az A szóban S i eltolással **illeszkedik** az A -ra), ha $A[i + 1..i + m] = S$, ahol $m = |S|$.

Mintaillesztési probléma

Bemenet: $A, S \subseteq \Sigma^*$, A a szöveg, S minta

Kimenet: A legkisebb (összes) olyan i , amelyre $A[i + 1..i + m] = S$, ahol $m = |S|$ a minta hossza

Naiv algoritmus:

1. $n := |A|;$
2. $m := |S|;$
3. **for** $i := 0$ **to** $n - m$ **do**
4. **if** $A[i + 1..i + m] = S[1..m]$ **then begin**
5. $WriteLn(i + 1);$ **Exit**
6. **End;**

A naiv algoritmus futási ideje: $T_{lr}(n, m) = O(n \cdot m)$

Mintaillesztés véges determinisztikus automatával

Véges determinisztikus automata olyan $M = (Q, q_0, F, \Sigma, \delta)$ rendezett ötös, ahol

- Q véges halmaz, az **állapotok halmaza**,
- $q_0 \in Q$ a **kezdőállapot**,
- $F \subseteq Q$ a **végállapotok** (elfogadó állapotok) halmaza,
- Σ véges halmaz, a **bemeneti jelek** halmaza,
- $\delta : Q \times \Sigma \rightarrow Q$, az automata **átmenetfüggvénye**

Jelölje $\delta^* : Q \times \Sigma^* \rightarrow Q$ a δ átmentfüggvény kiterjesztését szavakra:

$$\delta^*(q, w) = \begin{cases} q, & \text{ha } w = \lambda, \\ \delta(\delta^*(q, x_1 \dots x_{n-1}), x_n), & \text{ha } w = x_1 \dots x_n \end{cases}$$

Automata működése - algoritmus

Legyen $M = (Q, q_0, F, \Sigma, \delta)$ olyan véges determinisztikus automata, amely a $\Sigma^* S$ nyelvet ismeri fel, azaz $L(M) = \Sigma^* S$. Mivel $\Sigma^* S$ az összes olyan szavak halmaza, amelyek az S mintára végződnek, ezért $A_i \in L(M)$ akkor és csak akkor, ha az S minta $i - m$ eltolással előfordul az A szövegben. \Rightarrow

Ha van olyan $M = (Q, q_0, F, \Sigma, \delta)$ véges determinisztikus automatánk, amely a $\Sigma^* S$ nyelvet ismeri fel, azaz $L(M) = \Sigma^* S$, akkor az algoritmus mintaillesztést valósít meg:

1. $q := 0$;
2. $n := |A|$;
3. $m := |S|$;
4. **for** $i := 1$ **to** n **do begin**
5. $q := \delta(q, A[i])$;
6. **if** $q \in F$ **then begin** $\{S = A[i - m + 1..i]\}$
7. $WriteLn(i - m + 1)$;
8. **Exit**;
9. **end**;
10. **end**;

Az algoritmus futási ideje, eltekintve az automata előállítási költségétől $O(n)$.

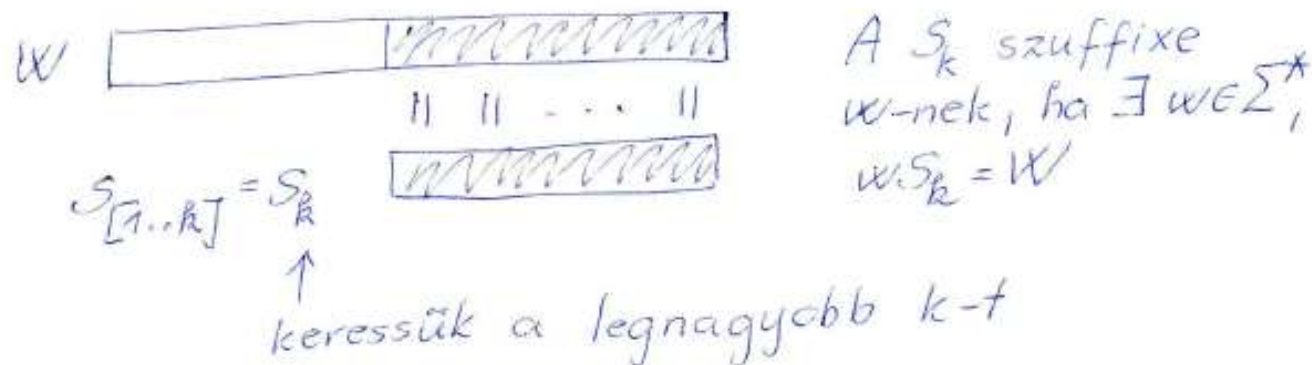
Átmenetfüggvény

Hogyan és mekkora költséggel lehet előállítani az M automatát (átmenetfüggvényét)?

Legyen $\sigma : \Sigma^* \rightarrow \{0, 1, \dots, m\}$ az S minta **szuffix-függvénye**

$$\sigma(W) = \max\{k : S_k \sqsupseteq W\}$$

Az S minta akkor és csak akkor illeszkedik az A szövegre i eltolással, ha $\sigma(A_{i+m}) = m$.



pl: $W = abaacbabaa$
 $\underbrace{abaac} b \dots$
 $S_k \Rightarrow k = 4$

Lemma 1

Ha $q = \sigma(W)$, akkor $\sigma(Wx) = \sigma(S_q x)$ minden $W \in \Sigma^*$ és $x \in \Sigma$ -ra.

Bizonyítás:

Legyen $u = \sigma(Wx)$ és $r = \sigma(S_q x) \Rightarrow$

a legnagyobb olyan u , amelyre $S_u \sqsupseteq Wx \Rightarrow$

$S_{u-1} \sqsupseteq W \Rightarrow$

σ definíciója miatt $u - 1 \leq q \Rightarrow$

$S_{u-1} \sqsupseteq S_q \Rightarrow$

$S_u = S_{u-1}x \sqsupseteq S_q x$

de r a legnagyobb olyan index, hogy $S_r \sqsupseteq S_q x$, tehát $u \leq r$.

Másrészt $S_r \sqsupseteq S_q x \sqsupseteq Wx \Rightarrow r \leq u$.



$$u = r$$

Mire van szükségünk

Olyan automatát akarunk szerkeszteni, amelynek δ átmenetfüggvénye a σ függvényt számítja:

$$\delta^*(0, W) = \sigma(W)$$

Legyen $Q = \{0, 1, \dots, m\}$, ahol $m = |S|$.

Minden $0 \leq q \leq m$ állapotra és $x \in \Sigma$ jelre:

$$\delta(q, x) = \begin{cases} q + 1, & \text{ha } x = S[q + 1], \\ \max\{k : S_k \sqsupseteq S_q x\}, & \text{ha } x \neq S[q + 1] \end{cases}$$

Nyilvánvaló, hogy $\delta(q, x) = \sigma(S_q x)$.

Példa:

$$S = ababc \quad \Sigma = \{a, b, c\}$$

$$\textcircled{1} \quad \begin{cases} \delta(0, a) = 1 & \begin{cases} a \text{?} \checkmark \\ a \text{?} \checkmark \end{cases} \\ \delta(0, b) = 0 & \begin{cases} b \# \\ a \# \end{cases} \\ \delta(0, c) = 0 & \begin{cases} c \# \\ a \# \end{cases} \end{cases}$$

$$\textcircled{2} \quad \begin{cases} \delta(1, a) = 1 & \begin{cases} aa \text{''} \\ a \dots \end{cases} \\ \delta(1, b) = 2 & \begin{cases} ab \text{''} \\ ab \dots \end{cases} \\ \delta(1, c) = 0 & \begin{cases} ac \\ -a \dots \end{cases} \end{cases} \quad x \neq s(q+1)$$

$$\textcircled{3} \quad \begin{cases} \delta(2, a) = 3 & \begin{cases} aba \text{''} \\ ab \text{''} a \text{''} \dots \end{cases} \\ \delta(2, b) = 0 & \begin{cases} abb \\ -a \dots \end{cases} \\ \delta(2, c) = 0 & \begin{cases} abc \\ -a \dots \end{cases} \end{cases}$$

$$\textcircled{4} \quad \begin{cases} \delta(3, a) = 1 & \begin{cases} abaa \text{''} \\ a \dots \end{cases} \\ \delta(3, b) = 4 & \begin{cases} abab \text{''} \\ ab \text{''} ab \text{''} \dots \end{cases} \\ \delta(3, c) = 0 & \begin{cases} abac \\ a \dots \end{cases} \end{cases} \quad x \neq s(q+1)$$

$$\textcircled{5} \quad \delta(4, a) = 3 \quad \begin{cases} ababa \text{''} \\ ab \text{''} a \text{''} \dots \end{cases} \quad x \neq s(q+1)$$

$$\delta(4, b) = 0 \quad \begin{cases} ababb \\ a \dots \end{cases}$$

$$\delta(4, c) = \boxed{5} \quad \begin{cases} ababc \text{''} \\ ab \text{''} abc \text{''} \end{cases}$$

Lemma 2

Minden $W \in \Sigma^*$ szóra $\delta^*(0, W) = \sigma(W)$.

Bizonyítás:

Bizonyítás W hossza szerinti indukcióval.

$W = \lambda$ esetén $\delta^*(0, W) = 0 = \sigma(W)$.

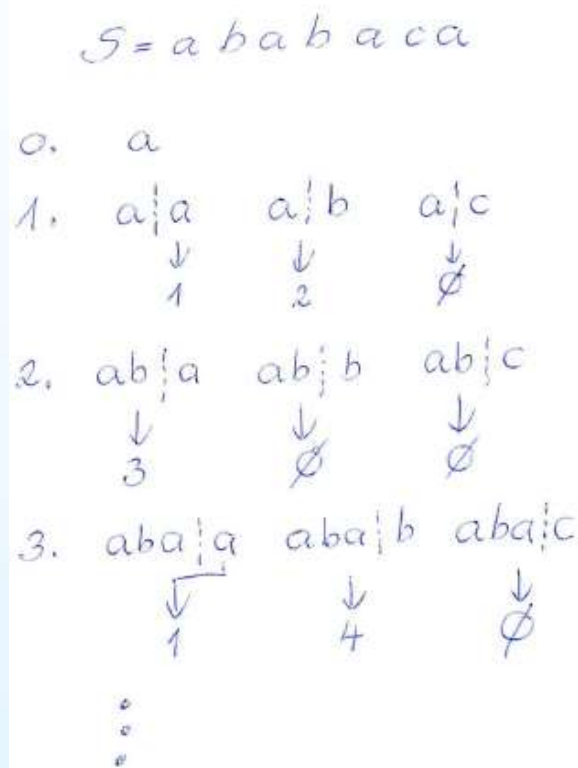
Legyen $W = Vx$, $V \in \Sigma^*$, $x \in \Sigma$ és tegyük fel, hogy $\delta^*(0, V) = \sigma(V) = q$.

$$\begin{aligned}\delta^*(0, Vx) &= \delta(\delta^*(0, V), x) \\ &= \delta(\sigma(V), x) \\ &= \delta(q, x) \\ &= \sigma(S_q x) \\ &= \sigma(Vx)\end{aligned}$$

Az átmenetfüggvény kiszámítása

1. $m := |S|$;
2. **for** $q := 0$ **to** m **do**
3. **for** $x \in \Sigma$ **do begin**
4. $k := \min(m + 1, q + 2)$;
5. **repeat**
6. $k := k - 1$;
7. **until** $S_k \supseteq S_q x$;
8. $\delta(q, x) := k$;
9. **end**;
10. **end**;

Az algoritmus futási ideje $O(m^3|\Sigma|)$.



	a	b	c	S
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

a sikeres illesztésnek megfelelő bejegyzések

Az automata működése:

A szöveg összes $A[i]$ jele alatt feltüntettük az automata $\sigma(A_i)$ állapotát, amelybe az A_i prefix feldolgozása után kerül.

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$A[i]$	a	b	a	b	a	b	a	c	a	b	a	b	a	c	a

$\sigma(A_i)$	0	1	2	3	4	5	4	5	6	7	2	3	4	5	6	7
---------------	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$\max\{k : S_k \supseteq A_i\}$

A minta 2-szer fordul elő a szövegben.

Az átmenetfüggvény kiszámítására szolgáló eljárás működése:

$$S = ababaca$$

$$\Sigma = \{a, b, c\} \quad m = |S| = 7$$

1. $q_i = \phi$
 $x = a \quad k := \min(\overset{m+1}{8}, \overset{q+2}{2}) = 2$

$$k = 1$$

$$S_1 \supset S_0 x ? \quad \delta(0, a) := 1$$

$$a \supset a \quad \checkmark$$

$$x = b \quad k := \min(8, 2) = 2$$

$$k = 1$$

$$S_1 \supset S_0 b \quad k = 0$$

$$a \supset b \quad \times \quad S_0 \supset b \quad \checkmark$$

$$\delta(0, b) := 0$$

$$x = c \quad k = 0$$

$$\delta(0, c) := 0$$

2. $q_i = 1$

$$x = a \quad k := \min(8, 3) = 3$$

$$k = 2$$

$$S_2 \supset S_1 a$$

$$ab \supset aa \quad \times$$

$$k = 1$$

$$S_1 \supset S_1 a \quad \delta(1, a) := 1$$

$$a \supset aa$$

$$x = b \quad k = 3$$

$$k = 2$$

$$S_2 \supset S_1 b \quad \delta(1, b) := 2$$

$$ab \supset ab \quad \checkmark$$

$$x = c \quad k = 3$$

$$k = 2$$

$$S_2 \supset S_1 c$$

$$ab \supset ac \quad \times$$

$$k = 1$$

$$S_1 \supset S_1 c \quad \delta(1, c) := 0$$

$$a \supset ac \quad \times$$

3. $q_i = 2$
 $x = a \quad k := \min(8, 4) = 4$

$$k = 3$$

$$S_3 \supset S_2 a \quad \delta(2, a) := 3$$

$$aba \supset aba \quad \checkmark$$

$$x = b \quad k = 4$$

$$k = 3$$

$$S_3 \supset S_2 b$$

$$aba \supset abb \quad \times$$

$$k = 2$$

$$S_2 \supset S_2 b$$

$$ab \supset abb \quad \times$$

$$k = 1$$

$$S_1 \supset S_2 b \quad \delta(2, b) := 0$$

$$a \supset abb \quad \times$$

$$x = c \quad \delta(2, c) := 0$$

4. $q_i = 3$

$$x = a \quad k := \min(8, 5) = 5$$

$$k = 4$$

$$S_4 \supset S_3 a$$

$$abab \supset abaa \quad \times$$

$$k = 3$$

$$S_3 \supset S_3 a$$

$$aba \supset abaa \quad \times$$

$$k = 2$$

$$ab \supset abaa \quad \times$$

$$k = 1$$

$$a \supset abaa \quad \checkmark$$

$$\delta(3, a) := 1$$

⋮