

Proiectarea algoritmilor: Căutare peste șiruri

Ștefan Ciobâcă, Dorel Lucanu

Faculty of Computer Science
Alexandru Ioan Cuza University, Iași, Romania
dlucanu@info.uaic.ro

PA 2015/2016

1 Algoritmul Knuth-Morris-Pratt

2 Expresii regulate

Plan

1 Algoritmul Knuth-Morris-Pratt

2 Expresii regulate

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				=										
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				=	=									
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				=	=	=								
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				=	=	=	=							
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				=	=	=	=	=						
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Algoritmul naiv¹

<i>b</i>	<i>a</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>a</i>	<i>b</i>
				=	=	=	=	=	≠					
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>				

¹Exemplu din [CLRS]

Intuiția²

b	a	c	b	a	b	a	b	a	a	b	c	b	a	b
				=	=	=	=	=	≠					
				a	b	a	b	a	c	a				
					a	b	a	b	a	c	a			
						a	b	a	b	a	c	a		
							a	b	a	b	a	c	a	

²Exemplu din [CLRS]

Intuiția³

?	?	?	?	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	?	?	?	?	?	?
				=	=	=	=	=	≠					
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	?	?				
					<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	?	?	?			
						<i>a</i>	<i>b</i>	<i>a</i>	?	?	?	?		

³Exemplu din [CLRS]

Intuiția³

?	?	?	?	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	?	?	?	?	?	?
				=	=	=	=	=	≠					
				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	?	?				
					<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	?	?	?			
						<i>a</i>	<i>b</i>	<i>a</i>	?	?	?	?		

Pentru pattern-ul *ababaca*, dacă la o poziție i se potrivesc exact 5 caractere, nu există nicio șansă ca pattern-ul să se potrivească la poziția $i + 1$.

³Exemplu din [CLRS]

Ideea

?	?	?	?	x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	?	?	?	?	?	?	?	?
				=	=	=	=	=	=	=	≠							
				x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	?							
								=	=	=								
								x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	?			

Ideea

?	?	?	?	x ₁	x ₂	x ₃	x ₄	x ₁	x ₂	x ₃	?	?	?	?	?	?	?	?
				=	=	=	=	=	=	=	≠							
				x ₁	x ₂	x ₃	x ₄	x ₁	x ₂	x ₃	?							
								=	=	=								
								x ₁	x ₂	x ₃	x ₄	x ₁	x ₂	x ₃	?			

Ideea

?	?	?	?	x_1	...	x_k	...	x_1	...	x_k	?	?	?	?
				=	=	=	=	=	=	=	≠			
				x_1	...	x_k	...	x_1	...	x_k	?			

Ideea

?	?	?	?	x_1	...	x_k	...	x_1	...	x_k	?	?	?	?
				=	=	=	=	=	=	=	≠			
				x_1	...	x_k	...	x_1	...	x_k	?			

Ne interesează cea mai mare valoare a lui k astfel încât $x_1 \dots x_k$ să fie atât prefix cât și sufix al părții din pattern care s-a potrivit.

Notății

- 1 $s \sqsubseteq t$ dacă s este prefix al lui t

Notății

- ① $s \sqsubseteq t$ dacă s este prefix al lui t
- ② $s \sqsupseteq t$ dacă s este sufix al lui t

Notății

- ① $s \sqsubseteq t$ dacă s este prefix al lui t
- ② $s \sqsupseteq t$ dacă s este sufix al lui t
- ③ Exemple: $\epsilon \sqsubseteq aba$

Notății

- ① $s \sqsubseteq t$ dacă s este prefix al lui t
- ② $s \sqsupseteq t$ dacă s este sufix al lui t
- ③ Exemple: $\epsilon \sqsubseteq aba$, $aba \sqsubseteq ababa$

Notății

- ① $s \sqsubseteq t$ dacă s este prefix al lui t
- ② $s \sqsupseteq t$ dacă s este sufix al lui t
- ③ Exemple: $\epsilon \sqsubseteq aba$, $aba \sqsubseteq ababa$, $aba \sqsupseteq aba$

Notății

- ① $s \sqsubseteq t$ dacă s este prefix al lui t
- ② $s \sqsupseteq t$ dacă s este sufix al lui t
- ③ Exemple: $\epsilon \sqsubseteq aba$, $aba \sqsubseteq ababa$, $aba \sqsubseteq aba$, $aa \sqsupseteq abaa$.

Funcția eșec

❶ $t[0..m-1]$ (patternul)

Funcția eșec

- 1 $t[0..m-1]$ (patternul)
- 2 $s[0..n-1]$ (textul (subiectul))

Funcția eșec

- 1 $t[0..m-1]$ (patternul)
- 2 $s[0..n-1]$ (textul (subiectul))
- 3 Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$

Funcția eșec

- ① $t[0..m-1]$ (patternul)
- ② $s[0..n-1]$ (textul (subiectul))
- ③ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ④ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1						

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0					

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0	0				

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0	0	1			

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0	0	1	2		

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0	0	1	2	3	

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0	0	1	2	3	0

Funcția eșec

- ❶ $t[0..m-1]$ (patternul)
- ❷ $s[0..n-1]$ (textul (subiectul))
- ❸ Funcția eșec $f : \{0, \dots, m-1\} \rightarrow \mathbb{N}$
- ❹ $f(i) =$ cel mai lung prefix propriu al $t[0..i-1]$ care este și sufix al $t[0..i-1]$
- ❺ Exemplu:

i	0	1	2	3	4	5	6
$t[i]$	a	b	a	b	a	c	a
$f(i)$	-1	0	0	1	2	3	0

Exemplu pe tablă.

Algoritmul KMP

1 compute $f[0..m-1]$;

(vom vedea mai târziu cum)

Algoritmul KMP

- ❶ compute $f[0..m-1]$; (vom vedea mai târziu cum)
- ❷ $i = 0$; (încerc să găsesc o potrivire pe poziția $i = 0$)

Algoritmul KMP

- ① compute $f[0..m-1]$; (vom vedea mai târziu cum)
- ② $i = 0$; (încerc să găsesc o potrivire pe poziția $i = 0$)
- ③ $k = 0$; (pentru moment sunt $k = 0$ poziții care se potrivesc)

Algoritmul KMP

- ❶ `compute f[0..m-1];` (vom vedea mai târziu cum)
- ❷ `i = 0;` (încerc să găsesc o potrivire pe poziția $i = 0$)
- ❸ `k = 0;` (pentru moment sunt $k = 0$ poziții care se potrivesc)
- ❹ `while (i + m <= n)`

Algoritmul KMP

- ① `compute f[0..m-1];` (vom vedea mai târziu cum)
- ② `i = 0;` (încerc să găsesc o potrivire pe poziția $i = 0$)
- ③ `k = 0;` (pentru moment sunt $k = 0$ poziții care se potrivesc)
- ④ `while (i + m <= n)`
 - ① `if k == m` (s-au potrivit toate cele m caractere)

Algoritmul KMP

- ① compute $f[0..m-1]$; (vom vedea mai târziu cum)
- ② $i = 0$; (încerc să găesc o potrivire pe poziția $i = 0$)
- ③ $k = 0$; (pentru moment sunt $k = 0$ poziții care se potrivesc)
- ④ while ($i + m \leq n$)
 - ① if $k == m$ (s-au potrivit toate cele m caractere)
 - return i ;

Algoritmul KMP

- ① `compute f[0..m-1];` (vom vedea mai târziu cum)
- ② `i = 0;` (încerc să gădesc o potrivire pe poziția $i = 0$)
- ③ `k = 0;` (pentru moment sunt $k = 0$ poziții care se potrivesc)
- ④ `while (i + m <= n)`
 - ① `if k == m` (s-au potrivit toate cele m caractere)
 - `return i;`
 - ② `else`

Algoritmul KMP

- ① compute $f[0..m-1]$; (vom vedea mai târziu cum)
- ② $i = 0$; (încerc să găesc o potrivire pe poziția $i = 0$)
- ③ $k = 0$; (pentru moment sunt $k = 0$ poziții care se potrivesc)
- ④ while ($i + m \leq n$)
 - ① if $k == m$ (s-au potrivit toate cele m caractere)
 - return i ;
 - ② else
 - ① if $s[i + k] == t[k]$ (încă un caracter se potrivește)

Algoritmul KMP

```

❶ compute f[0..m-1];                                (vom vedea mai târziu cum)
❷ i = 0;                                              (încerc să găesc o potrivire pe poziția  $i = 0$ )
❸ k = 0;      (pentru moment sunt  $k = 0$  poziții care se potrivesc)
❹ while (i + m <= n)
    ❶ if k == m                                     (s-au potrivit toate cele  $m$  caractere)
        • return i;
    ❷ else
        ❶ if s[i + k] == t[k]                       (încă un caracter se potrivește)
            k = k + 1;

```

Algoritmul KMP

- ① compute $f[0..m-1]$; (vom vedea mai târziu cum)
- ② $i = 0$; (încerc să găsesc o potrivire pe poziția $i = 0$)
- ③ $k = 0$; (pentru moment sunt $k = 0$ poziții care se potrivesc)
- ④ while ($i + m \leq n$)
 - ① if $k == m$ (s-au potrivit toate cele m caractere)
 - return i ;
 - ② else
 - ① if $s[i + k] == t[k]$ (încă un caracter se potrivește)
 - $k = k + 1$;
 - ② else (nepotrivire)

Algoritmul KMP

```

❶ compute  $f[0..m-1]$ ;                                (vom vedea mai târziu cum)
❷  $i = 0$ ;                                              (încerc să găsesc o potrivire pe poziția  $i = 0$ )
❸  $k = 0$ ;                                              (pentru moment sunt  $k = 0$  poziții care se potrivesc)
❹ while ( $i + m \leq n$ )
    ❶ if  $k == m$                                        (s-au potrivit toate cele  $m$  caractere)
        • return  $i$ ;
    ❷ else
        ❶ if  $s[i + k] == t[k]$                          (încă un caracter se potrivește)
             $k = k + 1$ ;
        ❷ else                                         (nepotrivire)
             $i = i + k - f[k]$ ;

```

Algoritmul KMP

```

❶ compute f[0..m-1];                                (vom vedea mai târziu cum)
❷ i = 0;                                              (încerc să găesc o potrivire pe poziția i = 0)
❸ k = 0;                                              (pentru moment sunt k = 0 poziții care se potrivesc)
❹ while (i + m <= n)
    ❶ if k == m                                       (s-au potrivit toate cele m caractere)
        • return i;
    ❷ else
        ❶ if s[i + k] == t[k]                       (încă un caracter se potrivește)
            k = k + 1;
        ❷ else                                       (nepotrivire)
            i = i + k - f[k]; k = max(f[k], 0);

```

Algoritmul KMP

```

1 compute f[0..m-1];                                (vom vedea mai târziu cum)
2 i = 0;                                              (încerc să găsesc o potrivire pe poziția  $i = 0$ )
3 k = 0;                                              (pentru moment sunt  $k = 0$  poziții care se potrivesc)
4 while (i + m <= n)
    1 if k == m                                     (s-au potrivit toate cele  $m$  caractere)
        • return i;
    2 else
        1 if s[i + k] == t[k]                     (încă un caracter se potrivește)
            k = k + 1;
        2 else                                     (nepotrivire)
            i = i + k - f[k]; k = max(f[k], 0);
5 return -1;

```


Algoritmul KMP

```

① compute f[0..m-1];                                (vom vedea mai târziu cum)
② i = 0;                                              (încerc să găsesc o potrivire pe poziția  $i = 0$ )
③ k = 0;                                              (pentru moment sunt  $k = 0$  poziții care se potrivesc)
④ while (i + m <= n)
    ① if k == m                                       (s-au potrivit toate cele  $m$  caractere)
        • return i;
    ② else
        ① if s[i + k] == t[k]                       (încă un caracter se potrivește)
            k = k + 1;
        ② else                                       (nepotrivire)
            i = i + k - f[k]; k = max(f[k], 0);
⑤ return -1;

```

Timpul de execuție

<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>
-1	0	0	1	2	3	0

Timpul de execuție

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

- ❶ Observație: pentru orice k , $f[k] < k$.

Timpul de execuție

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

- 1 Observație: pentru orice k , $f[k] < k$.
- 2 Observație la fiecare iterație a buclei `while`, valoarea expresiei $2i + k$ crește cu cel puțin o unitate (vezi pe tablă).

Timpul de execuție

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

- 1 Observație: pentru orice k , $f[k] < k$.
- 2 Observație la fiecare iterație a buclei `while`, valoarea expresiei $2i + k$ crește cu cel puțin o unitate (vezi pe tablă).
- 3 Observație: ce valoarea poate avea $2i + k$ (maxim)?

Timpul de execuție

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

- 1 Observație: pentru orice k , $f[k] < k$.
- 2 Observație la fiecare iterație a buclei `while`, valoarea expresiei $2i + k$ crește cu cel puțin o unitate (vezi pe tablă).
- 3 Observație: ce valoare poate avea $2i + k$ (maxim)?
- 4 Concluzie: câte iterații ale buclei `while` sunt executate?

Algoritmul KMP (prezentare alternativă)

```

KMP(s, n, t, m, f) {
    i = 0;
    k = 0;
    while (i < n) {
        while (k != -1) && (p[k] != s[i])
            k = f[k];
        if (k == m-1)
            return i-m+1; /* gasit p in s */
        else {
            i = i+1;
            k = k+1;
        }
    }
    return -1; /* p nu apare in s */
}

```

Algoritmul KMP (prezentare alternativă)

```

KMP(s, n, t, m, f) {
    i = 0;
    k = 0;
    while (i < n) {
        while (k != -1) && (p[k] != s[i])
            k = f[k];
        if (k == m-1)
            return i-m+1; /* gasit p in s */
        else {
            i = i+1;
            k = k+1;
        }
    }
    return -1; /* p nu apare in s */
}

```

Corelație: variabila i din acest program “ține minte” valoarea $i + k$ din programul precedent.

Calculul funcției eșec

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

Calculul funcției eșec

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

O implementare naivă poate avea complexitatea $O(m^3)$ (exercițiu pentru acasă).

Calculul funcției eșec

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

O implementare naivă poate avea complexitatea $O(m^3)$ (exercițiu pentru acasă).

Dacă presupunem că $f[0..i-1]$ a fost deja calculat, cum calculăm eficient $f[i]$?

Calculul funcției eșec

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

O implementare naivă poate avea complexitatea $O(m^3)$ (exercițiu pentru acasă).

Dacă presupunem că $f[0..i-1]$ a fost deja calculat, cum calculăm eficient $f[i]$?

- 1 Dacă $t[i-1] = t[f[i-1]]$, atunci $f[i] = f[i-1] + 1$.

Calculul funcției eșec

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

O implementare naivă poate avea complexitatea $O(m^3)$ (exercițiu pentru acasă).

Dacă presupunem că $f[0..i-1]$ a fost deja calculat, cum calculăm eficient $f[i]$?

- 1 Dacă $t[i-1] = t[f[i-1]]$, atunci $f[i] = f[i-1] + 1$.
- 2 (altfel,) dacă $t[i-1] = t[f[f[i-1]]]$, atunci $f[i] = f[f[i-1]] + 1$.

Calculul funcției eșec

a	b	a	b	a	c	a
-1	0	0	1	2	3	0

O implementare naivă poate avea complexitatea $O(m^3)$ (exercițiu pentru acasă).

Dacă presupunem că $f[0..i-1]$ a fost deja calculat, cum calculăm eficient $f[i]$?

- ❶ Dacă $t[i-1] = t[f[i-1]]$, atunci $f[i] = f[i-1] + 1$.
- ❷ (altfel,) dacă $t[i-1] = t[f[f[i-1]]]$, atunci $f[i] = f[f[i-1]] + 1$.
- ❸ etc.

Calculul funcției eșec

❶ $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)

Calculul funcției eșec

- ① $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ② $k = -1;$

Calculul funcției eșec

- ① $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ② $k = -1;$
- ③ for ($i = 1; i < m; ++i$) (calculez $f[i]$)

Calculul funcției eșec

- ① $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ② $k = -1;$
- ③ for ($i = 1; i < m; ++i$) (calculez $f[i]$)
 (știi că $t[0..i-2]$ are un prefix propriu de lungime k care este și sufix)

Calculul funcției eșec

- ① $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ② $k = -1;$
- ③ for ($i = 1; i < m; ++i$) (calculez $f[i]$)
 (știu că $t[0..i-2]$ are un prefix propriu de lungime k care este și sufix)
 - ① while $k \neq -1 \ \&\& \ t[k] \neq t[i-1]$

Calculul funcției eșec

- ① $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ② $k = -1;$
- ③ for ($i = 1; i < m; ++i$) (calculez $f[i]$)
 (știi că $t[0..i-2]$ are un prefix propriu de lungime k care este și sufix)
 - ① while $k \neq -1 \ \&\& \ t[k] \neq t[i-1]$
 $k = f[k]$ (iau un prefix mai mic)

Calculul funcției eșec

- ① $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ② $k = -1;$
- ③ for ($i = 1; i < m; ++i$) (calculez $f[i]$)
 (știi că $t[0..i-2]$ are un prefix propriu de lungime k care este și sufix)
 - ① while $k \neq -1 \ \&\& \ t[k] \neq t[i-1]$
 $k = f[k]$ (iau un prefix mai mic)
 - ② $k = k + 1;$ (adaug caracterul $t[k] = t[i-1]$)

Calculul funcției eșec

- ❶ $f[0] = -1;$ (deoarece $t[0..-1] = \epsilon$ nu are niciun prefix propriu)
- ❷ $k = -1;$
- ❸ for ($i = 1; i < m; ++i$) (calculez $f[i]$)
 (știu că $t[0..i-2]$ are un prefix propriu de lungime k care este și sufix)
 - ❶ while $k \neq -1 \ \&\& \ t[k] \neq t[i-1]$
 $k = f[k]$ (iau un prefix mai mic)
 - ❷ $k = k + 1;$ (adaug caracterul $t[k] = t[i-1]$)
 - ❸ $f[i] = k;$

Plan

1 Algoritmul Knuth-Morris-Pratt

2 Expresii regulate

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută.

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*.

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ε , *empty sunt expresii regulate*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*
- *dacă e_1, e_2 sunt expresii regulate, atunci $e_1 e_2$ și $e_1 + e_2$ sunt expresii regulate;*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*
- *dacă e_1, e_2 sunt expresii regulate, atunci $e_1 e_2$ și $e_1 + e_2$ sunt expresii regulate;*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*
- *dacă e_1, e_2 sunt expresii regulate, atunci $e_1 e_2$ și $e_1 + e_2$ sunt expresii regulate;*
- *dacă e este expresie regulată, atunci (e) și e^* sunt expresii regulate.*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ϵ , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*
- *dacă e_1, e_2 sunt expresii regulate, atunci $e_1 e_2$ și $e_1 + e_2$ sunt expresii regulate;*
- *dacă e este expresie regulată, atunci (e) și e^* sunt expresii regulate.*

Definiție

În această secțiune considerăm cazul când "pattern"-ul constituie doar o specificație a ceea ce se caută în sensul că el desemnează o mulțime de șiruri pentru care se caută. Numim o astfel de specificație "*pattern*" *generalizat*. Un alt mod de a specifica "pattern"-uri generalizate îl constituie expresiile regulate.

Definiție

Mulțimea expresiilor regulate peste alfabetul Σ este definită recursiv astfel:

- ε , *empty sunt expresii regulate*
- *orice caracter din Σ este o expresie regulată;*
- *dacă e_1, e_2 sunt expresii regulate, atunci $e_1 e_2$ și $e_1 + e_2$ sunt expresii regulate;*
- *dacă e este expresie regulată, atunci (e) și e^* sunt expresii regulate.*

Arborele sintactic abstract: **pe** **tabla**.

Legătura cu pachetul <regex> din C++

<regex>	expresia regulata
[abc]	$a + b + c$
\d sau [[:digit:]]	$0 + 1 + \dots + 9$
[[:digit:]]*	$(0 + 1 + \dots + 9)^*$
[[:digit:]]+	$(0 + 1 + \dots + 9)(0 + 1 + \dots + 9)^*$

Limbajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) =$

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$

Limbaul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)),

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) =$

Limbaul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$

Limbajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter

Limbajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;

Limbaul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$ atunci $L(e) = L(e_1) \cup L(e_2)$;

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$ atunci $L(e) = L(e_1) \cup L(e_2)$;
- dacă $e = e_1^*$

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$ atunci $L(e) = L(e_1) \cup L(e_2)$;
- dacă $e = e_1^*$ atunci $L(e) = \bigcup_k L(e_1^k)$, unde $L(e_1^0) = \{\varepsilon\}$, $L(e_1^{k+1}) = L(e_1^k)L(e_1)$;

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(empty) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$ atunci $L(e) = L(e_1) \cup L(e_2)$;
- dacă $e = e_1^*$ atunci $L(e) = \bigcup_k L(e_1^k)$, unde $L(e_1^0) = \{\varepsilon\}$, $L(e_1^{k+1}) = L(e_1^k)L(e_1)$;
- dacă $e = (e_1)$

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(\text{empty}) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$ atunci $L(e) = L(e_1) \cup L(e_2)$;
- dacă $e = e_1^*$ atunci $L(e) = \bigcup_k L(e_1^k)$, unde $L(e_1^0) = \{\varepsilon\}$, $L(e_1^{k+1}) = L(e_1^k)L(e_1)$;
- dacă $e = (e_1)$ atunci $L(e) = L(e_1)$.

Exemplu: Fie alfabetul $A = \{a, b, c\}$. Avem $L(a(b + a)c) =$

Limbaajul definit de o expresie regulată

Definiție

Mulțimea de șiruri (limbaajul) $L(e)$ definit de o expresie regulată e este definit recursiv astfel:

- $L(\varepsilon) = \{\varepsilon\}$ (ε este șirul vid (de lungime zero)), $L(\text{empty}) = \emptyset$
- dacă e este un caracter atunci $L(e) = \{e\}$;
- dacă $e = e_1 e_2$ atunci $L(e) = L(e_1)L(e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$;
- dacă $e = e_1 + e_2$ atunci $L(e) = L(e_1) \cup L(e_2)$;
- dacă $e = e_1^*$ atunci $L(e) = \bigcup_k L(e_1^k)$, unde
 $L(e_1^0) = \{\varepsilon\}$, $L(e_1^{k+1}) = L(e_1^k)L(e_1)$;
- dacă $e = (e_1)$ atunci $L(e) = L(e_1)$.

Exemplu: Fie alfabetul $A = \{a, b, c\}$. Avem $L(a(b + a)c) = \{abc, aac\}$ și $L((ab)^*) = \{\varepsilon, ab, abab, ababab, \dots\} = \{(ab)^k \mid k \geq 0\}$.

sfex

Automatul asociat unei expresii regulate

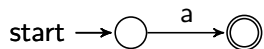
– cazul de baza

e este o litera (un simbol) $a \in \Sigma$

Automatul asociat unei expresii regulate

– cazul de baza

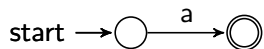
e este o litera (un simbol) $a \in \Sigma$



Automatul asociat unei expresii regulate

– cazul de baza

e este o litera (un simbol) $a \in \Sigma$

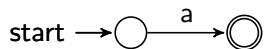


e este ε

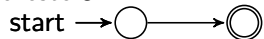
Automatul asociat unei expresii regulate

– cazul de baza

e este o litera (un simbol) $a \in \Sigma$



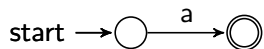
e este ε



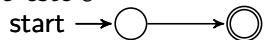
Automatul asociat unei expresii regulate

– cazul de baza

e este o litera (un simbol) $a \in \Sigma$



e este ε

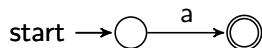


e este *empty*

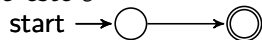
Automatul asociat unei expresii regulate

– cazul de baza

e este o litera (un simbol) $a \in \Sigma$



e este ε



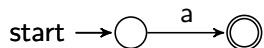
e este *empty*



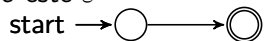
Automatul asociat unei expresii regulate

– cazul de baza

e este o litera (un simbol) $a \in \Sigma$



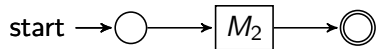
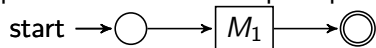
e este ε



e este *empty*



pentru cazul inductiv presupunem:



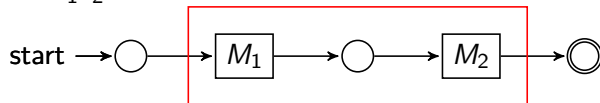
Automatul asociat unei expresii regulate

Automatul asociat unei expresii regulate

$$e = e_1 e_2:$$

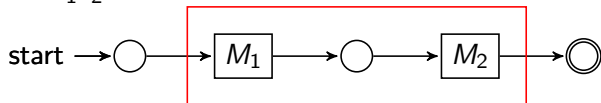
Automatul asociat unei expresii regulate

$e = e_1 e_2$:



Automatul asociat unei expresii regulate

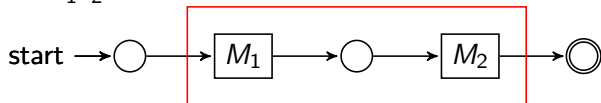
$e = e_1 e_2$:



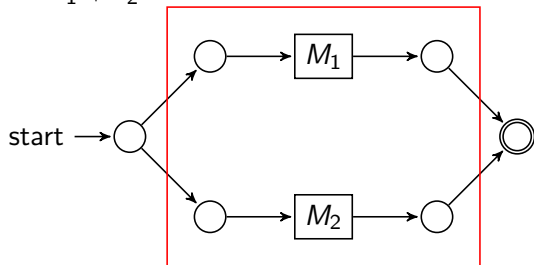
$e = e_1 + e_2$:

Automatul asociat unei expresii regulate

$e = e_1 e_2$:



$e = e_1 + e_2$:

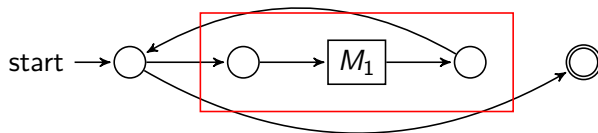


Automatul asociat unei expresii regulate

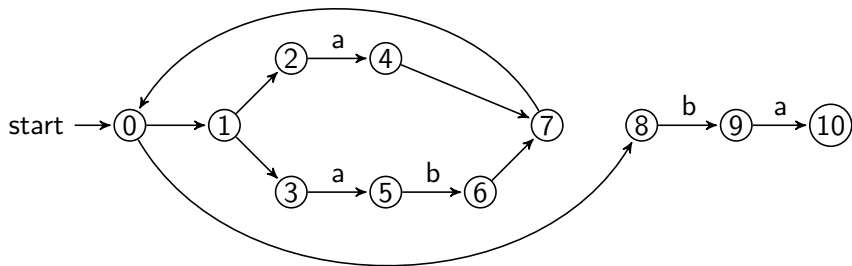
$$e = e_1^*:$$

Automatul asociat unei expresii regulate

$e = e_1^*$:



Exemplu



Detaliile procesului de construcție pe tablă

Automate nedeterminate

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite:

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$,

Automate nedeterminate

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări,

Automate nedeterminate

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul,

Automate nedeterminate

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile,

Automate nedeterminate

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială,

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală
- dacă $M(e)$ este automatul asociat lui e , atunci $L(M(e)) = L(e)$

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală
- dacă $M(e)$ este automatul asociat lui e , atunci $L(M(e)) = L(e)$
- tranzițiile neetichetate se numesc și ε -tranziții (sau spontane)

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală
- dacă $M(e)$ este automatul asociat lui e , atunci $L(M(e)) = L(e)$
- tranzițiile neetichetate se numesc și ε -tranziții (sau spontane)
- automatul construit direct din definiție este în general nedeterminist (și neminimal)

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală
- dacă $M(e)$ este automatul asociat lui e , atunci $L(M(e)) = L(e)$
- tranzițiile neetichetate se numesc și ε -tranziții (sau spontane)
- automatul construit direct din definiție este în general nedeterminist (și neminimal)
- costisitor de aplicat în practică

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală
- dacă $M(e)$ este automatul asociat lui e , atunci $L(M(e)) = L(e)$
- tranzițiile neetichetate se numesc și ε -tranziții (sau spontane)
- automatul construit direct din definiție este în general nedeterminist (și neminimal)
- costisitor de aplicat în practică
- se poate construi un automat echivalent determinist?

Automate nedeterministe

- Automatele asociate expresiilor regulate sunt cazuri particulare de automate finite: $M = (Q, \Sigma, \delta, q_0, Q_f)$, unde Q este mulțimea de stări, Σ alfabetul, $\delta : Q \times A \rightarrow Q$ tranzițiile, $q_0 \in Q$ starea inițială, $Q_f \subseteq Q$ stările finale
- **limbajul acceptat** $L(M)$ este mulțimea de cuvinte ce descriu parcursuri de la starea inițială la o stare finală
- dacă $M(e)$ este automatul asociat lui e , atunci $L(M(e)) = L(e)$
- tranzițiile neetichetate se numesc și ε -tranziții (sau spontane)
- automatul construit direct din definiție este în general nedeterminist (și neminimal)
- costisitor de aplicat în practică
- se poate construi un automat echivalent determinist?
- răspunsul este afirmativ (automatele finite nedeterministe au aceeași putere de acceptare ca și cele deterministe), dar cu anumite costuri (a se vedea slide-urile următoare)

Construcția unui automat determinist echivalent

Fie N un automat nedeterminist cu mulțimea de stări Q . Construim un automat determinist D astfel:

- mulțimea de stări este $\mathcal{P}(Q)$ (număr exponențial de stări!!!)
- există tranziție etichetată cu a de la Q_1 la Q_2 dacă și numai dacă Q_2 este mulțimea tuturor stărilor q_2 cu proprietatea că există $q_1 \in Q_1$ și tranziție etichetată cu a de la q_1 la q_2 în N
- starea inițială a lui D este $\{q_0\}$, unde q_0 este starea inițială a lui N
- o submulțime Q_f este stare finală dacă și numai dacă include o stare finală q_f a lui N

Exemplu pe tablă.

Construcția de mai sus se poate îmbunătăți utilizând un algoritm bazat pe derivatele Brzozowski.

Derivativele Brzozowski

Derivativele unei expresii regulate (Brzozowski, 1964):

$$\delta_a(\text{empty}) = \text{empty}$$

$$\varepsilon?(\text{empty}) = \text{empty}$$

$$\delta_a(\varepsilon) = \text{empty}$$

$$\varepsilon?(\varepsilon) = \varepsilon$$

$$\delta_a(b) = \begin{cases} \varepsilon & , b = a \\ \text{empty} & , b \neq a \end{cases}$$

$$\varepsilon?(b) = \text{empty}$$

$$\delta_a(e_1 e_2) = \delta_a(e_1) e_2 + \varepsilon?(e_1) \delta_a(e_2)$$

$$\varepsilon?(e_1 e_2) = \varepsilon?(e_1) \varepsilon?(e_2)$$

$$\delta_a(e_1 + e_2) = \delta_a(e_1) + \delta_a(e_2)$$

$$\varepsilon?(e_1 + e_2) = \varepsilon?(e_1) + \varepsilon?(e_2)$$

$$\delta_a(e^*) = \delta_a(e) e^*$$

$$\varepsilon?(e^*) = \varepsilon$$

Extensia la cuvinte: $\delta_\varepsilon(e) = e$, $\delta_{wa}(e) = \delta_a(\delta_w(e))$

Simplificări

concatenarea și $+$ sunt asociative, $+$ este și comutativă

$$e + e = e$$

$$e + \textit{empty} = \textit{empty} + e = e$$

$$e \textit{ empty} = \textit{empty} e = \textit{empty}$$

$$e\varepsilon = \varepsilon e = e$$

Proprietatea fundamentală

Theorem (Brzozowski)

Mulțimea derivatelor unei expresii $\{\delta_w(e) \mid w \in A^\}$ este finită.*

Exemplu:

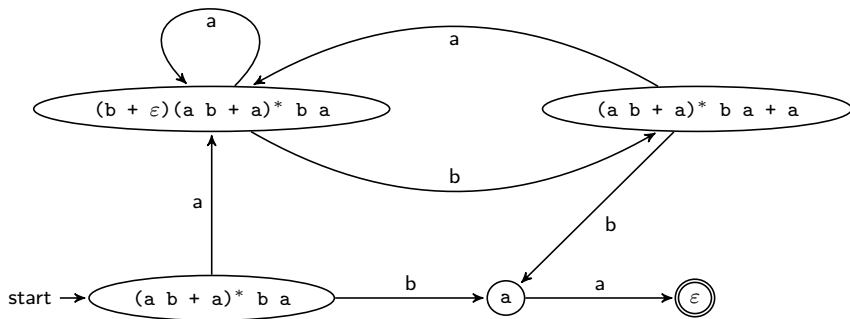
$$\begin{aligned} &\{\delta_w((ab + a)^* b a) \mid w \in A^*\} = \\ &\{(b + \varepsilon)((a b + a)^* b a), (a b + a)^* b a, a, \varepsilon, \text{empty}\} \end{aligned}$$

Construcția automatului (Brzozowski)

- mulțimea de stări este mulțimea derivatelor
- există tranziție etichetată cu a de la q_1 la q_2 dacă și numai dacă q_1 corespunde unei derivate $\delta_w(e)$ și q_2 corespunde derivatei $\delta_{wa}(e)$ pentru un $w \in A^*$;
- starea inițială este $e = \delta_\varepsilon(e)$
- o stare q este finală (de acceptare) dacă și numai dacă corespunde unei derivate $\delta_w(e)$ și $\varepsilon?(\delta_w(e)) = \varepsilon$.

Exemplu

$$e = (a b + a)^* b a$$



Construcția unui automat determinist

(Berry, Setti: From Regular Expressions To Deterministic Automata, 1986)

O **continuare** a lui a în e este orice expresie $\delta_{wa}(e) \neq \text{empty}$.

- se marchează simbolurile din e ca fiind distincte; fie e' expresia obținută (de exemplu $e = (ab + b)^*ba$ este transformată în $e' = (a_1b_2 + b_3)^*b_4a_5$)
- se construiește automatul M' pentru e' urmând ideea din algoritmul lui Brzozowski:
 - M' are o stare pentru fiecare continuare a unui simbol marcat în e'
 - există tranziție de la q_1 la q_2 dacă și numai dacă q_1 corespunde unei continuări C , C, C poate genera un cuvânt care începe cu a ($\delta_{wa}(e) \neq \text{empty}$) și q_2 corespunde continuării lui a
 - starea inițială este e'
 - q este o stare finală (de acceptare) dacă și numai dacă ea corespunde unei continuări C și $\varepsilon?(C) = \varepsilon$
- se elimină mărcile din M'
- se determinizează M' construind M ale cărui stări sunt submulțimi de stări ale lui M'

Construcții mai performante

- utilizând funcțiile `first` și `follow` (Berry, Setti, 1986)
- paralelizare (Myer, A Four Russians Algorithm for Regular Expression Pattern Matching)
- o altă construcție pentru automatul nedeterminist este Glushkov-McNaughton-Yamada (1960-1961), care poate fi și paralelizată (Navarro & Raffinot, 2004)

Complexitatea căutării cu expresii regulate

Presupunem că lungimea expresiei regulate este m (numărul de caractere fără operatori) și $m_{\Sigma} = |\Sigma \cup \{., +, *\}|$.

Theorem (Thomson, 1968)

Problema căutării cu expresii regulate poate fi rezolvată în timpul $O(mn)$ cu automate nedeterminate și spațiu $O(m)$.

Theorem (Kleene, 1956)

Problema căutării cu expresii regulate poate fi rezolvată în timpul $O(n + 2^{m_{\Sigma}})$ cu automate deterministe și spațiu $O(2^{m_{\Sigma}})$.

Theorem (Myers, 1992)

Problema căutării cu expresii regulate poate fi rezolvată în timpul $O(mn / \log n)$ cu automate deterministe și spațiu $O(mn / \log n)$.