

Algoritmo de força bruta

```
pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
```

Algoritmo de força bruta

```
pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
```

Algoritmo de força bruta

```
pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a
```

Algoritmo de força bruta

```
pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a
4 a b a b b a b a b b a
5 a b a b b a b a b b a
6 a b a b b a b a b b a
7 a b a b b a b a b b a
8 a b a b b a b a b b a
9 a b a b b a b a b b a
10 a b a b b a b a b b a
11 a b a b b a b a b b a
12 a b a b b a b a b b a
```

Algoritmo de força bruta

i	j	i+j	0	1	2	3	4	5	6	7	8	9	10
			A	B	A	C	A	D	A	B	R	A	C
0	2	2	A	B	R	A							
1	0	1	A	B	R	A							
2	1	3	A	B	R	A							
3	0	3	A	B	R	A							
4	1	5		A	B	R	A						
5	0	5			A	B	R	A					
6	4	10				A	B	R	A				

txt →

← pat

entries in red are mismatches

entries in gray are for reference only

entries in black match the text

↑ match

return i when j is M

Brute-force substring search

Algoritmo de força bruta

Devolve a primeira de ocorrências de pat em txt.

```
public static
int search(String pat, String txt) {
    int i, n = txt.length();
    int j, m = pat.length();
    for (i = 0; i <= n-m; i++) {
        for (j = 0; j < m; j++)
            if (txt.charAt(i+j) != pat.charAt(j))
                break;
        if (j == m) return i;
    }
    return n;
}
```

Algoritmo de força bruta

Relação invariante: no início de cada iteração do "for (j= 0; ...)" vale que

$$(i0) \text{ pat}[0..j-1] = \text{txt}[i..i+j-1]$$

Pior caso

pat = a a a a a a a a a a b

```

0 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a a a a a a a a a a a a a a a a a a a a a a a txt
0 a a a a a a a a a a b
1  a a a a a a a a a a b
2   a a a a a a a a a a b
3    a a a a a a a a a a b
4     a a a a a a a a a a b
5      a a a a a a a a a a b
6       a a a a a a a a a a b
7        a a a a a a a a a a b
8         a a a a a a a a a a b
9          a a a a a a a a a a b
10           a a a a a a a a a a b
11            a a a a a a a a a a b
12             a a a a a a a a a a b
    
```

Melhor caso

pat = b a a a a a a a a a a

```

0 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a a a a a a a a a a a a a a a a a a a a a a a txt
0 b a a a a a a a a a a
1  b a a a a a a a a a a
2   b a a a a a a a a a a
3    b a a a a a a a a a a
4     b a a a a a a a a a a
5      b a a a a a a a a a a
6       b a a a a a a a a a a
7        b a a a a a a a a a a
8         b a a a a a a a a a a
9          b a a a a a a a a a a
10           b a a a a a a a a a a
11            b a a a a a a a a a a
12             b a a a a a a a a a a
    
```

Consumo de tempo

Consumo de tempo da função search().

linha todas as execuções da linha

$$\begin{aligned}
 1-2 &= 1 \\
 3 &= n - m + 1 \\
 4 &\leq (n - m + 1)(m + 1) \\
 5 &\leq (n - m + 1)m \\
 6 &\leq (n - m + 1) \\
 7 &= n - m \\
 8-9 &= 1
 \end{aligned}$$

$$\begin{aligned}
 \text{total} &< 3(n - m + 2) + 2(n - m + 1)(m + 1) \\
 &= O((n - m + 1)m)
 \end{aligned}$$

Pior caso

i	j	i+j	0	1	2	3	4	5	6	7	8	9
			txt → A A A A A A A A A A B									
0	4	4	A	A	A	A	B ← pat					
1	4	5		A	A	A	A	B				
2	4	6			A	A	A	A	B			
3	4	7				A	A	A	A	B		
4	4	8					A	A	A	A	B	
5	5	10						A	A	A	A	B

Brute-force substring search (worst case)

Conclusões

O consumo de tempo de search() no pior caso é $O((n - m + 1)m)$.

O consumo de tempo de search() no melhor caso é $O(n - m + 1)$.

Isto significa que no pior caso o consumo de tempo é essencialmente proporcional a $m \cdot n$.

Em geral o algoritmo é rápido e faz não mais que $1.1 \times n$ comparações.

Algoritmo de força bruta: versão alternativa

```

public static
int search(String pat, String txt) {
    int i, n = txt.length();
    int j, m = pat.length();
    for (i=0, j=0; i < n && j < m; i++) {
        if(txt.charAt(i)==pat.charAt(j))j++;
        else {
            i -= j; // retrocesso
            j = 0;
        }
    }
    if (j == m) return i - m;
    return n; }

```

Navigation icons

Algoritmo força bruta: direita para esquerda

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
0 a b a b b a b a b b a

```

Navigation icons

Algoritmo força bruta: direita para esquerda

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a

```

Navigation icons

Algoritmo força bruta: direita para esquerda

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a

```

Navigation icons

Algoritmo força bruta: direita para esquerda

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a

```

Navigation icons

Algoritmo força bruta: direita para esquerda

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a
4 a b a b b a b a b b a

```

Navigation icons

Algoritmo força bruta: direita para esquerda

pat = a b a b b a b a b b a

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a
4 a b a b b a b a b b a
5 a b a b b a b a b b a

```

Algoritmo força bruta: direita para esquerda

pat = a b a b b a b a b b a

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a
4 a b a b b a b a b b a
5 a b a b b a b a b b a
6 a b a b b a b a b b a

```

Navigation icons

Navigation icons

Algoritmo força bruta: direita para esquerda

pat = a b a b b a b a b b a

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b b a txt
0 a b a b b a b a b b a
1 a b a b b a b a b b a
2 a b a b b a b a b b a
3 a b a b b a b a b b a
4 a b a b b a b a b b a
5 a b a b b a b a b b a
6 a b a b b a b a b b a
7 a b a b b a b a b b a
8 a b a b b a b a b b a
9 a b a b b a b a b b a
10 a b a b b a b a b b a
11 a b a b b a b a b b a
12 a b a b b a b a b b a

```

Navigation icons

Força bruta: direita para esquerda

Devolve a primeira de ocorrências de pat em txt.

```

public static
int search(String pat, String txt) {
    int i, n = txt.length();
    int j, m = pat.length();
    for (i = 0; i <= n-m; i+=1 /*skip*/) {
        for (j = m-1; j >= 0; j--)
            if(txt.charAt(i+j)!=pat.charAt(j))
                break;
        if (j == -1) return i;
    }
    return n;
}

```

Navigation icons

Próximos passos

Existe algoritmo mais rápido que o força bruta?

Existe algoritmo que faz apenas n comparações entre caracteres?

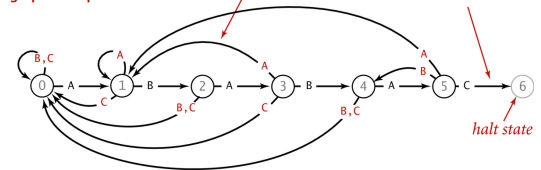
Existe algoritmo que faz menos que n comparações?

Algoritmo KMP para busca de substring

internal representation

	j	0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6

graphical representation



DFA corresponding to the string A B A B A C

Navigation icons

Navigation icons

Algoritmo de força bruta

```

P = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a T
0 a b a b
1 a
2 a b
3 a b a b b
4 a
5 a b a b b a b a b b
6 a
7 a b a
8 a
9 a
10 a b a b
11 a
12 a b a b b a b a b b a
    
```

Navigation icons

Algoritmo de força bruta: versão alternativa

```

public static
int search(String pat, String txt) {
    int i, n = txt.length();
    int j, m = pat.length();
    for(i=0, j=0; i < n && j < m; i++) {
        if(txt.charAt(i)==pat.charAt(j))j++;
        else {
            i -= j; // retrocesso
            j = 0;
        }
    }
    if (j == m) return i - m;
    return n; }
    
```

Navigation icons

Ideia básica do algoritmo

```

P = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a T
0 a
    
```

Navigation icons

Ideia básica do algoritmo

```

P = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a T
0 a
1 a b
    
```

Navigation icons

Ideia básica do algoritmo

```

P = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a T
0 a
1 a b
2 a b a
    
```

Navigation icons

Ideia básica do algoritmo

```

P = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a T
0 a
1 a b
2 a b a
3 a
    
```

Navigation icons

Ideia básica do algoritmo

$P = a b a b b a b a b b a$
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b b a b a b b a T

0 a
 1 a b
 2 a b a
 3 a
 4 a b



Ideia básica do algoritmo

$P = a b a b b a b a b b a$
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b b a b a b b a T

0 a
 1 a b
 2 a b a
 3 a
 4 a b
 5 a b a



Ideia básica do algoritmo

$P = a b a b b a b a b b a$
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b b a b a b b a T

0 a
 1 a b
 2 a b a
 3 a
 4 a b
 5 a b a
 6 a b a b



Ideia básica do algoritmo

$P = a b a b b a b a b b a$
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b b a b a b b a T

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



Ideia básica do algoritmo

$P = a b a b b a b a b b a$
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b b a b a b b a T

0 a
 1 a b
 2 a b a
 3 a
 4 a b
 5 a b a
 6 a b a b
 7 a b a
 8 a b a b



Ideia básica do algoritmo

$P = a b a b b a b a b b a$
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b b a b a b b a T

0 a
 1 a b
 2 a b a
 3 a
 4 a b
 5 a b a
 6 a b a b
 7 a b a
 8 a b a b
 9 a b a b b



Ideia básica do algoritmo

$P = a b a b b a b a b b a$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b a b b a T

```

0 a
1 a b
2 a b a
3   a
4   a b
5   a b a
6   a b a b
7     a b a
8     a b a b
9     a b a b b
10    a b a b b a

```



Ideia básica do algoritmo

$P = a b a b b a b a b b a$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b a b b a T

```

0 a
1 a b
2 a b a
3   a
4   a b
5   a b a
6   a b a b
7     a b a
8     a b a b
9     a b a b b
10    a b a b b a
11    a b a b b a b

```



Ideia básica do algoritmo

$P = a b a b b a b a b b a$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b b a b a b b a T

```

0 a
1 a b
2 a b a
3   a
4   a b
5   a b a
6   a b a b
7     a b a
8     a b a b
9     a b a b b
10    a b a b b a
11    a b a b b a b
12    a b a b b a b a

```



Ideia básica do algoritmo

$P = a b a b b a b a b b a$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b b a b a b b a T

```

0 a
1 a b
2 a b a
3   a
4   a b
5   a b a
6   a b a b
7     a b a
8     a b a b
9     a b a b b
10    a b a b b a
11    a b a b b a b
12    a b a b b a b a
13    a b a b b a b a b

```



Ideia básica do algoritmo

$P = a b a b b a b a b b a$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b b a b a b b a T

```

0 a
1 a b
2 a b a
3   a
4   a b
5   a b a
6   a b a b
7     a b a
8     a b a b
9     a b a b b
10    a b a b b a
11    a b a b b a b
12    a b a b b a b a
13    a b a b b a b a b
14    a b a b b a b a b a

```



Ideia básica do algoritmo

$P = a b a b b a b a b b a$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b b a b a b b a T

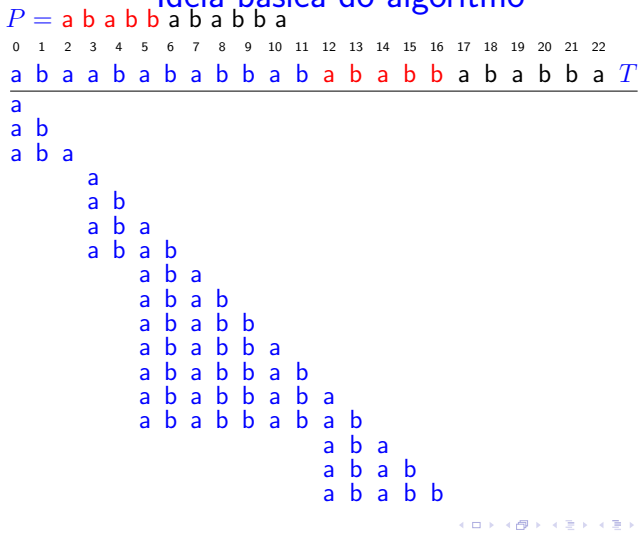
```

0 a
1 a b
2 a b a
3   a
4   a b
5   a b a
6   a b a b
7     a b a
8     a b a b
9     a b a b b
10    a b a b b a
11    a b a b b a b
12    a b a b b a b a
13    a b a b b a b a b
14    a b a b b a b a b a
15    a b a b b a b a b a b

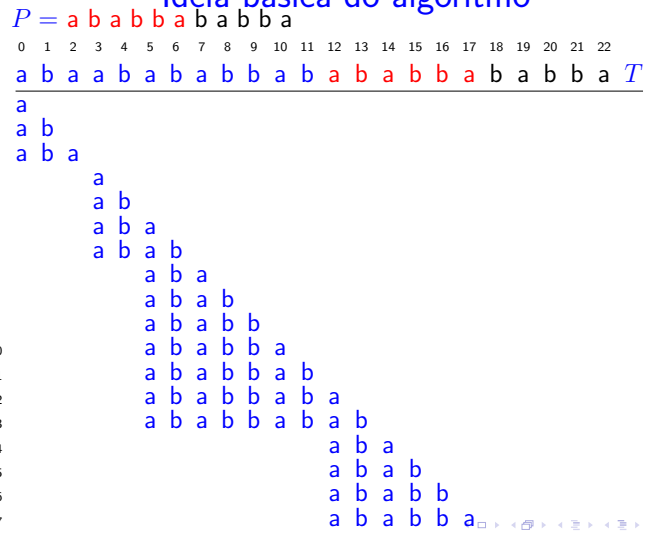
```



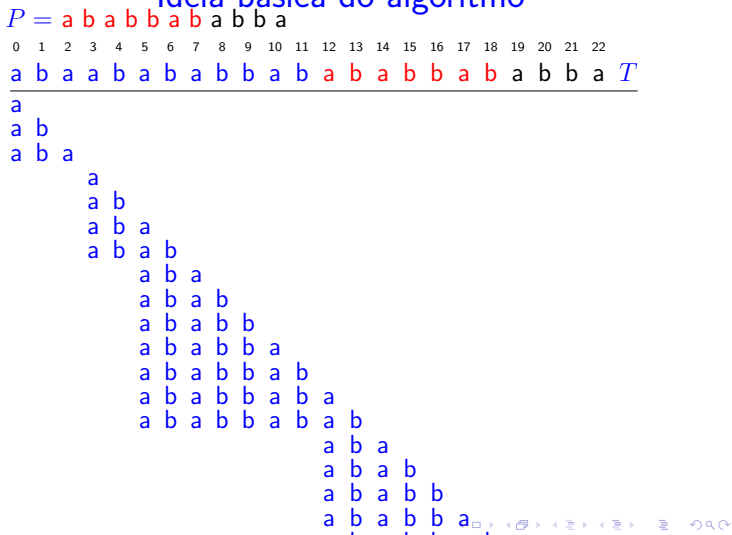
Ideia básica do algoritmo



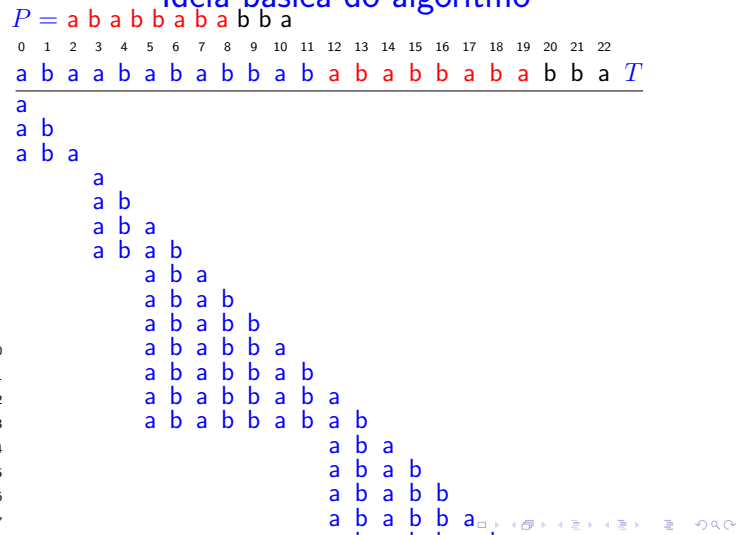
Ideia básica do algoritmo



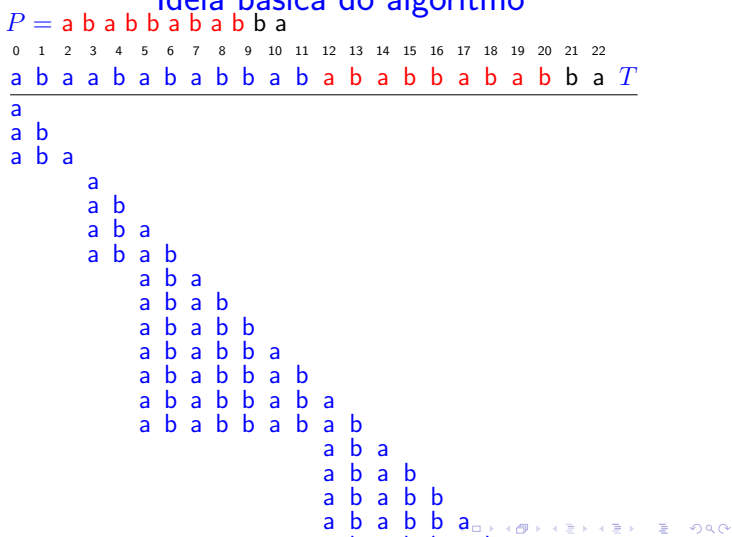
Ideia básica do algoritmo



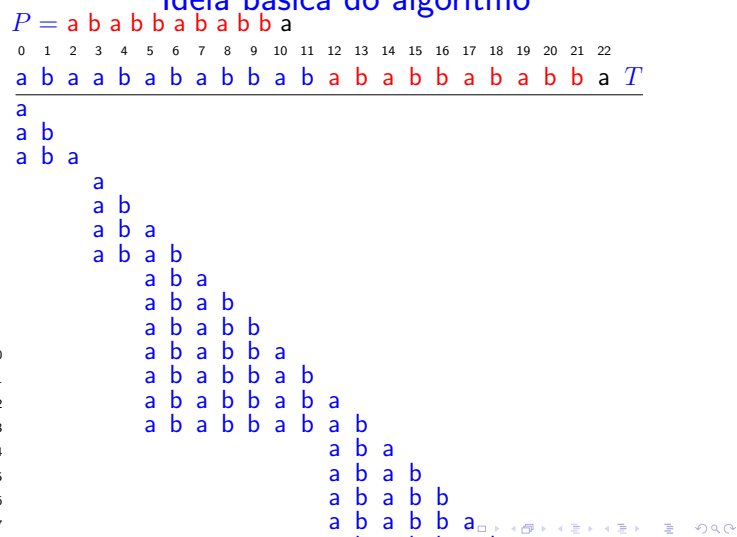
Ideia básica do algoritmo



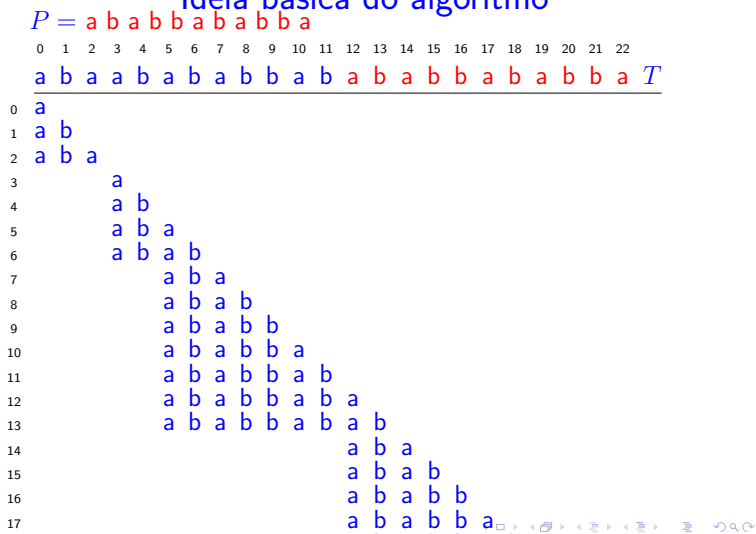
Ideia básica do algoritmo



Ideia básica do algoritmo



Ideia básica do algoritmo



Ideia geral

Quando encontramos um conflito entre $txt[i]$ e $pat[j]$, **não** é necessário retroceder i e passar a comparar $txt[i-j+1..]$ com $pat[0..]$.

Basta:

encontrar o comprimento do maior prefixo de $pat[0..]$ que é sufixo de $txt[..i]$,

ou seja,

encontrar o maior k tal que $pat[0..k-1]$ é igual a $txt[i-k+1..i]$ que é igual a $pat[j-k+1..j-1]+txt[i]$,

e passar a comparar $txt[i+1..]$ com $pat[k..]$.

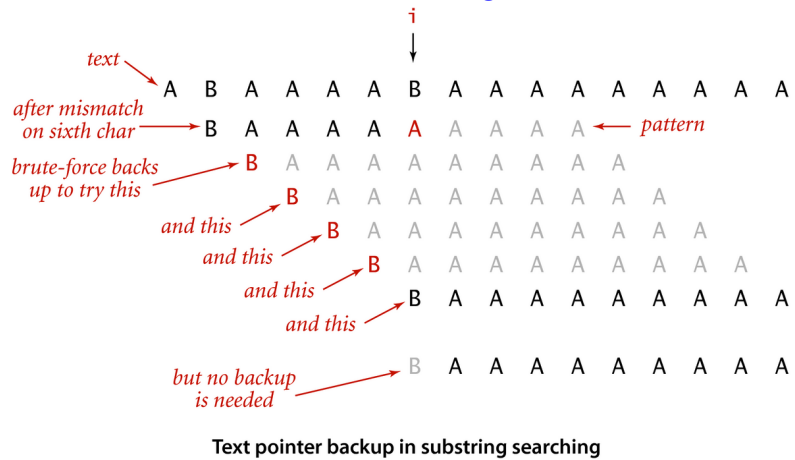
Algoritmo KMP

Examina os caracteres de txt um a um, da esquerda para a direita, **sem nunca retroceder**.

Em cada iteração, o algoritmo sabe qual posição k de pat deve ser emparelhada com a próxima posição $i+1$ de txt .

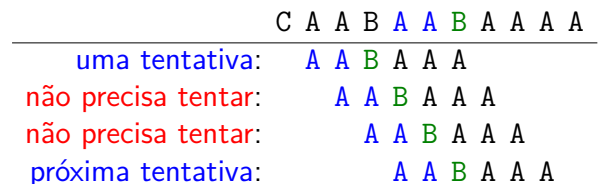
Ou seja, no fim de cada iteração, o algoritmo sabe qual índice k deve fazer o papel de j na próxima iteração.

Ideia básica do algoritmo



Ideia geral

Exemplo: texto CAABAABAAAA e padrão AABAAA: depois do conflito entre $txt[6]$ e $pat[5]$, não precisamos retroceder no texto: podemos continuar e comparar $txt[7..]$ com $pat[3..]$:



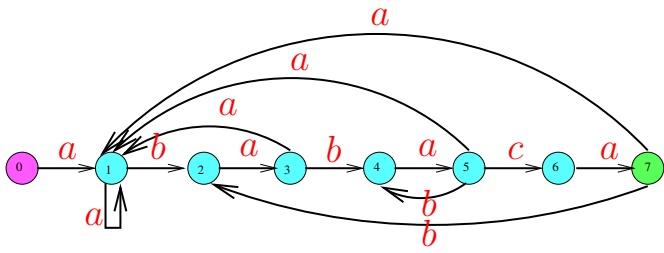
Algoritmo KMP

O algoritmo KMP usa uma tabela $dfa[][]$ que armazena os índices mágicos k .

O nome da tabela deriva da expressão *deterministic finite-state automaton*.

As colunas da tabela são indexadas pelos índices $0..m-1$ do padrão e as linhas são indexadas pelo alfabeto, que é o conjunto de todos os caracteres do texto e do padrão.

Autômato de estados determinístico (DFA)



0..7 = conjunto de estados

$\Sigma = \{a, b, c\}$ = alfabeto

δ = função de transição

0 é estado inicial e 7 é estado final

Autômato finito determinístico (DFA)

O algoritmo KMP simula o funcionamento do autômato de estados.

O autômato começa no estado 0 e examina os caracteres do texto, um de cada vez, da esquerda para a direita, mudando para um novo estado cada vez que lê um caractere do texto.

Se atingir o estado m , dizemos que o autômato reconheceu ou aceitou o padrão.

Se chegar ao fim do texto sem atingir o estado m , sabemos que o padrão não ocorre no texto.

Autômato finito determinístico (DFA)

Um autômato finito é formado uma 5-upla $(Q, \Sigma, \delta, q_0, F)$, onde

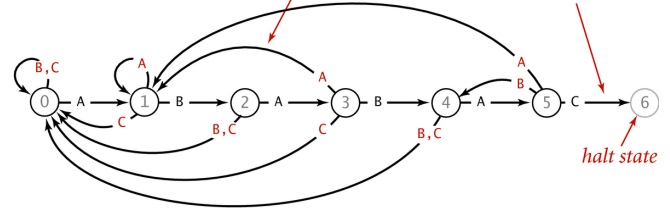
- ▶ Q é um conjunto finitos de estados,
- ▶ Σ é um conjunto finito chamado alfabeto,
- ▶ $\delta : Q \times \Sigma \rightarrow Q$ é a função de transição,
- ▶ $q_0 \in Q$ é o estado inicial, e
- ▶ $F \subseteq Q$ é o conjunto de aceitação.

Exemplo: pat = ABABAC

internal representation

j	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

graphical representation



Autômato finito determinístico (DFA)

O autômato está no estado j se acabou de casar os j primeiros caracteres do padrão com um segmento do texto, ou seja, se acabou de casar $pat[0..j-1]$ com $txt[i-j..i-1]$.

Para cada estado j , a transição que corresponde ao caractere $pat[j]$ é de casamento e leva ao estado $j+1$.

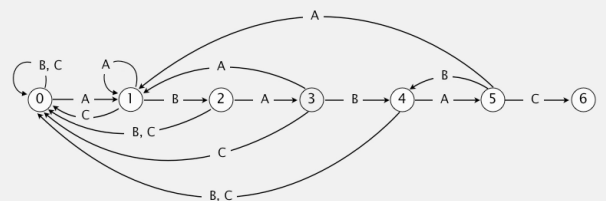
Todas as outras transições que começam no estado j são de conflito e levam a um estado $\leq j$.

O autômato de estados é uma ideia muito importante em compilação, na teoria da computação, etc.

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A

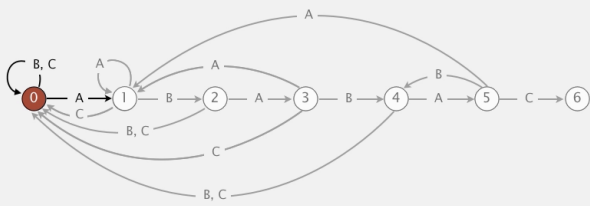
pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6



Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A
↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

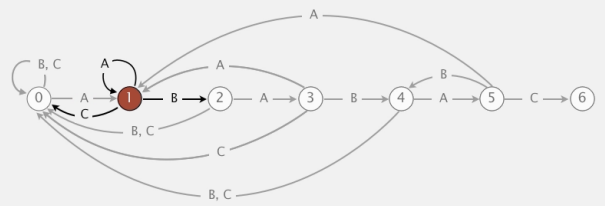


4

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A
↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

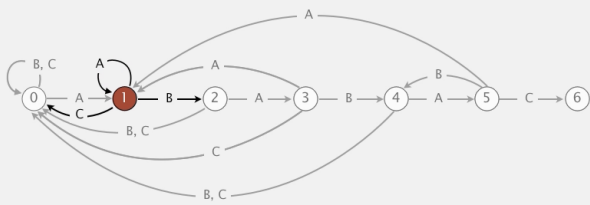


5

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A
↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

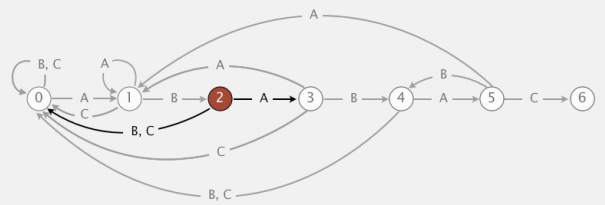


6

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A
↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

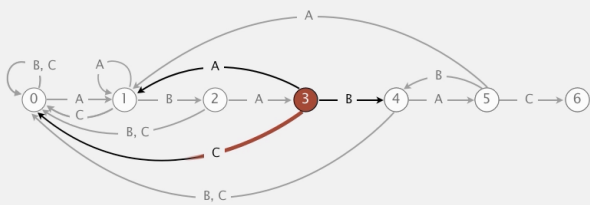


7

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A
↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

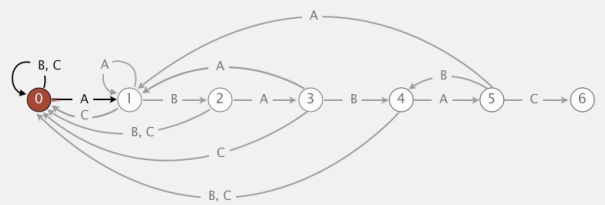


8

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A
↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

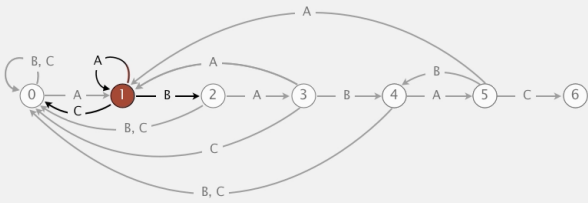


9

Knuth-Morris-Pratt demo: DFA simulation

A A B A C **A** A B A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

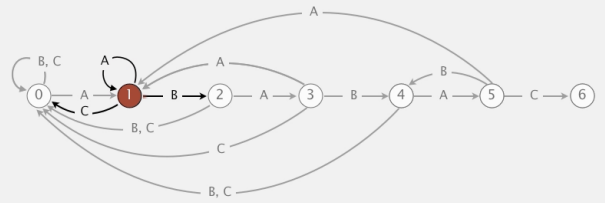


10

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A **A** B A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

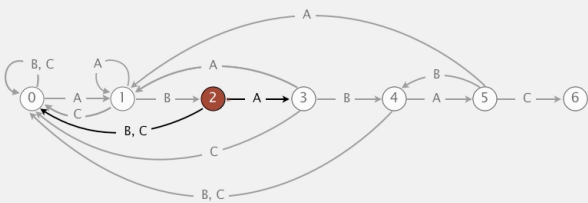


11

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A **A** B A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

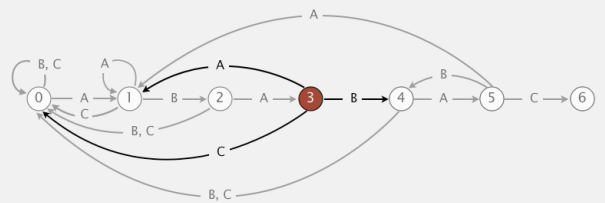


12

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A **B** A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

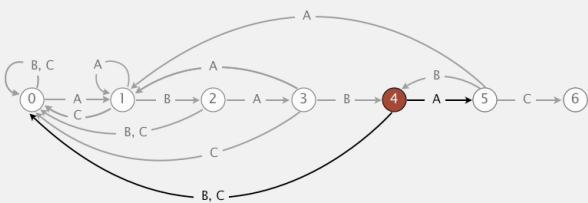


13

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A **B** A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

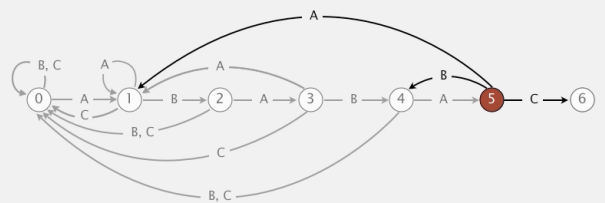


14

Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A **B** A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

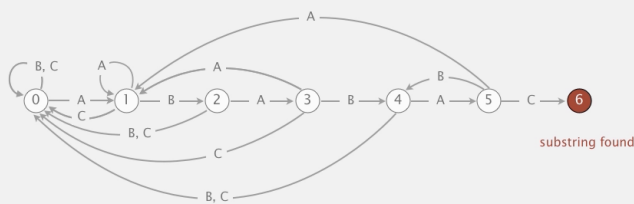


15

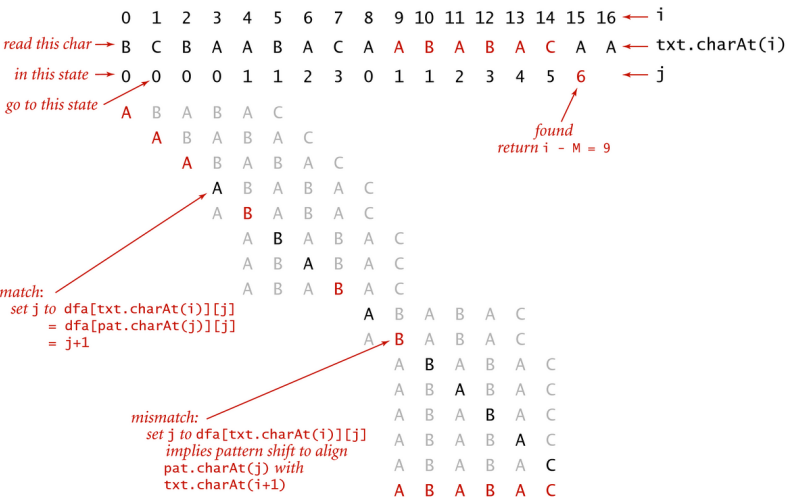
Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6



16



Trace of KMP substrings search (DFA simulation) for A B A B A C

Autômato de estados determinístico (DFA)

A tabela `dfa[][]` representa uma máquina imaginária conhecida como **autômato de estados** (*deterministic finite-state automaton, DFA*).

Os estados do autômato correspondem aos índices $0 \dots m-1$ de `pat`.

Também há um estado *final* `m`.

Para cada estado e cada caractere do **alfabeto**, há uma **transição** que leva desse estado a um outro.

Algoritmo KMP

Retorna a posição a partir de onde `pat` ocorre em `txt` se `pat` não ocorre em `txt` retorna `n`.

```
public int search(String txt) {
    int i, n = txt.length();
    int j, m = pat.length();

    for (i = 0, j = 0; i < n && j < m; i++)
        j = dfa[txt.charAt(i)][j];

    if (j == m) return i - m;
    return n;
}
```

Invariantes

O método `search()` de **KMP** tem os seguintes **invariantes**.

Imediatamente antes do teste `i < n && j < m` vale que:

- ▶ `pat` não ocorre em `txt[0 .. i-1]`;
- ▶ `pat[0 .. k]` é diferente de `txt[i-k .. i]` para todo `k` no conjunto $j+1 \dots m-1$; e
- ▶ `pat[0 .. j-1]` é igual a `txt[i-j .. i-1]`.

Construção do DFA

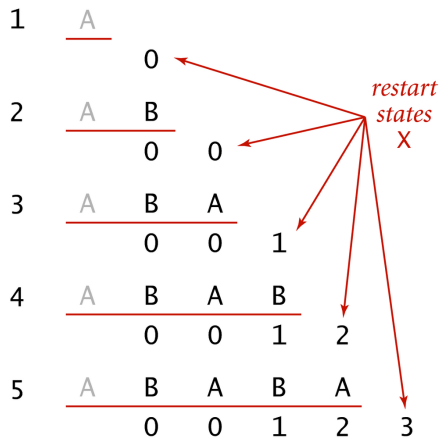
Para construir a tabela `dfa[][]` que representa o autômato podemos pré-processar o padrão `pat` desde que o **alfabeto** de `txt` seja conhecido.

Para qualquer caractere `c` do **alfabeto** e qualquer `j` em $0 \dots m-1$, o valor de `dfa[c][j]` é

o comprimento do **maior prefixo** de `pat[0 .. j]` que é **sufixo** de `pat[0 .. j-1]+c`.

Uma implementação literal dessa definição faria cerca de Rm^3 comparações entre caracteres para calcular a tabela `dfa[][]`, sendo `R` o número de caracteres do **alfabeto**.

Exemplo: padrão ABABAC e alfabeto A B C



DFA simulations to compute

Knuth-Morris-Pratt demo: DFA construction

Include one state for each character in pattern (plus accept state).

pat.charAt(j)	0	1	2	3	4	5
	A	B	A	B	A	C
dfa[][j]	A					
	B					
	C					

Constructing the DFA for KMP substring search for A B A B A C

Knuth-Morris-Pratt demo: DFA construction

Include one state for each character in pattern (plus accept state).

pat.charAt(j)	0	1	2	3	4	5
	A	B	A	B	A	C
dfa[][j]	A					
	B					
	C					

Constructing the DFA for KMP substring search for A B A B A C



Knuth-Morris-Pratt demo: DFA construction

Match transition. If in state j and next char $c == \text{pat.charAt}(j)$, go to $j+1$.

first j characters of pattern have already been matched next char matches now first $j+1$ characters of pattern have been matched

pat.charAt(j)	0	1	2	3	4	5
	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	2		4		
	C					6

Constructing the DFA for KMP substring search for A B A B A C



Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c != \text{pat.charAt}(j)$.

pat.charAt(j)	0	1	2	3	4	5
	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	2		4		
	C					6

Constructing the DFA for KMP substring search for A B A B A C



Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c != \text{pat.charAt}(j)$.

pat.charAt(j)	0	1	2	3	4	5
	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0				6

Constructing the DFA for KMP substring search for A B A B A C



Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0				6

Constructing the DFA for KMP substrng search for A B A B A C



21

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0	0			6

Constructing the DFA for KMP substrng search for A B A B A C



21

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0	0			6

Constructing the DFA for KMP substrng search for A B A B A C



22

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0	0	0		6

Constructing the DFA for KMP substrng search for A B A B A C



22

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0	0	0		6

Constructing the DFA for KMP substrng search for A B A B A C



23

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	5		
	B	0	2	4		
	C	0	0	0		6

Constructing the DFA for KMP substrng search for A B A B A C



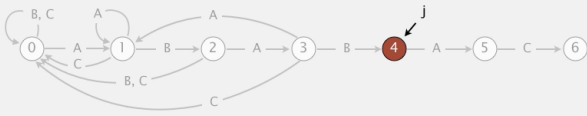
23

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

pat.charAt(j)	0	1	2	3	4	5
A	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	
	C	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



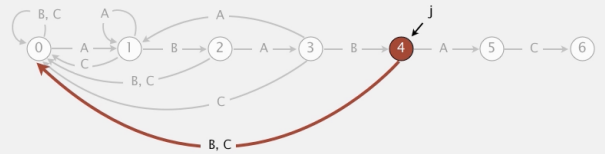
24

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

pat.charAt(j)	0	1	2	3	4	5
A	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	0

Constructing the DFA for KMP substring search for A B A B A C



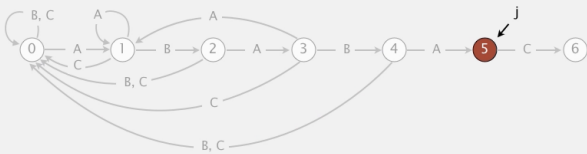
24

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

pat.charAt(j)	0	1	2	3	4	5
A	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



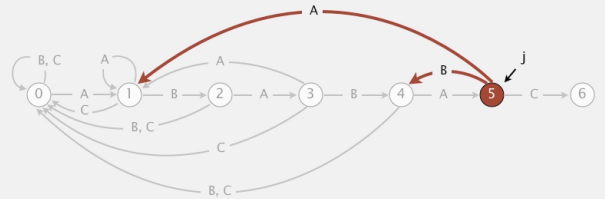
25

Knuth-Morris-Pratt demo: DFA construction

Mismatch transition: back up if $c \neq \text{pat.charAt}(j)$.

pat.charAt(j)	0	1	2	3	4	5
A	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C

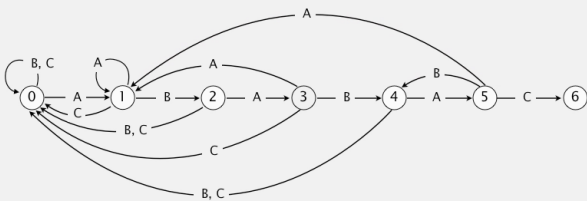


25

Knuth-Morris-Pratt demo: DFA construction

pat.charAt(j)	0	1	2	3	4	5
A	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



26

Knuth-Morris-Pratt demo: DFA construction in linear time

Include one state for each character in pattern (plus accept state).

pat.charAt(j)	0	1	2	3	4	5
A	A	B	A	B	A	C
dfa[][j]	A					
	B					
	C					

Constructing the DFA for KMP substring search for A B A B A C



28

Knuth-Morris-Pratt demo: DFA construction in linear time

Match transition. For each state j , $\text{dfa}[\text{pat.charAt}(j)][j] = j+1$.

↑
first j characters of pattern have already been matched

↑
now first $j+1$ characters of pattern have been matched

pat.charAt(j)	0	1	2	3	4	5
A	1		3		5	
B		2		4		
C						6

Constructing the DFA for KMP substring search for A B A B A C



29

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For state 0 and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][0] = 0$.

pat.charAt(j)	0	1	2	3	4	5
A	1		3		5	
B	0	2		4		
C	0					6

Constructing the DFA for KMP substring search for A B A B A C



30

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For state 0 and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][0] = 0$.

pat.charAt(j)	0	1	2	3	4	5
A	1		3		5	
B	0	2		4		
C	0					6

Constructing the DFA for KMP substring search for A B A B A C



30

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of empty string
↓
0

pat.charAt(j)	0	1	2	3	4	5
A	1		3		5	
B	0	2		4		
C	0					6

Constructing the DFA for KMP substring search for A B A B A C



31

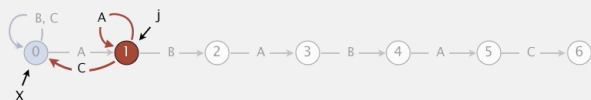
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of empty string
↓
0

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3		5	
B	0	2		4		
C	0	0				6

Constructing the DFA for KMP substring search for A B A B A C



31

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of empty string
↓
0

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3		5	
B	0	2		4		
C	0	0				6

Constructing the DFA for KMP substring search for A B A B A C



31

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
A	1	1	3		5	
B	0	2		4		
C	0	0				6

Constructing the DFA for KMP substrng search for A B A B A C



32

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
A	1	1	3		5	
B	0	2	0	4		
C	0	0	0			6

Constructing the DFA for KMP substrng search for A B A B A C



32

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
A	1	1	3		5	
B	0	2	0	4		
C	0	0	0			6

Constructing the DFA for KMP substrng search for A B A B A C



32

Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
A	1	1	3		5	
B	0	2	0	4		
C	0	0	0			6

Constructing the DFA for KMP substrng search for A B A B A C



33

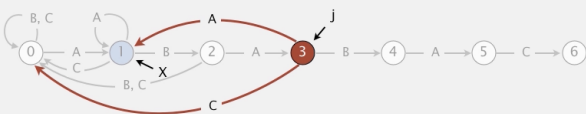
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
A	1	1	3	1	5	
B	0	2	0	4		
C	0	0	0	0		6

Constructing the DFA for KMP substrng search for A B A B A C



33

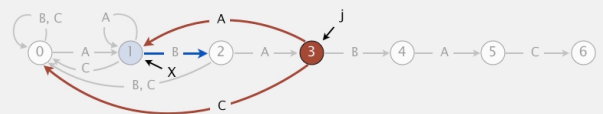
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
A	1	1	3	1	5	
B	0	2	0	4		
C	0	0	0	0		6

Constructing the DFA for KMP substrng search for A B A B A C



33

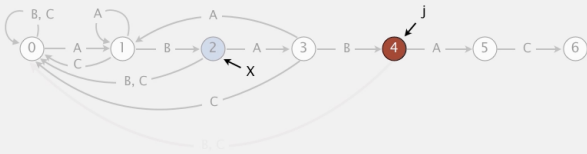
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A B
↓

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	
B	0	2	0	4		
C	0	0	0	0	6	

Constructing the DFA for KMP substrng search for A B A B A C



34

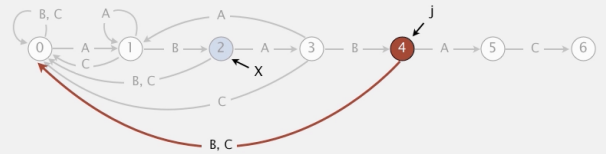
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A B
↓

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	
B	0	2	0	4	0	
C	0	0	0	0	0	6

Constructing the DFA for KMP substrng search for A B A B A C



34

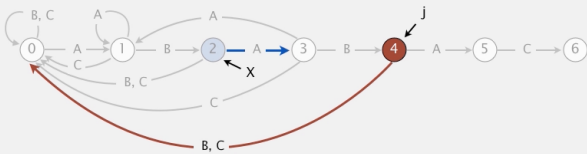
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A B
↓

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	
B	0	2	0	4	0	
C	0	0	0	0	0	6

Constructing the DFA for KMP substrng search for A B A B A C



34

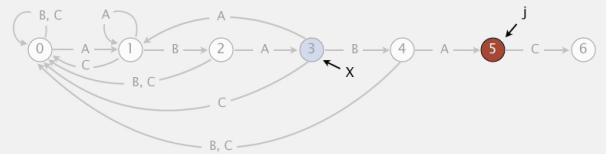
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A B A
↓

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	
B	0	2	0	4	0	
C	0	0	0	0	0	6

Constructing the DFA for KMP substrng search for A B A B A C



35

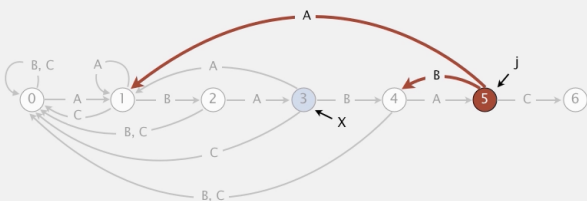
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A B A
↓

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6

Constructing the DFA for KMP substrng search for A B A B A C



35

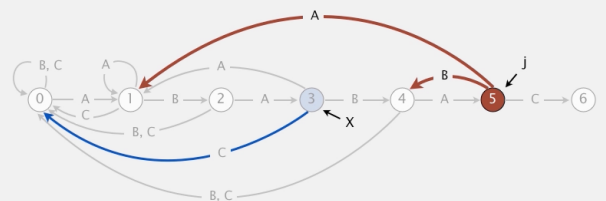
Knuth-Morris-Pratt demo: DFA construction in linear time

Mismatch transition. For each state j and char $c \neq \text{pat.charAt}(j)$, set $\text{dfa}[c][j] = \text{dfa}[c][X]$; then update $X = \text{dfa}[\text{pat.charAt}(j)][X]$.

X = simulation of B A B A
↓

pat.charAt(j)	0	1	2	3	4	5
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6

Constructing the DFA for KMP substrng search for A B A B A C

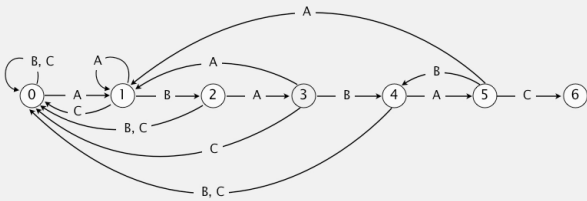


35

Knuth-Morris-Pratt demo: DFA construction in linear time

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3	1	5	1
	B	0	2	0	4	0
	C	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



37

Construção da DFA

Trecho de código do KMP que constrói a DFA.

```
dfa[pat.charAt(0)][0] = 1;
for (int j = 1, X = 0; j < m; j++) {
    for (int c = 0; c < R; c++)
        // copie casos de conflito
        dfa[c][j] = dfa[c][X];

    // defina casos de casamento
    dfa[pat.charAt(j)][j] = j+1;

    // atualize estado de reinício
    X = dfa[pat.charAt(j)][X];
}
```

Navigation icons

Construção da DFA: programação dinâmica

$dfa[c][j] = \text{maior } k \text{ tal que}$
 $pat[0..k-1] = pat[j-k+1..j-1]+c$

Para $j = 0$:

$dfa[c][0] = 1$, se $pat[0] = c$
 0 , se $pat[0] \neq c$

Para $j > 0$:

$dfa[c][j] = dfa[c][j-1]+1$, se $pat[j] = c$
 $dfa[c][X]$, se $pat[j] \neq c$,
 onde X é o maior valor tal que
 $pat[0..X]=pat[..j-1]+c$.

Navigation icons

KMP: construtor

```
public KMP(String pat) {
    this.pat = pat;
    int m = pat.length();
    dfa = new int[R][m];
    dfa[pat.charAt(0)][0] = 1;
    for (int j = 1, X = 0; j < m; j++){
        // calcule dfa[][j]
        for (int c = 0; c < R; c++){
            dfa[c][j] = dfa[c][X];
            dfa[pat.charAt(j)][j] = j+1;
            X = dfa[pat.charAt(j)][X];
        }
    }
}
```

Navigation icons

Classe KMP: esqueleto

```
public class KMP {
    private final int R = 256;
    private String pat;
    // dfa[][] representa o autômato
    private int[][] dfa;
    public KMP(String pat) { ... }
    public int search(String txt) {...}
}
```

Navigation icons

KMP: search()

```
public int search(String txt) {
    inti, n = txt.length();
    int j, m = pat.length();

    for (i = 0, j = 0; i < n && j < m; i++){
        j = dfa[txt.charAt(i)][j];

        if (j == m) return i - m;
        returnn;
    }
}
```

Navigation icons

Consumo de tempo

O consumo de tempo do algoritmo **KMP** é $O(m + n)$.

Proposição. O algoritmo **KMP** examina não mais que $m + n$ caracteres.

Se levarmos em conta o tamanho do **alfabeto**, R , o consumo de tempo para construir o DFA é mR .

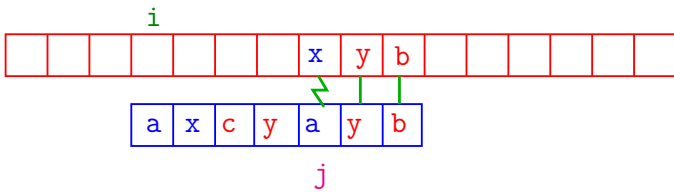
Boyer-Moore



Fonte: ADS: Boyer Moore String Search

Primeiro algoritmo de Boyer-Moore

O **primeiro algoritmo** de R.S. Boyer e J.S. Moore (1977) é baseado na seguinte heurística.



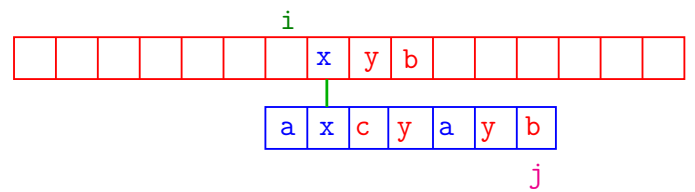
Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
```

Primeiro algoritmo de Boyer-Moore

O **primeiro algoritmo** de R.S. Boyer e J.S. Moore (1977) é baseado na seguinte heurística.



Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
```

Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
3     a n d a n d o
```

Navigation icons

Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
3     a n d a n d o
4       a n d a n d o
```

Navigation icons

Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
3     a n d a n d o
4       a n d a n d o
5         a n d a n d o
```

Navigation icons

Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
3     a n d a n d o
4       a n d a n d o
5         a n d a n d o
6           a n d a n d o
```

Navigation icons

Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
3     a n d a n d o
4       a n d a n d o
5         a n d a n d o
6           a n d a n d o
7             a n d a n d o
```

Navigation icons

Boyer-Moore

pat = a n d a n d o

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2   a n d a n d o
3     a n d a n d o
4       a n d a n d o
5         a n d a n d o
6           a n d a n d o
7             a n d a n d o
8               a n d a n d o
```

Navigation icons

Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

1 a b a b b a b a b b a



Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

1 a b a b b a b a b b a

2 a b a b b a b a b b a



Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

1 a b a b b a b a b b a

2 a b a b b a b a b b a

3 a b a b b a b a b b a



Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

1 a b a b b a b a b b a

2 a b a b b a b a b b a

3 a b a b b a b a b b a

4 a b a b b a b a b b a



Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

1 a b a b b a b a b b a

2 a b a b b a b a b b a

3 a b a b b a b a b b a

4 a b a b b a b a b b a

5 a b a b b a b a b b a



Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

1 a b a b b a b a b b a

2 a b a b b a b a b b a

3 a b a b b a b a b b a

4 a b a b b a b a b b a

5 a b a b b a b a b b a

6 a b a b b a b a b b a



Boyer-Moore

pat = a b a b b a b a b b a

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

a b a a b a b a b b a b a b a b b a b b a txt

```

1 a b a b b a b b a
2 a b a b b a b b a
3 a b a b b a b a b b a
4 a b a b b a b a b b a
5 a b a b b a b a b b a
6 a b a b b a b a b b a
7 a b a b b a b a b b a
8 a b a b b a b a b b a
9 a b a b b a b a b b a
10 a b a b b a b a b b a
11 a b a b b a b a b b a
12 a b a b b a b a b b a
13 a b a b b a b a b b a

```

Bad-character heuristic

Ideia (“*bad-character heuristic*”): calcular um deslocamento de modo que `txt[j]` fique emparelhado com a última ocorrência do caractere `txt[j]` em `pat`.

Suponha que o conjunto `a` que pertencem todos os elementos de `pat` e de `txt` é conhecido de antemão. Este conjunto é o **alfabeto** do problema.

Suponha que o alfabeto é o conjunto de todos os 256 caracteres.

Bad-character heuristic

Para implementar essa ideia fazemos um pré-processamento de `pat`, determinando para cada símbolo `x` do alfabeto a posição de sua última ocorrência em `pat`.

0	1	2	3	4	5	6
a	n	d	a	n	d	o

0	...	'a'	'b'	'c'	'd'	'n'	'o'	'p'	...	255
-1	...	3	-1	-1	5	4	6	-1

Classe BoyerMoore: esqueleto

```

public class BoyerMoore {
    private final int R; // tam. alfabeto
    // pulo bad-character
    private int[] right;
    // padrão como array ou String
    private char[] pattern;
    private String pat;
    public BoyerMoore(String pat) {...}
    public BoyerMoore(char[] pattern,
        int R) {...}
    public int search(String txt) {...}
    public int search(char[] text) {...}
}

```

BoyerMoore: construtor

```

public BoyerMoore(String pat) {
    this.R = 256;
    this.pat = pat;
    // última ocorrência de c em pat
    right = new int[R];
    for (int c = 0; c < R; c++)
        right[c] = -1;
    for (int j = 0; j < pat.length(); j++)
        right[pat.charAt(j)] = j;
}

```

BoyerMoore: search()

Recebe strings `pat` e `txt` com $m \geq 1$ e $n \geq 0$, e retorna o índice da primeira ocorrência de `pat` em `txt`. Se `pat` não ocorre em `txt`, retorna `n`.

```

public int search(String txt) {
    int n = txt.length();
    int m = pat.length();
    int skip;

```

BoyerMoore: search()

```

for (int i = 0; i <= n-m; i += skip) {
    skip = 0;
    for (int j = m-1; j >= 0; j--) {
        if(pat.charAt(j)!=txt.charAt(i+j)){
            int r= right[txt.charAt(i+j)]
            skip = Math.max(1,j-r);
            break;
        }
    }
    if (skip == 0) return i; // achou
}
return n; // não achou
}

```

Melhor caso

pat = a b c d e

```

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? txt
1 a b c d e

```

Melhor caso

pat = a b c d e

```

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? txt
1 a b c d e
2     a b c d e
3         a b c d e

```

Pior caso

pat = b a a a a a a a a a

```

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
a a a a a a a a a a a a a a a a a a a a a a a txt
1 b a a a a a a a a a
2  b a a a a a a a a a
3   b a a a a a a a a a
4    b a a a a a a a a a
5     b a a a a a a a a a
6      b a a a a a a a a a
7       b a a a a a a a a a
8        b a a a a a a a a a
9         b a a a a a a a a a
10          b a a a a a a a a a
11           b a a a a a a a a a
12            b a a a a a a a a a
13             a b a a a a a a a a

```

Melhor caso

pat = a b c d e

```

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? txt
1 a b c d e
2     a b c d e

```

Melhor caso

pat = a b c d e

```

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? txt
1 a b c d e
2     a b c d e
3         a b c d e
4             a b c d e

```

Melhor caso

pat = a b c d e

```

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? x ? ? ? ? txt
1 a b c d e
2     a b c d e
3       a b c d e
4         a b c d e
5           a b c ...
  
```

Navigation icons

Bad-character heuristic: variante

O primeiro algoritmo de R.S. Boyer e J.S. Moore (1977) é baseado na seguinte heurística.

Diagram illustrating the bad-character heuristic. A text window of size k is shown with the character 'x' at position k . Below it, a pattern window of size k is shown with the characters 'x a x y a y b'. The character 'x' in the text window matches the character 'x' in the pattern window.

Navigation icons

Bad-character heuristic: variante

pat = a n d a n d o

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
  
```

Navigation icons

Conclusões

O consumo de tempo do algoritmo BoyerMoore no pior caso é $O((n - m + 1)m)$.

O consumo de tempo do algoritmo BoyerMoore no melhor caso é $O(n/m)$.

Isto significa que no pior caso o consumo de tempo é essencialmente proporcional a mn e no melhor caso o algoritmo é **sublinear**.

Navigation icons

Bad-character heuristic: variante

O primeiro algoritmo de R.S. Boyer e J.S. Moore (1977) é baseado na seguinte heurística.

Diagram illustrating the bad-character heuristic. A text window of size k is shown with the character 'x' at position k . Below it, a pattern window of size k is shown with the characters 'x a x y a y b'. The character 'x' in the text window matches the character 'x' in the pattern window.

Navigation icons

Bad-character heuristic: variante

pat = a n d a n d o

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2     a n d a n d o
  
```

Navigation icons

Bad-character heuristic: variante

```

pat = a n d a n d o
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2     a n d a n d o
3     a n d a n d o

```



Bad-character heuristic: variante

```

pat = a n d a n d o
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2     a n d a n d o
3         a n d a n d o
4             a n d a n d o

```



Bad-character heuristic: variante

```

pat = a n d a n d o
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2     a n d a n d o
3         a n d a n d o
4             a n d a n d o
5                 a n d a n d o

```



Bad-character heuristic: variante

```

pat = a n d a n d o
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31
as andorinhas andam andando alto txt
1 a n d a n d o
2     a n d a n d o
3         a n d a n d o
4             a n d a n d o
5                 a n d a n d o
6                     a n d a ...

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a
3     a b a b b a b a b b a

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a
3     a b a b b a b a b b a
4       a b a b b a b a b b a

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a
3     a b a b b a b a b b a
4       a b a b b a b a b b a
5         a b a b b a b a b b a

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a
3     a b a b b a b a b b a
4       a b a b b a b a b b a
5         a b a b b a b a b b a
6           a b a b b a b a b b a

```



Bad-character heuristic: variante

```

pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a
3     a b a b b a b a b b a
4       a b a b b a b a b b a
5         a b a b b a b a b b a
6           a b a b b a b a b b a
7             a b a b b a b a b b a

```



Bad-character heuristic: variante

```

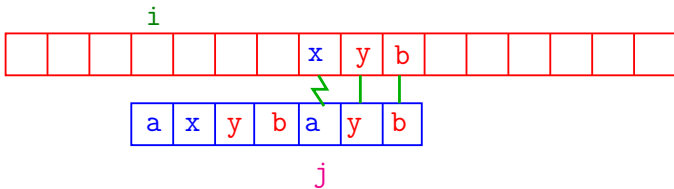
pat = a b a b b a b a b b a
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
a b a a b a b a b b a b a b a b b a b a b b a txt
1 a b a b b a b a b b a
2   a b a b b a b a b b a
3     a b a b b a b a b b a
4       a b a b b a b a b b a
5         a b a b b a b a b b a
6           a b a b b a b a b b a
7             a b a b b a b a b b a
8               a b a b b a b a b b a

```



Segundo algoritmo de Boyer-Moore

O **segundo algoritmo** de R.S. Boyer e J.S. Moore (1977) é baseado na seguinte heurística.



Navigation icons: back, forward, search, etc.

Good suffix heuristic

Não precisa conhecer o **alfabeto** explicitamente.

A implementação deve começar com um **pré-processamento** de **pat**: para cada **j** em $0, 1, \dots, m - 1$ devemos calcular o maior **k** em $0, 1, \dots, m - 2$ tal que:

- ▶ $pat[j..m-1]$ é sufixo de $pat[0..k]$ ou
- ▶ $pat[0..k]$ é sufixo de $pat[j..m-1]$

Chamemos de $bm[j]$ esse valor **k**.

Navigation icons: back, forward, search, etc.

Good suffix heuristic

O vetor $bm[]$ pode ser calculado facilmente em tempo $O(m^3)$.

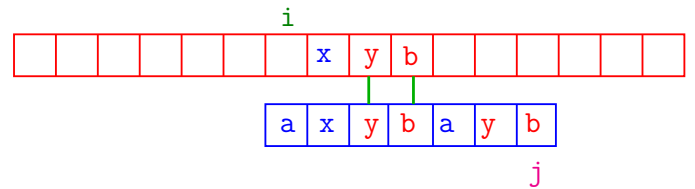
Com mais trabalho, o vetor $bm[]$ pode ser determinado em tempo $O(m)$. Veja **CLRS**, seção 34.4.

Veja também a página do professor Paulo Feofiloff: [Busca de palavras em um texto](#)

Navigation icons: back, forward, search, etc.

Segundo algoritmo de Boyer-Moore

O **segundo algoritmo** de R.S. Boyer e J.S. Moore (1977) é baseado na seguinte heurística.



Navigation icons: back, forward, search, etc.

Good suffix heuristic

Exemplo 1:

	0	1	2	3	4	5
pat	c	a	a	b	a	a
bm	-1	-1	-1	-1	2	4

Exemplo 2:

	0	1	2	3	4	5	6	7
pat	b	a	-	b	a	*	b	a
bm	1	1	1	1	1	1	4	4

Navigation icons: back, forward, search, etc.

Apêndice: Rabin-Karp



Fonte: [ADS: Boyer Moore String Search](#)

Referência: [Algoritmo de Rabin-Karp para busca de substrings \(PF\)](#)

Navigation icons: back, forward, search, etc.

Rabin-Karp

Criado por Richard M. Karp and Michael O. Rabin (1987).

O algoritmo também é conhecido como **busca por impressão digital** (*fingerprint search*).

Procura um segmento do texto que tenha o mesmo valor hash do padrão **pat**.

Usa hashing modular: módulo **Q**.

Se hash de **pat** é diferente do hash de todos os segmentos do texto então o padrão **não ocorre no texto**. A recíproca não vale: pode haver **colisão**.

Exemplo 1

		pat.charAt(j)																			
j		0	1	2	3	4															
		2	6	5	3	5	% 997 = 613														
		txt.charAt(i)																			
i		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15				
		3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3				
0		3	1	4	1	5	% 997 = 508														
1			1	4	1	5	% 997 = 201														
2				4	1	5	% 997 = 715														
3					1	5	% 997 = 971														
4						5	% 997 = 442														
5							9	2	6	5	3	% 997 = 929									
6	←							2	6	5	3	5	% 997 = 613								

Basis for Rabin-Karp substring search

Exemplo 2

Procurar o padrão **12345** nos primeiros 100 mil dígitos da expansão decimal de π .

```
31415926535897932384626433832795028841971693993751058209749445923078
16406286208998628034825342117067982148086513282306647093844609550582
23172535940812848111745028410270193852110555964462294895493038196442
88109756669334461284756482337867831652712019091456485669234603486104
54326648213393607260249141273724587006606315588174881520920962829254
09171536436789259036001133053054882046652138414695194151160943305727
03657595919530921861173819326117931051185480744623799627495673518857
52724891227938183011949129833673362440656643086021394946395224737190
7021796094370277053921717629317675238467481848766940513200056812714
5263560827785713427577896091736371787214684409012249534301465495653
71050792279689258923542019956112129021960864034418159813629774771309
9605187072113498999837297804995105973173281609631859502445945534690
83026425223082533446850352619311881710100031378387528865875332083814
20617177669147303598253490428755468731159562863882353787593751957781
85778053217122680661300192787661119590921642019893809525720106548586
32788659361533818279682303019520353018529689957736225994138912497217
75283479131515574857242454150695950829533116861727855889075098381754
63746493931925506040092770167113900984882401285836160356370766010471
01819429555961989467678374494482553797747268471040475346462080466842
59069491293313677028989152104752162056966024058038150193511253382430
03558764024749647326391419927260426992279678235478163600934172164121
9924586315030286182974557067498385054945885869269956909272107975093
02955321165344987202755960236480665499119881834797753566369807426542
5278625518184175746728909777279380008164706001614524919217321721477
23501414419735885481613611573525521334757418494684385233239073941433
3454776241686251898356948562099219222184272550254256887671790494601
6534680498862723279178608578438382796797668145410095388378636095068
00642251252051173929848960841284886269456042419652860222106611863067
```

Algoritmo de Horner

Para calcular o valor de **hash** de um string usamos o **algoritmo de Horner**:

```
private long hash(String key, int m) {
    long h = 0;
    for (int j = 0; j < m; j++)
        h = (h * R + key.charAt(j)) % Q;
    return h;
}
```

Exemplo

		pat.charAt(j)				
i		0	1	2	3	4
		2	6	5	3	5
0		2	% 997 = 2			
1		2	6	% 997 = (2*10 + 6) % 997 = 26		
2		2	2	6	5	% 997 = (26*10 + 5) % 997 = 265
3		3	2	6	5	3 % 997 = (265*10 + 3) % 997 = 659
4		4	2	6	5	3 5 % 997 = (659*10 + 5) % 997 = 613

Computing the hash value for the pattern with Horner's method

Ideia chave: hash de substrings consecutivos

Seja $t_i = \text{txt}[i]$ e $x_i =$ o inteiro $t_i t_{i-1} \dots t_{i+m-1}$

Assim,

$$x_{i-1} = t_{i-1}R^{m-1} + t_iR^{m-2} + \dots + t_{i+m-2}$$

$$x_i = \quad \quad \quad + t_iR^{m-1} + \dots + t_{i+m-2}R + t_{i+m-1}$$

Logo, $x_i = (x_{i-1} - t_{i-1}R^{m-1})R + t_{i+m-1}$

Portanto, o valor **hash**(x_i) = $x_i \% Q$ pode ser obtido a partir do valor de **hash**(x_{i-1}) = $x_{i-1} \% Q$ em **tempo constante**.

$$\text{hash}(x_i) = ((\text{hash}(x_{i-1}) - t_{i-1}R^{m-1})R + t_{i+m-1}) \% Q$$

Exemplo

i	...	2	3	4	5	6	7	...	
<i>current value</i>	1	4	1	5	9	2	6	5	↘ <i>text</i>
<i>new value</i>		4	1	5	9	2	6	5	↗
		4	1	5	9	2	<i>current value</i>		
	-	4	0	0	0				
			1	5	9	2	<i>subtract leading digit</i>		
				*	1	0	<i>multiply by radix</i>		
		1	5	9	2	0			
				+	6	<i>add new trailing digit</i>			
		1	5	9	2	6	<i>new value</i>		

Key computation in Rabin-Karp substring search

Exemplo: $Q=997$

$$\begin{aligned} (10000 + 535) \times 1000 &= (30 + 535) \times 3 \\ &= 565 \times 3 \\ &= 1695 \\ &= 698 \end{aligned}$$

$$\begin{aligned} 508 - 3 \times 10000 &= 508 - 3 \times (30) \\ &= 508 + 3 \times (-30) \\ &= 508 + 3 \times (997 - 30) \\ &= 508 + 3 \times 967 \\ &= 508 + 907 \\ &= 418 \end{aligned}$$

RabinKarp: construtor

```
public RabinKarp(String pat) {
    this.pat = pat;
    m = pat.length();
    Q = longRandomPrime();
    RM = 1;
    // calcula R^(m-1)%Q
    for (int i = 1; i <= m-1; i++)
        RM = (R * RM) % Q;
    patHash = hash(pat, m);
}
```

Implementação

Escolha Q igual a um primo grande para evitar a chance de colisão.

Evite *overflow* e números negativos:

- ▶ use um tipo-de-dados capaz de armazenar Q^2 ;
- ▶ na prática, escolha Q que caibe em um `int` ($2^{31} - 1$ é primo);
- ▶ faça as contas com `long`;
- ▶ tome o resto da divisão por Q depois de cada operação;
- ▶ some Q aos resultados intermediários quando necessário.

Classe RabinKarp: esqueleto

```
public class RabinKarp {
    private String pat;
    private long patHash; // hash do padrão
    private int m;
    private long Q;
    private int R = 256;
    private long RM;
    public RabinKarp(String pat) {...}
    private long hash(String key, int m) {}
    private int search(String txt) {...}
    public boolean check(String txt, int i)
}
```

RabinKarp: search()

```
private int search(String txt) {
    int n = txt.length();
    long txtHash = hash(txt, m);
    if (patHash == txtHash
        && check(txt, 0)) return 0;
    for (int i = 1; i <= n - m; i++) {
        txtHash = (txtHash + Q -
                    RM * txt.charAt(i-1) % Q) % Q;
        txtHash = (txtHash * R +
                    txt.charAt(i+m-1)) % Q;
        if (patHash == txtHash
            && check(txt, i)) return i;
    }
    return n; } // não achou
```

RabinKarp: check()

```
// versão Las Vegas
private boolean check(String txt, int i){
    for (int j = 0; j < m; j++){
        if (pat.charAt(j)!=txt.charAt(i+j))
            return false;
        return true;
    }
}

// versão Monte Carlo
private boolean check(String txt, int i){
    return true
}
```

Navigation icons

Monte Carlo versus Las Vegas

A versão **Monte Carlo** consome **tempo linear**, mas pode dar uma **resposta errada**, com baixíssima probabilidade.

A versão **Las Vegas** sempre dá a **resposta certa**, mas pode consumir **tempo não linear**, com baixíssima probabilidade.

Navigation icons

Implementações

Veja as implementações de busca de substrings:

- ▶ **glibc**: Implementation of strstr in glibc
- ▶ **cpython**: The stringlib Library
- ▶ **Boyer-Moore-Horspool algorithm**

Navigation icons

Exemplo

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
	3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3
0	3 % 997 = 3															
1	3	1	% 997 = (3*10 + 1) % 997 = 31													
2	3	1	4	% 997 = (31*10 + 4) % 997 = 314												
3	3	1	4	1	% 997 = (314*10 + 1) % 997 = 150											
4	3	1	4	1	5	% 997 = (150*10 + 5) % 997 = 508 ^{RM} ^R										
5		1	4	1	5	9	% 997 = ((508 + 3*(997 - 30))*10 + 9) % 997 = 201									
6			4	1	5	9	2	% 997 = ((201 + 1*(997 - 30))*10 + 2) % 997 = 715								
7				1	5	9	2	6	% 997 = ((715 + 4*(997 - 30))*10 + 6) % 997 = 971							
8					5	9	2	6	5	% 997 = ((971 + 1*(997 - 30))*10 + 5) % 997 = 442						
9						9	2	6	5	3	% 997 = ((442 + 5*(997 - 30))*10 + 3) % 997 = 929 ^{match}					
10	←	return i-M+1 = 6					2	6	5	3	5	% 997 = ((929 + 9*(997 - 30))*10 + 5) % 997 = 613				

Rabin-Karp substring search example

Navigation icons

Qual algoritmo é melhor

Força bruta é bom se o **padrão** e o **texto** não tiverem muitas auto-repetições.

KMP é rápido e tem a vantagem de nunca retroceder sobre o **texto**, o que é importante se o texto for dado como um fluxo contínuo (*streaming*).

BoyerMoore é provavelmente o mais rápido na prática.

RabinKarp é rápido mas pode dar resultados errados, com baixíssima probabilidade.

Navigation icons

Próximo passo

Que acontece se o **padrão** não é apenas uma string mas um **conjunto de strings** descrito por uma **expressão regular** como $A^*|(A^*BA^*BA^*)^*$ ou $((A^*B|AC)D)$, por exemplo?

Essa generalização do problema de busca é muito importante. A solução envolve o conceito de **autômato de estados não determinístico**.

Navigation icons