

MAC0323 Algoritmos e Estruturas de Dados II

Edição 2020 – 2

AULA 29

Mais busca de palavras

PF 13

<http://www.ime.usp.br/~pf/algoritmos/aulas/strma.html>

Busca de palavras em um texto

Problema: Dados $\text{pat}[0..m-1]$ e $\text{txt}[0..n-1]$, encontrar o número de ocorrências de pat em txt .

Exemplo: Para $n = 10$, $m = 4$, e

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

	0	1	2	3
pat	b	a	b	a

pat ocorre 2 vezes em txt .

Algoritmo força bruta: **direita** para esquerda

pat = a b a b 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 txt

0 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 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

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

Algoritmo força bruta: direita para esquerda

pat = a b a b 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 txt

0 a b a b b a b b a
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 b a

Algoritmo força bruta: direita para esquerda

pat = a b a b 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 txt

0 a b a b b a b b a
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 b a
4 a b a b 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 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 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**

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									
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	

Força bruta: direita para esquerda

Devolve a primeira de ocorrências de `pat` em `txt`.

```
int search(char *pat, char *txt) {
    int i, n = strlen(txt);
    int j, m = strlen(pat);
    for (i = 0; i <= n-m; i+=1 /* skip */) {
        for (j = m-1; j >= 0; j--)
            if(txt[i+j] != pat[j])
                break;
        if (j == -1) return i;
    }
    return n;
}
```

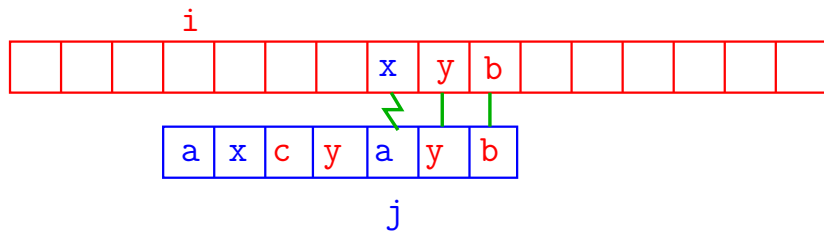
Boyer-Moore



Fonte: ADS: Boyer Moore String Search

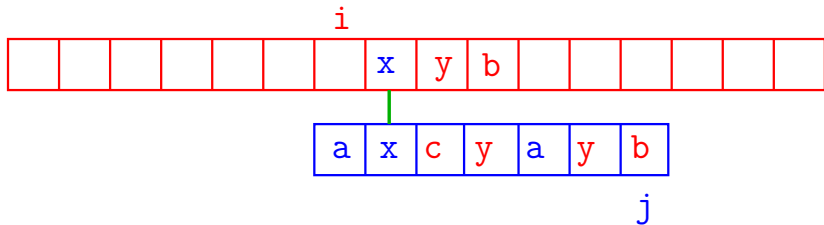
Primeiro algoritmo de Boyer-Moore

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



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

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

a s a n d o r i n h a s a n d a m a n d a n d o a l t o t x t

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

a s a n d o r i n h a s a n d a m a n d a n d o a l t o t x t

1 a n d a n d o

2 a n d a n d o

3 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

4 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
4 a n d a n d o
5 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

4 a n d a n d o

5 a n d a n d o

6 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

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

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

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 a b b a txt

1 a b a 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 a b b a txt

1 a b a 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 a b b a txt

1 a b a 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 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 a b b a txt

1 a b a 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 a b b a txt

1 a b a 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 a b b a txt

1 a b a 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

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 a b b a txt

1 a b a 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
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`.

Bad-character heuristic

Ideia (“*bad-character heuristic*”): calcular um **deslocamento** de modo que $\text{txt}[j]$ fique emparelhado com a **última ocorrência** do caractere $\text{txt}[j]$ 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
pat	a	n	d	a	n	d	o

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

Classe BoyerMoore: esqueleto

```
static int R;                /* tam. alfabeto */
static int *right;          /* pulo bad-character */
static char *pat;

void BoyerMooreInit(char *pattern) {...}
int search(char *txt) {...}
```

BoyerMoore: pré-processamento

```
void BoyerMooreInit(char *pattern) {
    int m = strlen(pattern);
    R = 256;
    pat = mallocSafe((m+1)*sizeof(char));
    strcpy(pat, pattern);

    /* última ocorrência de c em pat */
    right = mallocSafe(R*sizeof(int));
    for (int c = 0; c < R; c++)
        right[c] = -1;
    for (int j = 0; j < m; j++)
        right[pat[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`.

```
int search(char *txt) {  
    int n = strlen(txt);  
    int m = strlen(pat);  
    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[j] != txt[i+j]) {
            int r = right[txt[i+j]]
            skip = MAX(1, j-r);
            break;
        }
    }
    if (skip == 0) return i; /* achou */
}
return n; /* não achou */
}
```

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 txt

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
13 a b a 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																		

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

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 ...

Conclusões

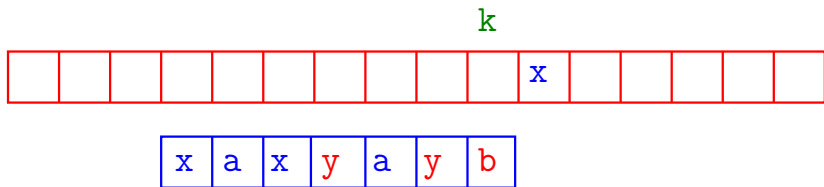
O consumo de tempo do algoritmo **BoyerMoore** no **pior caso** é $O((n - m + 1)m)$.

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**.

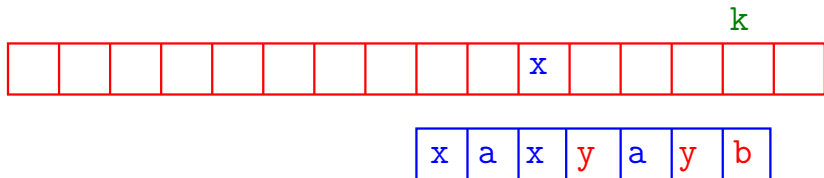
Bad-character heuristic: variante

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



Bad-character heuristic: variante

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



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

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

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
a s a n d o r i n h a s a n d a m a n d a n d o a l t o t x t

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 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 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 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 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 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 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

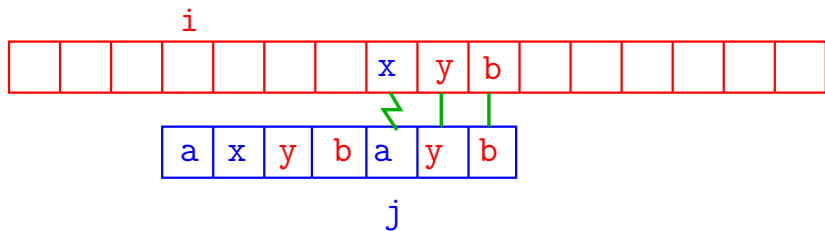
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

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 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

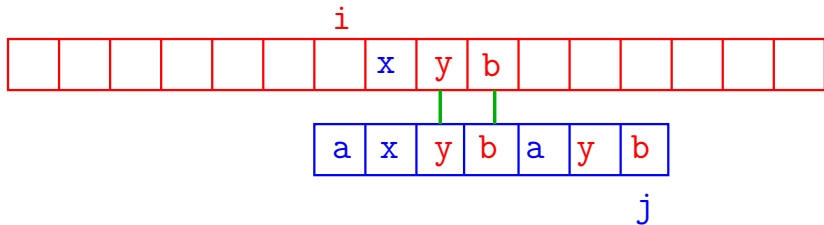
Segundo algoritmo de Boyer-Moore

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



Segundo algoritmo de Boyer-Moore

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



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 .

Good suffix heuristic

Exemplo 1:

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

Exemplo 2:

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

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](#)

Apêndice: Rabin-Karp



Fonte: ADS: Boyer Moore String Search

Referência:

Algoritmo de Rabin-Karp para busca de substrings (PF)

Rabin-Karp

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

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

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	9	% 997 = 201									
2			4	1	5	9	2	% 997 = 715								
3				1	5	9	2	6	% 997 = 971							
4					5	9	2	6	5	% 997 = 442						
5						9	2	6	5	3	% 997 = 929					
6	←						2	6	5	3	5	% 997 = 613				

match (with arrow pointing to the 613 result)

Exemplo 2

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

31415926535897932384626433832795028841971693993751058209749445923078
16406286208998628034825342117067982148086513282306647093844609550582
23172535940812848111745028410270193852110555964462294895493038196442
88109756659334461284756482337867831652712019091456485669234603486104
54326648213393607260249141273724587006606315588174881520920962829254
09171536436789259036001133053054882046652138414695194151160943305727
03657595919530921861173819326117931051185480744623799627495673518857
52724891227938183011949129833673362440656643086021394946395224737190
70217986094370277053921717629317675238467481846766940513200056812714
52635608277857713427577896091736371787214684409012249534301465495853
71050792279689258923542019956112129021960864034418159813629774771309
9605187072113499999837297804995105973173281609631859502445945534690
83026425223082533446850352619311881710100031378387528865875332083814
20617177669147303598253490428755468731159562863882353787593751957781
85778053217122680661300192787661119590921642019893809525721006548586
32788659361533818279682303019520353018529689957736225994138912497217
75283479131515574857242454150695950829533116861727855889075098381754
63746493931925506040092770167113900984882401285836160356370766010471
01819429555961989467678374494482553797747268471040475346462080466842
59069491293313677028989152104752162056966024058038150193511253382430
0355876402474964732639141199272604269922796782354781636009341172164121
992458623150302861829745557067498385054945885869269956909272107975093
02955321165344987202755960236480665499119881834797753566369807426542
5278625518184175746728909777279380008164706001614524919217321721477
23501414419735685481613611573525521334757418494684385233239073941433
34547762416862518983569485562099219222184272550254256887671790494601
65346680498862723279178608578438382796797668145410095388378636095068
00642251252051173929848960841284886269456042419652850222106611863067

Algoritmo de Horner

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

```
static long hash(char *key, int m) {  
    long h = 0;  
    for (int j = 0; j < m; j++)  
        h = (h * R + key[j]) % Q;  
    return h;  
}
```

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

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	6	5	% 997 = (26*10 + 5) % 997 = 265	
3	2	6	5	3	% 997 = (265*10 + 3) % 997 = 659
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 consecutivas

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

Assim,

$$\begin{aligned}x_{i-1} &= t_{i-1}R^{m-1} + t_iR^{m-2} + \dots + t_{i+m-2} \\x_i &= \quad \quad \quad + t_iR^{m-1} + \dots + t_{i+m-2}R + t_{i+m-1}.\end{aligned}$$

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

Portanto, o valor $\text{hash}(x_i) = x_i \% Q$ pode ser obtido a partir do valor de $\text{hash}(x_{i-1}) = x_{i-1} \% Q$ 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>new value</i>		4	1	5	9	2	6	5
		4	1	5	9	2	<i>current value</i>	
	-	4	0	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
(move right one position in the text)

Implementação

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.

Exemplo: $Q=997$

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

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

Classe RabinKarp: esqueleto

```
static char *pat;  
static long patHash; /* hash do padrão */  
static int m;  
static long Q;  
static int R = 256;  
static long RM;  
  
void RabinKarpInit(char *pattern) {...}  
static long hash(char *key, int m) {...}  
int search(char *txt) {...}  
bool check(char *txt, int i) {...}
```

RabinKarp: pré-processamento

```
void RabinKarpInit(char * pattern) {
    m = strlen(pattern);
    pat = mallocSafe((m+1)*sizeof(char));
    strcpy(pat, pattern);
    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);
}
```


RabinKarp: search()

```
int search(char *txt) {
    int n = strlen(txt);
    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[i-1] % Q) % Q;
        txtHash = (txtHash * R + txt[i+m-1]) % Q;
        if (patHash == txtHash && check(txt, i))
            return i; /* achou */
    }
    return n; /* não achou */
}
```

RabinKarp: check()

```
/* versão Las Vegas */  
static bool check(char *txt, int i) {  
    for (int j = 0; j < m; j++)  
        if (pat[j] != txt[i+j])  
            return false;  
    return true;  
}
```

```
/* versão Monte Carlo */  
static bool check(char *txt, int i) {  
    return true  
}
```

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															
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															
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

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.

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.

Implementações

Veja as implementações de busca de substrings:

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

Próximo passo

O 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**.