

MATA54 - Estruturas de Dados e Algoritmos II

Departamento de Ciência da Computação
Instituto de Matemática e Estatística
UFBA

Prof. George Lima, Ph.D.

Busca de padrões em texto

Motivação: busca de palavras

Aplicações: busca de uma sub-sequência de caracteres

- ▶ Texto: sequência de caracteres t de tamanho n
- ▶ Padrão: sequência de caracteres p de tamanho $m \leq n$
- ▶ Qual é a posição s em t onde p inicia?

Abordagens a serem estudadas

- ▶ Força-bruta
- ▶ KMP (Knuth-Morris-Pratt)
- ▶ BM (Boyer-Moore)
- ▶ RK (Rabin-Karp)

Força bruta: ilustração

Ex: $t = \text{"ababababababb"}$ e $p = \text{"ababb"}$

a	b	a	b	a	b	a	b	a	a	b	a	b	b
a	b	a	b	b									
→	a	b	a	b	b								
	→	a	b	a	b	b							
		→	a	b	a	b	b						
			→	a	b	a	b	b					
				→	a	b	a	b	b				
					→	a	b	a	b	b			
						→	a	b	a	b	b		
							→	a	b	a	b	b	
								→	a	b	a	b	b

Força bruta

Busca de $p[0, \dots, m-1]$ em $t[0, \dots, n-1]$, $m \leq n$.

```
1 int strMacher(char t[], char p[]) {
2     int i, j, m, n;
3
4     n = strlen(t);
5     m = strlen(p);
6     for (i = 0; i < n-m; i++) {
7         for (j = 0; j < m; j++) {
8             if (p[j] != t[i+j]) break;
9         }
10        if (j == m) return (i);
11    }
12    return (-1);
13 }
```

Complexidade $\Theta(nm)$ – exemplo?

É possível melhorar pré-processando o padrão p

Algoritmo KMP: ideia básica

Objetivo: evitar fazer retrocesso sobre o texto

O que se sabe quando j elementos do padrão foram comparados com o texto e o elemento $p[j + 1]$ difere do texto na posição i ?

- ▶ Que $t[i - j \dots i] = p[0 \dots j]$ (obs: sequência iniciando de 0)
 - ▶ Se há prefixo de $p[0 \dots j]$ com tamanho k igual a seu sufixo, então ao avançar p sobre t , pode-se mover $j - k$ posições para que o prefixo emparelhe às posições do texto já comparadas ao seu sufixo.
- ▶ Que pode-se evitar retrocessos sobre o texto
- ▶ Que o deslocamento do padrão sobre o texto é função do conhecimento sobre o padrão

t:	a	b	a	b	a	b	a	b	a	a	b	a	b	b
p:					a	b	a	b	a	a	b	a	b	b
					a	b	a	b	a	a	b	a	b	b
					<hr/>									
					a	b	a	b	a	a	b	a	b	b
					<hr/>									
					→									
					a	b	a	b	a	a	b	a	b	b
					<hr/>									

KMP: ilustração

Ex: $t = \text{"ababababababb"}$ e $p = \text{"ababb"}$

a	b	a	b	a	b	a	b	a	a	b	a	b	b
a	b	a	b	b									
	→	a	b	a	b	b							
		→	a	b	a	b	b						
			→	a	b	a	b	b					
				→	a	b	a	b	b				
					→	a	b	a	b	b			
						→	a	b	a	b	b		
							→	a	b	a	b	b	

KMP: pré-processamento do padrão

Objetivo

Para cada posição j de p , encontrar o tamanho do maior sufixo que é igual a um prefixo de $p[0 \dots j-1]$

$$\pi[j] = \begin{cases} \text{máximo } k > 0 & \text{tal que } p[j-k \dots j-1] = p[0 \dots k-1] \\ 0 & \text{se tal } k \text{ não existe e } j > 0 \\ -1 & \text{se } j = 0 \end{cases}$$

Ex: $p = \text{"ababbababa"}$

p	a	b	a	b	b	a	b	a	b	a
j	0	1	2	3	4	a	b	a	b	a
π	-1	0	0	1	2	0	1	2	3	4

KMP: pré-processamento do padrão

```
1 void getPi(char p[], int pi[]) {  
2  
3     int i, j, m;  
4  
5     m = strlen(p);  
6     if (m > 0) pi[0] = -1;  
7     if (m > 1) pi[1] = 0;  
8  
9     for(i = 2; i < m; i++) { // m - 2 iteracoes  
10         j = pi[i-1];  
11         while (j >= 0 && p[i-1] != p[j]) // < m iteracoes  
12             j = pi[j];  
13  
14         pi[i] = j+1;  
15  
16     }  
17     return;  
18 }
```

Complexidade: $\Theta(m)$

KMP: busca

```
1 int kmpMatcher(char t[], char p[], int pi[]) {
2
3     int i = 0, j = 0, m, n;
4
5     m = strlen(p);
6     n = strlen(t);
7     while (j < m && i < n) { // no max. 2n iteracoes
8         if (t[i] == p[j]) {
9             i++; j++;
10        } else {
11            j = pi[j];
12            if (j <= 0) {
13                j = 0; i++;
14            }
15        }
16    }
17    if (j == m) return (i - m);
18    else return (-1);
19 }
```

KMP: Complexidade

Pior caso

- ▶ Para cada bloco de m caracteres do texto, há não mais que $2m$ comparações
- ▶ Menos para o último bloco, para o qual há até m comparações
 - ▶ Ex.: $t: "aa \cdots a"$ e $p: "aa \cdots b"$. A cada comparação com $p[m-1] = b$, avança-se o padrão sobre o texto de apenas uma posição.

Então, o número total de comparações é limitado a

$$\left\lfloor \frac{n-m}{m} \right\rfloor 2m + m \leq 2(n-m) + m < 2n$$

- ▶ O pré-processamento do padrão requer $\Theta(m)$ operações
- ▶ Complexidade do pré-processamento e busca é $\Theta(n + m)$

Algoritmo Boyer-Moore: ideia básica

Objetivo: minimizar retrocessos sobre o texto e maximizar deslocamentos do padrão

- ▶ Comparações do padrão com o texto da direita para a esquerda
- ▶ Deslocamentos para a direita

Se $p[m - j - 1] \neq t[i - j]$? (deslocamento do caractere ruim)

- ▶ Pode-se avançar $m - k$ posições sobre o texto; k ou é a posição de $t[i]$ mais à direita em p ou é m se $t[i] \neq p$

t:	a	b	a	b	a	b	a	b	a	a	b	a	b	b
p:	a	b	a	a										
		→	a	b	a	a								
				→	a	b	a	a						
						→	a	b	a	a				

Algoritmo Boyer-Moore: ideia básica

Objetivo: minimizar retrocessos sobre o texto e maximizar deslocamentos do padrão

- ▶ Comparações do padrão com o texto da direita para a esquerda
- ▶ Deslocamentos para a direita

Se $p[j] \neq t[i]$ e $p[m - j + 1 \dots m - 1] = t[i + 1 \dots i + m - j - 1]$?
(deslocamento do bom sufixo)

- ▶ Seja $s = p[m - j + 1 \dots m - 1]$. Pode-se avançar o p sobre t de forma a alinhar alguma ocorrência mais à esquerda de s com $t[i + 1 \dots i + m - j - 1]$; Caso s não exista em $p[0 \dots j]$, o deslocamento do padrão é máximo.

t:	a	b	a	b	a	b	b	a	b	a	a	b	a
p:	a	a	b	a	b								
		<hr/>											
				→	a	a	b	a	b	a			
								→	a	a	b	a	a
													b

Pré-processamento: Deslocamento do caractere ruim

```
1
2 #define ASIZE 256 // para alfabeto = asc II
3
4 void badC(char p[], int badCShift[]) {
5
6     int i, m;
7
8     m = strlen(p);
9     for (i = 0; i < ASIZE; i++)
10         badCShift[i] = m;
11     for (i = 0; i < m - 1; i++)
12         badCShift[(int) p[i]] = m - i - 1;
13 }
```

- ▶ Complexidade: $\Theta(m + |\Sigma|)$; $|\Sigma|$: tamanho do alfabeto
- ▶ O algoritmo pode ser melhorado para $\Theta(m)$ com uso de hashing

Boyer-Moore: busca

Usando apenas o deslocamento do caractere ruim:

Boyer-Moore-Horspool

```
1 int bmMatcher(char t[], char p[]) {
2
3     int i, j, m, n, badCShift[ASIZE];
4
5     m = strlen(p);
6     n = strlen(t);
7     badC(p, badCShift);
8
9     i = m - 1;
10    while(i < n) {
11        j = 0;
12        while (j < m && t[i-j] == p[m-j-1]) j++;
13        if (j == m) return (i - m + 1);
14        i = i + badCShift[(int) t[i]];
15    }
16    return (-1);
17 }
```

Boyer-Moore (busca): Complexidade

Pior caso

Há retorciosos, o que leva a $\Theta(nm)$ comparações: igual a força-bruta Ex.: $t: "aa \cdots a"$ e $p: "ba \cdots a"$.

Melhor caso: considerando padrão não encontrado

Avanços de padrão sobre o texto ocorre com deslocamento máximo, levando a $\Theta(\frac{n}{m})$ comparações

Ex.: $t: "aa \cdots a"$ e $p: "bb \cdots b"$.

Caso médio*: depende do tamanho do alfabeto Σ

Sob a hipótese de uniformidade da ocorrência dos símbolos no texto e no padrão, o número médio de comparações é

$$O\left(\frac{n}{|\Sigma|}\right)$$

* R. A. Baeza-Yates, M. Regnier. "Average running time of the Boyer-Moore-Horspool algorithm". Theoretical Computer Science 92(992):19-31, 1992.

Algoritmo Rabin-Karp: ideia básica

Objetivo: comparar padrão com o texto apenas quando necessário

- Uso de função uma hashing $h(.)$ sobre o padrão e sobre cada bloco de m caracteres do texto

Se $h(p) = h(t[i \dots i + m - 1])$, verifica-se se $t[i \dots i + m - 1]$ é igual ao padrão. Caso contrário, avança-se uma posição no texto.

	$h(p) \neq h(t)$													
t:	a	b	a	b	a	b	a	b	a	a	b	a	b	b
p:	a	b	a	a										
→	a	b	a	a										
→		a	b	a	a									
→			a	b	a	a								
→				a	b	a	a							
→					a	b	a	a						
→						a	b	a	a					
→							a	b	a	a				

Construção da função hashing

Símbolos do alfabeto Σ são codificados numericamente

- ▶ Existem $|\Sigma|$ valores numéricos no alfabeto (ex.: ASCII)
- ▶ Padrão é então constituído de m dígitos numa base \mathbf{b} ,
 $p = d_0 d_{m-2} \cdots d_{m-1}$. Exemplo: $\Sigma = \{0, 1, \dots, 9\}$, $\mathbf{b} = 10$
- ▶ Com q escolhido de forma apropriada, durante o pré-processamento do padrão,

$$h(p) = \overbrace{\left(\sum_{i=0}^{m-1} p[i] \times \mathbf{b}^{m-i-1} \right)}^{\text{pode ser muito grande}} \bmod q$$
$$h(t[0 \cdots m-1]) = \overbrace{\left(\sum_{i=0}^{m-1} t[i] \times \mathbf{b}^{m-i-1} \right)}^{\text{pode ser muito grande}} \bmod q$$

Controlando valores numéricos

Observações

- ▶ Módulo da soma é igual ao módulo da soma dos módulos:

$$h(x) = \left(\sum_{i=0}^{m-1} x[i] \times \mathbf{b}^{m-i-1} \right) \bmod q = \left(\sum_{i=0}^{m-1} (x[i] \times \mathbf{b}^{m-i-1} \bmod q) \right) \bmod q$$

- ▶ q deve ser tal que o valor $q \times \mathbf{b}$ caiba numa palavra de computador (ex.: 64 bits)
- ▶ q deve ser preferencialmente primo para diminuir o número de colisões
- ▶ Valor de $h(x)$ é computado em $\Theta(m)$

Construindo a solução de Rabin-Karp

Primeira tentativa:

1. $v_p \leftarrow h(p)$ $\Theta(m)$
2. $v_t \leftarrow h(t[0, \dots, m-1])$ $\Theta(m)$
3. $s \leftarrow 0$
4. Enquanto $s \leq m - n$ $\Theta(n - m)$
 - 4.1 Se $v_p = v_t$
 - 4.1.1 Se $p = t[s, \dots, s + m - 1]$ $\Theta(m)$ somente se $v_p = v_t$
Retorna s
 - 4.2 $s \leftarrow s + 1$
 - 4.3 $v_t \leftarrow h(t[s, \dots, s + m - 1])$ $\Theta(m)$
5. Retorna -1

Pior caso: $\Theta(m(m - n))$ independentemente de v_p e v_t !!!

Solução: calcular v_t no passo 4.3 em $\Theta(1)$!!!

Construindo a solução de Rabin-Karp

Melhorando os cálculos para $h(t[s, \dots, s + m - 1])$

1. $v_p \leftarrow h(p)$ $\Theta(m)$
2. $v_t \leftarrow h(t[0, \dots, m - 1])$ $\Theta(m)$
3. $w \leftarrow b^{m-1} \bmod q$
4. $s \leftarrow 0$
5. Enquanto $s \leq m - n$ $\Theta(n - m)$
 - 5.1 Se $v_p = v_t$
 - 5.1.1 Se $p = t[s, \dots, s + m - 1]$
Retorna s
 - 5.2 $s \leftarrow s + 1$
 - 5.3 Se $s \leq m - n$
 - 5.3.1 $v_t \leftarrow (v_t - t[s - 1] \times w) \times b + t[s + m - 1]) \bmod q$ $\Theta(1)$
6. Retorna -1

Passo 5.3.1 calcula v_t usando seu valor da iteração anterior

Rabin-Karp

```
1 int rkMatcher(char t[], char p[]) {
2     long m, n, w = 1, vp = 0, vt = 0, s = 0;
3     const long q = 7919; // q primo
4
5     m = strlen(p); n = strlen(t);
6     for(int i = 1; i < m; i++) // equiv B^{m-1}
7         w = (w * B) % q;
8     for(int i = 0; i < m; i++) { // hash inicial
9         vp = (vp * B + p[i]) % q;
10        vt = (vt * B + t[i]) % q;
11    }
12    while(s <= n - m) {
13        if (vp == vt && strncmp(p, t+s, m) == 0)
14            return s;
15        s++;
16        if (s <= n - m) {
17            vt = (B*(vt - t[s-1] * w) + t[s + m - 1])%q;
18            if (vt < 0) vt = vt + q;
19        }
20    }
21    return -1;
22 }
```

Rabin-Karp: Complexidade

- ▶ Apenas pré-processamento: $\Theta(m)$
- ▶ Pior caso: **colisão com probabilidade 1**. Em todas as iterações, $h(p) = h(t[s, \dots, s + m - 1])$:
 $\Theta(m) + \Theta((n - m + 1)m) = O(nm)$
- ▶ Melhor caso: **colisão com probabilidade 0**. Com padrão não encontrado, em todas as iterações
 $h(p) \neq h(t[s, \dots, s + m - 1])$: $\Theta(m + n - m + 1) = \Theta(n)$
- ▶ Caso médio*: **Considere que $\Pr\{\text{Colisao}\} \approx \frac{1}{m}$** . Para cada valor de $s = 0, 1, \dots, m - n$, a verificação custa $\Theta(m)$, então, em média,

$$\Theta(m) + \sum_{s=0}^{n-m} \frac{1}{m} \Theta(m) = \Theta(m) + \Theta(n - m + 1) = \Theta(m + n)$$

Resumo

- ▶ Casamento de padrões em texto, um problema recorrente em várias aplicações
- ▶ Força-bruta: solução simples, mas com custo maior que o necessário
- ▶ KMP: evita retrocessos no texto com base no conhecimento sobre o padrão
- ▶ Boyer-Moore: Considera alfabeto e compara o padrão da direita para esquerda. Pode haver retrocessos sobre o texto, mas com baixa probabilidade para alfabetos grandes
- ▶ Rabin-Karp: Uso de hashing. Desempenho depende da função hashing escolhida. Pode ser estendido para casamento de padrões em duas dimensões.