# Pesquisa e Publicação de Informação Indexação e Procura (cont.)

#### Nuno D. Mendes

Licenciatura em Sistemas e Tecnologias de Informação

18 Mai 2012 ISEGI – UNL



## Motivação

- ► A procura de um termo no texto revela-se necessária quando
  - ▶ Não há um índice construído na colecção de documentos
  - Existe um índice, mas a sua granularidade não permite determinar a posição exacta das ocorrências
  - ► Existe um índice, mas pretende-se fazer a procura de um padrão (termo com erros, ou expressão regular)



## Motivação

- ► A procura de um termo no texto revela-se necessária quando
  - ► Não há um índice construído na colecção de documentos
  - Existe um índice, mas a sua granularidade não permite determinar a posição exacta das ocorrências
  - Existe um índice, mas pretende-se fazer a procura de um padrão (termo com erros, ou expressão regular)

### Tipos de procura

- Procura sequencial
  - ► Força bruta
  - ► Algoritmo Knuth-Morris-Pratt
  - ► Algoritmo Boyer-Moore



#### Motivação

- ► A procura de um termo no texto revela-se necessária quando
  - ► Não há um índice construído na colecção de documentos
  - Existe um índice, mas a sua granularidade não permite determinar a posição exacta das ocorrências
  - ► Existe um índice, mas pretende-se fazer a procura de um padrão (termo com erros, ou expressão regular)

#### Tipos de procura

- Procura sequencial
  - ► Força bruta
  - ► Algoritmo Knuth-Morris-Pratt
  - ► Algoritmo Boyer-Moore
- 2 Procura aproximada (emparelhamento de padrões)
  - ▶ Procura com erros
  - ► Procura de expressões regulares



Formalização do problema

## Definição

O problema da procura de um padrão  ${\cal P}$  num texto  ${\cal T}$  pode ser formalizado da seguinte forma:

Seja  $\Sigma$  um alfabeto de todos os caracteres que ocorrem em T. Procurar o padrão P no texto  $T \in \Sigma^*$  consiste em determinar todas as ocorrências de P em T. Se a procura for exacta ou aproximada, então  $P \in \Sigma^*$ , se se tratar da procura de uma expressão regular, então o padrão P exprime um conjunto de elementos de  $\Sigma^*$ , i.e.,  $P \in 2^{\Sigma^*}$ , segundo uma expressão regular.



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo

P = abracadabra
T = abracabracadabra



Procura sequencial

## Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra

abracabracadabra



Procura sequencial

## Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra

abracabracadabra



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo

P = abracadabra

= abracabracadabra



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra
T = abracabracadabra



\_

Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra

T = abracabracadabra



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra
T = abracabracadabra



Procura seguencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo

$$P = abracadabra$$



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra

T = a<mark>b</mark>racabracadabra



Procura sequencial

## Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo

$$P = abracadabra$$

T = ab<mark>r</mark>acabracadabra



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo

 $P = \frac{ab}{ab}$ racadabra

T = abr<mark>ac</mark>abracadabra



Procura sequencial

## Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra

T = abra<mark>c</mark>abracadabra



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

#### Exemplo



Procura sequencial

#### Força Bruta

## Definição

É o algoritmo mais simples de procura. Consiste em testar o padrão P para todas as posição do texto T. Dado que existem O(n) posições no texto e que em cada posição é preciso testar a ocorrência de um padrão de tamanho O(m), este algoritmo tem a complexidade O(nm), no pior caso.

## Exemplo

P = abracadabra
T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

$$N = 00....$$

P = abracadabra

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

#### Exemplo

$$N = 000...$$

 $P = \underline{\underline{a}}\overline{b}racadabra$ 

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

#### Exemplo

$$N = 0000.....$$
 $P = \underline{a}b\overline{r}$ acadabra
 $T = abracabracadabra$ 



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

#### Exemplo

$$N = 00001....$$

 $P = \underline{\mathbf{a}} \mathbf{br} \overline{\mathbf{a}} \mathbf{cadabra}$ 

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

$$N = 000010....$$
  
 $P = abracadabra$ 

 $T = \frac{-}{\text{abracabracadabra}}$ 



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

#### Exemplo



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

$$N = 000010101...$$
 $P = \underline{a}bracad\overline{a}bra$ 
 $T = abracabracadabra$ 



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 0000101012..

= abracadabra

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 00001010123.

 $' = abracad\overline{abra}$ 

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234  $P = \underline{abra} \underline{cadabra}$  T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234
P = abracadabra
T = abracadracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234 P = abracadabraT = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234
P = abracadabra
T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234P = abracadabra

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234

<sup>D</sup> = abracadabra

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234

P = abracadabra

T = abracabracadabra



#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

#### Exemplo

N = 000010101234

P = abracadabra

T = abracabracadabra

O emparelhamento falha na posição 7 depois de exactamente 7 caracteres emparelhados. Podemos avançar 7-N[7]=6 posições. Tendo começado na posição 0, avançamos para a posição 6. (No caso da falha ocorrer na posição 0, avanamos sempre uma posição).

Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

N = 000010101234

P = abracadabra

T = abracabracadabra



Procura sequencial

#### Knuth-Morris-Pratt

## Definição

Evita a comparação de posição do texto com o padrão, usando informação sobre tentativas anteriores. Depende da construção de uma tabela N, baseada na estrutura do padrão P. O algoritmo permite uma complexidade O(n).

## Exemplo

Cada posição do texto é comparada, no máximo, 2 vezes com o padrão, logo temos uma complexidade O(2n) = O(n).



### Boyes-Moore

## Definição

Família de algoritmos que utiliza uma tabela semelhante a N, mas considera uma janela do tamanho do padrão e realiza o emparelhamento de trás para a frente. Toma o máximo de duas regras heurísticas<sup>a</sup> para decidir o melhor deslocamento no texto. 1) a **regra do bom sufixo**, que indica a penúltima ocorrência no padrão do sufixo já emparelhado; 2) a **regra do mau caracter**, que indica a primeira posição à esquerda do sufixo já emparelhado que corresponde ao caracter que causou a falha. Em média tem uma complexidade menor do que o algoritmo de Knut-Morris-Pratt (porque nem todas as posições são comparadas), mas no pior caso tem a complexidade da abordagem de força bruta.

### llustração

P = abracadabra
T = abracababracadabra
Bom sufixo abracababracadabra
Mau caracter abracababracadabra

<sup>&</sup>lt;sup>a</sup>sub-aproximações ao deslocamento óptimo que não perdem ocorrências do padrão

## Procura com Programação Dinâmica

## Definição

Procura o padrão P no texto T a menos de e erros (inserções/delições e não-emparelhamento), usando uma recorrência para preencher uma tabela, C.

$$\begin{split} &C[0,j] = 0, \quad \forall j \\ &C[i,0] = i, \quad \forall i \\ &C[i,j] = \left\{ \begin{array}{ll} C[i-1,j-1] & \text{se} & P[i] = T[j] \\ 1 + \min(C[i-1,j], C[i,j-1], C[i-1,j-1]) & \text{se} & P[i] \neq T[j] \end{array} \right. \end{split}$$

### Ilustração

## Expressão Regular

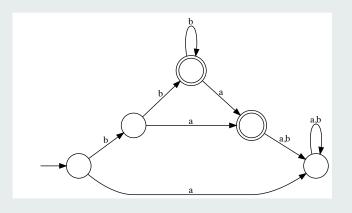
Uma expressão regular denota um conjunto de sequências de símbolos de um alfabeto  $\Sigma$  e é definida indutivamente da seguinte forma.  $\alpha$  é uma expressão regular sse:

- ▶  $\alpha = a$ , e  $a \in \Sigma$
- $\alpha=\beta\gamma$ , e  $\beta,\gamma$  são expressões regulares ( $\beta$  concatenado com  $\gamma$ )
- $\alpha = \beta^*$ , e  $\beta$  é uma expressão regular (0 ou mais ocorrências de  $\beta$ )
- $\alpha = (\beta)$ , e  $\beta$  é uma expressão regular
- $\alpha = \beta | \gamma$  e  $\beta, \gamma$  são expressões regulares ( $\beta$  ou  $\gamma$ )



## Autómato

Autómato determinista para a expressão regular definida sobre  $\Sigma = \{a,b\}$ :  $bb^*(b|b^*a)$ 



.