



Universidade de Brasília
Instituto de Ciências Exatas
Departamento de Ciência da Computação

Definições do Projeto Final

da disciplina Projeto e Análise de Algoritmos

Alexandre M. Kaihara
18/0029690

Felipe X. B. Silva
18/0016326

Gabriel R. D. Macêdo
15/0126808

Guilherme C. Suzuki
18/0032518

Italo Franklin
13/0115428

Jaqueline G. Coelho
15/0131283

Pedro L. C. Rocha
18/0054635

1 Ambiente e ferramentas

A seguir serão descritas as ferramentas necessárias para a execução do projeto, além do ambiente no qual o algoritmo irá operar.

- Linguagens: [Python 3](#), versão 3.8.10, e C++20;
- *Webcrawler*: [Scrapy](#), ferramenta *opensource* criada em Python;
- Servidor: [Apache HTTPD](#), API de servidor *Web opensource*;

Note que todas ferramentas citadas podem ser utilizadas: A licença do Apache HTTPD é do tipo GPLv3 e a do Scrapy é do tipo BSD.

2 Algoritmos utilizados

O algoritmo, em sua totalidade, consiste na junção de 3 etapas, executadas na respectiva ordem: indexação, busca e ranqueamento. Por isso, serão descritos os algoritmos a serem utilizados em cada etapa e, também, serão analisadas as complexidades de tempo de cada um deles.

2.1 Indexação

A abordagem usada é baseada na [indexação inversa](#). Como será feita a busca de palavras em documentos, então os documentos serão indexados por suas palavras, retornando um subconjunto de documentos (com possíveis intersecções entre esses subconjuntos), em uma estrutura de dados de consulta rápida.

Algorithm 1 Algoritmo de indexação inversa simples

```

1:  $Docs :=$  conjunto de páginas
2: function INDEX( $Docs$ )
3:    $Index \leftarrow \emptyset$ 
4:    $id \leftarrow 0$ 
5:   for  $D \in Docs$  do
6:     for  $s \in D.palavras$  do
7:       if  $Index[s] = \emptyset$  then
8:          $Index[s].id \leftarrow id$ 
9:          $id \leftarrow id + 1$ 
10:      end if
11:       $Index[s] \leftarrow Index[s] \sqcup \{D.id\}$ 
12:    end for
13:  end for
14:  return  $Index$ 
15: end function

```

No algoritmo 1, $Index$ refere-se a uma estrutura de dados que indexa conjuntos de documentos dada uma palavra como índice, operação representada por $Index[string]$. $Index.id$ é um identificador que será usado em outra etapas. \sqcup representa a operação de união disjunta entre conjuntos. $Docs$ é uma lista contendo todos documentos extraídos pelo *WebCrawler* e cada documento D possui uma forma de acessar todas as suas palavras na lista $D.palavras$ e um id único $D.id$.

Para o algoritmo 1, o número de passos pode ser dado por:

$$C_{indexing} = \sum_{i=1}^{Docs.n} \sum_{j=1}^{Docs[i].n} C_{append} \quad (1)$$

Sendo C_{append} um custo da operação $Index[s] \leftarrow Index[s] \sqcup \{D.id\}$, no momento em que $Index$ possui tamanho n . Note que, a seguinte inequação é satisfeita:

$$\begin{aligned}
\sum_{i=1}^{Docs.n} \sum_{j=1}^{Docs[i].n} C_{append} &\leq Docs.n \cdot \max_{D \in Docs} (D.n \cdot C_{append}^{\max}) \\
&= 1 \cdot (Docs.n \cdot C_{append}^{\max} \cdot \max_{D \in Docs} D.n)
\end{aligned}$$

Por isso, temos que, para o pior caso, $C_{indexing} \in O(Docs.n \cdot C_{append}^{\max} \cdot \max D.n)$, para a constante 1. Ou seja, o limitante superior é o produto entre número de documentos

analisados, o maior número de palavras que um documento pode apresentar e o custo máximo da operação C_{append} . As etapas posteriores requerem que cada lista contida em *Index* esteja ordenada, o que pode ser feito ao usar uma árvore binária para permitir inserções de custo $\Theta(\log n)$. Dessa forma, a complexidade do algoritmo no pior caso deve ser da ordem de $O(Docs.n \cdot \log_2 N \cdot \max D.n)$, em que N é o tamanho da maior árvore contida em *Index* após o final do algoritmo.

2.2 Pré-processamento das páginas

Para o algoritmo 3, será necessário processar cada documento conforme o algoritmo 2, o qual cria uma lista de identificadores para cada palavra do documento, consultando a estrutura *Index* construída no algoritmo 1.

Algorithm 2 Algoritmo de pré-processamento das páginas

```

Docs := conjunto de páginas
1: function MAKE_ID_LIST(Docs)
2:   for  $D \in Docs$  do
3:      $D.L \leftarrow \emptyset$  ▷ Lista de ids para o documento  $d$ 
4:     for  $s \in D$  do
5:        $D.L.append(Index[s])$  ▷ Adiciona id na lista
6:     end for
7:   end for
8: end function

```

O algoritmo possui complexidade semelhante ao do algoritmo 1, já que cada palavra de cada documento é explorada da mesma maneira. A diferença é que o custo da operação *append* é constante, mas ainda é preciso considerar o custo da consulta $Index[s]$. Logo, a complexidade do pior caso é $O(Docs.n \cdot \log_2 N \cdot \max D.n)$, em que N é o tamanho da maior lista contida em *Index*.

2.3 Pesquisa

2.3.1 Palavras simples

A busca de uma palavra depende da estrutura de dados usada na etapa de indexação. Para evitar fazer pesquisas com strings, o que pode ser lento, podemos utilizar uma função de *hashing* que converta cada palavra para um número inteiro em um intervalo pré-definido da seguinte forma:

$$\begin{aligned}
 hash(s) &= (s[0] + s[1] \cdot p + s[2] \cdot p^2 + \dots + s[m-1] \cdot p^{m-1}) \mod k \\
 &= \sum_{i=0}^{m-1} s[i] \cdot p^i \mod k
 \end{aligned}$$

Em que p é um número primo próximo da quantidade de caracteres possíveis, m é o tamanho da string e k é o tamanho do intervalo. Esta abordagem, na forma mais básica, permite colisões entre palavras. Dizemos que uma colisão entre duas palavras A e B ocorre se $hash(A) = hash(B)$, e a chance disso ocorrer, se A e B são strings aleatórias, é dada por $\frac{1}{m}$. Assim, podemos diminuir essa chance simplesmente aumentando m , ou podemos ainda usar 2 funções diferentes de hash ao mesmo tempo, o que diminuiria a chance para $\frac{1}{m^2}$. O pior caso e o caso médio de busca com técnicas de *hashing* tradicionais é da ordem de $O(\log_2 n)$, em que n é o número de elementos inseridos na tabela *hash*.

Também é possível utilizar uma árvore binária balanceada, a qual obtém resultados, no pior caso, da ordem de $O(m + \log_2 n)$, em que m é o tamanho da *string* e n é a quantidade de nós da árvore.

2.4 Pesquisa de palavras compostas

A busca de palavras compostas é feita por meio da busca da primeira palavra simples que compõe a entrada, o que fornece uma lista de documentos contendo tal palavra simples. Depois, cada documento é examinado a fim de determinar se a palavra composta existe no documento, utilizando uma versão adaptada do algoritmo de Knuth–Morris–Pratt.

Algorithm 3 Algoritmo de Knuth–Morris–Pratt

```

S := Lista de ids de strings a ser buscada
Text := Lista de ids de strings
1: function KMP(S,Text)
2:   Positions  $\leftarrow \emptyset$ 
3:   matches  $\leftarrow 0$ 
4:    $\Pi \leftarrow Preprocess(S)$ 
5:   for  $i \in \{1, 2, \dots, Text.n\}$  do
6:     while  $matches > 0 \wedge S[matches + 1] \neq Text[i]$  do
7:        $matches \leftarrow \Pi[matches]$ 
8:     end while
9:     if  $S[matches + 1] = Text[i]$  then
10:       $matches \leftarrow matches + 1$ 
11:    end if
12:    if  $matches = S.n$  then
13:       $Positions \leftarrow Positions \cup \{i - S.n + 1\}$ 
14:       $matches \leftarrow \Pi[matches]$ 
15:    end if
16:  end for
17:  return Positions
18: end function

```

O algoritmo 3 opera em complexidade $\Theta(Text.n + C_{preprocess})$, em que $C_{preprocess}$ é o custo de pré-processamento de Π , visto que o algoritmo explora cada palavra contida em *Text* apenas uma vez.

Algorithm 4 Algoritmo de pré-processamento do KMP

$S :=$ Lista de ids de strings a ser processada

```
1: function PREPROCESS( $S$ )
2:    $\Pi \leftarrow \{0\}$   $\triangleright \Pi[1] = 0$ 
3:    $length \leftarrow 0$ 
4:   for  $matches \in \{2, 3, \dots, S.n\}$  do
5:     while  $length > 0 \wedge S[length + 1] \neq S[matches]$  do
6:        $length \leftarrow \Pi[length]$ 
7:     end while
8:     if  $S[length + 1] = S[matches]$  then
9:        $length \leftarrow length + 1$ 
10:    end if
11:     $\Pi[matches] = length$ 
12:  end for
13: end function
```

A complexidade dessa etapa (algoritmo 4) é da ordem de $\Theta(S.n)$, pois cada palavra contida em S é processada uma única vez em tempo constante. Portanto, a complexidade total do algoritmo (3,4) é da ordem de $\Theta(Text.n + S.n)$.

Note que, ao invés de procurar por encaixes a nível de caracteres, são procurados encaixes a nível de palavras (identificadas por números inteiros). Por isso, é necessário um modo de converter cada palavra para seu id correspondente, o que pode ser feito ao consultar a estrutura *Index*. Nesse caso, a complexidade do pior caso ao levar em conta o cálculo de cada id é $O(S.n \cdot C_{Index} + Text.n + S.n)$, em que C_{Index} é o custo da consulta do id de uma palavra à estrutura *Index*.

2.5 Operações entre conjuntos

Para pesquisas que usam expressões que descrevem um conjunto usando os operadores \cup (União), \cap (Intersecção) e \setminus (Diferença). Cada uma delas pode ser realizada conforme os seguintes pseudocódigos.

2.5.1 União

Para suportar a operação *OR* entre dois termos de pesquisa, é suficiente fazer a união dos conjuntos de documentos de cada termo. Uma forma de fazer essa união está descrita no algoritmo 5. O algoritmo pressupõe que os conjuntos iniciais, A e B , estão ordenados, e adiciona sempre o menor dos elementos restantes para o conjunto de retorno, C . Quando os conjuntos A e B possuem um elemento em comum, apenas um deles é adicionado. Como o algoritmo consiste, essencialmente, na iteração nos conjuntos de entrada, sua complexidade é sempre $\Theta(A.n + B.n)$, onde $A.n$ é o tamanho do conjunto A e $B.n$ é o tamanho do conjunto B .

Algorithm 5 Algoritmo de União entre conjuntos

Require: $A[i] < A[i + 1], \forall i \in \{1, 2, \dots, A.n\}$

Require: $B[j] < B[j + 1], \forall j \in \{1, 2, \dots, B.n\}$

Ensure: $C[k] < C[k + 1], \forall k \in \{1, 2, \dots, C.n\}$

A,B := Lista de elementos

```
1: function UNION(A,B)
2:    $C \leftarrow \emptyset$ 
3:    $i, j \leftarrow 1$ 
4:   while  $(i \leq A.n) \wedge (j \leq B.n)$  do
5:     if  $A[i] = B[j]$  then
6:        $C \leftarrow C \cup \{A[i]\}$ 
7:        $i \leftarrow i + 1$ 
8:        $j \leftarrow j + 1$ 
9:     else if  $A[i] < B[j]$  then
10:       $C \leftarrow C \cup \{A[i]\}$ 
11:       $i \leftarrow i + 1$ 
12:     else
13:       $C \leftarrow C \cup \{B[j]\}$ 
14:       $j \leftarrow j + 1$ 
15:     end if
16:   end while
17:   while  $i \leq A.n$  do
18:      $C \leftarrow C \cup \{A[i]\}$ 
19:      $i \leftarrow i + 1$ 
20:   end while
21:   while  $j \leq B.n$  do
22:      $C \leftarrow C \cup \{B[j]\}$ 
23:      $j \leftarrow j + 1$ 
24:   end while
25:   return C
26: end function
```

2.5.2 Intersecção

O algoritmo para suportar a operação *AND* é similar ao algoritmo da operação *OR*. No entanto, aqui fazemos a intersecção, mostrada no algoritmo 6, ao invés da união. Para isso, adicionamos um elemento para o conjunto *C* somente se ele é o menor elemento restante tanto no conjunto *A* quanto no conjunto *B*. No pior caso, é preciso passar por todos os elementos de ambos os conjuntos. Portanto, o tempo máximo de execução do algoritmo é da ordem de $\Theta(A.n + B.n)$.

Algorithm 6 Algoritmo de Intersecção entre conjuntos

Require: $A[i] < A[i + 1], \forall i \in \{1, 2, \dots, A.n\}$

Require: $B[j] < B[j + 1], \forall j \in \{1, 2, \dots, B.n\}$

Ensure: $C[k] < C[k + 1], \forall k \in \{1, 2, \dots, C.n\}$

$A, B :=$ Lista de elementos

```
1: function INTERSECTION(A,B)
2:    $C \leftarrow \emptyset$ 
3:    $i, j \leftarrow 1$ 
4:   while  $(i \leq A.n) \wedge (j \leq B.n)$  do
5:     if  $e := A[i] = B[j]$  then
6:        $C \leftarrow C \cup \{e\}$ 
7:        $i, j \leftarrow (i + 1), (j + 1)$ 
8:     else if  $A[i] < B[j]$  then
9:        $i \leftarrow i + 1$ 
10:    else  $\triangleright A[i] > B[j]$ 
11:       $j \leftarrow j + 1$ 
12:    end if
13:  end while
14:  return C
15: end function
```

2.5.3 Diferença

Para dar suporte à operação *NOT*, isto é, a exclusão de um termo dos resultados, primeiro é feita a pesquisa com os demais termos, e depois são retiradas as páginas em que o termo excluído aparece. Em outras palavras, seja *A* o conjunto de páginas resultantes da pesquisa sem levar o termo excluído em consideração, e seja *B* o conjunto de todas as páginas onde o termo excluído aparece, o resultado final da busca é dado por $A \setminus B$. Esta operação está descrita no algoritmo 7. Ele é similar aos algoritmos das operações *AND* e *OR*, mas ao contrário deles, o menor elemento restante só é adicionado ao conjunto *C* se ele é o menor elemento restante de *A*, mas não de *B*. Como os conjuntos estão em ordem crescente, é garantido que se ele não for o menor elemento restante de *B*, ele não está nesse conjunto. No pior caso, é preciso passar por todo o conjunto *A* e *B*, e assim a sua complexidade também é $\Theta(A.n + B.n)$.

Algorithm 7 Algoritmo de Diferença entre conjuntos

Require: $A[i] < A[i + 1], \forall i \in \{1, 2, \dots, A.n\}$ **Require:** $B[j] < B[j + 1], \forall j \in \{1, 2, \dots, B.n\}$ **Ensure:** $C[k] < C[k + 1], \forall k \in \{1, 2, \dots, C.n\}$

A,B := Lista de elementos

```
1: function DIFFERENCE(A,B)  $\triangleright A \setminus B$ 
2:    $C \leftarrow \emptyset$ 
3:    $i, j \leftarrow 1$ 
4:   while  $(i \leq A.n) \wedge (j \leq B.n)$  do
5:     if  $A[i] = B[j]$  then
6:        $i, j \leftarrow (i + 1), (j + 1)$ 
7:     else if  $A[i] < B[j]$  then
8:        $C \leftarrow C \cup \{A[i]\}$ 
9:        $i \leftarrow i + 1$ 
10:    else
11:       $j \leftarrow j + 1$ 
12:    end if
13:  end while
14:  while  $i \leq A.n$  do
15:     $C \leftarrow C \cup \{A[i]\}$ 
16:     $i \leftarrow i + 1$ 
17:  end while
18:  return C
19: end function
```
