

Trabalho Prático 0: Trie

Francis Carlos dos Santos

Matrícula: 2012022167

1 Introdução

Pesquisa e recuperação de palavras é uma tarefa normalmente necessária em Ciência da Computação, principalmente nos últimos anos com o aumento do número de dados que podem ser acessados por meio da internet. Para um melhor proveito de toda a informação disponível na internet várias estruturas de dados foram criadas afim de facilitar essa tarefa.

Nesse trabalho o problema a ser resolvido é, dado um dicionário, descobrir com qual frequência aparecem as palavras presentes no dicionário. Será utilizada a estrutura Trie, nome que ajuda bastante a entender seu propósito pois vem do inglês *reTRIEval*, que significa recuperação. A ideia principal dessa estrutura é evitar pesquisar a palavra inteira, ou seja implementando uma estrutura onde as palavras podem ser alcançadas de acordo com seu prefixo, o que faz desse tipo de estrutura potencialmente mais eficiente que estruturas como o *Hash* e árvores binárias [1].

2 Solução do Problema

A figura (1) mostra a estrutura de dados utilizada para armazenar as palavras.

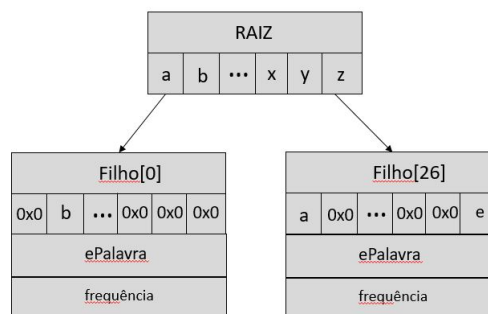


Figura 1: Estrutura de dados utilizada para armazenar o dicionário.

A figura acima ilustra como funciona a estrutura de dados proposta para resolver o problema. Cada célula receberá uma única letra. A estrutura contém um vetor apontador *filho* para uma estrutura do próprio tipo, este vetor contém vinte e sete posições. Uma variável booleana *ePalavra* que é verdadeira quando a letra inserida nesse nó representa

o fim de uma palavra do dicionário e um inteiro *freq* que é incrementada apenas quando *ePalavra = true* para contar quantas vezes a palavra aparece no texto.

Na prática cada letra do alfabeto é representada por um index, obtido através da equação mostrada abaixo.

Para auxiliar na solução do problema, foi criada também uma estrutura **Dict**, que basicamente tem a função de armazenar as palavras do dicionário e no final receber da função **encontraFreq** a frequência com que cada palavra aparece no texto. O algoritmo abaixo mostra a implementação das estruturas utilizadas.

Algoritmo 1: Estruturas utilizadas

```
1  typedef struct trie
2      bool ePalavra
3      int freq
4      struct trie *filho[TAM_ALFABETO]
5      TrieNo;

6  typedef struct
7      char string[MAX_TAM]
8      short int freq
9      Dict;
```

$$\text{index} = \text{letra} - \text{"a"} \quad (1)$$

Equação para obtenção do index de cada letra.

Como a letra 'a' é representada pelo número 97 na tabela ASCII, com essa equação é possível representar todas as letras do alfabeto como números de 0-27, logo cada índice do vetor apontador *filho* se torna uma letra quando inicializado, isso elimina a necessidade de utilizar comparação entre caracteres.

Após ler o dicionário do arquivo de entrada, todas as palavras são inseridas na árvore. No início todas as posições do vetor apontador *filho* são inicializadas com **NULL**, caso o prefixo não exista na árvore um novo nó é criado, e o ponteiro passa a apontar para o próximo filho caso contrário é feita a verificação se a letra representa uma folha e aponta para o próximo nó.

A busca é realizada sequencialmente, para cada palavra lida no texto é verificado se o caractere existe na árvore, caso exista é verificado se a letra em questão representa uma folha e o ponteiro é incrementado, caso contrário a busca termina. Os algoritmos 2 e 3 representam as duas operações descritas anteriormente.

Algoritmo 2: inserePalavra(*TrieRaiz,palavra)

```
1 // Insere letra por letra a palavra na arvore
2 // Caso o prefixo ja exista, caminha na arvore até encontrar a posição final
3 // toda letra final de palavra é marcada por ePalavra = true
4 ptr* = TrieRaiz
5 for i = 0 to k do
6     index = palavra[i] - 'a'
7     if ptr->filho[index] == NULL then // cria e inicializa celula
8         ptr->filho[index] = criaNo();
9         ptr = ptr->filho[index];
10        ptr->freq = 0;
11        if i == tamanhoPalavra-1 then Verifica se a letra e uma folha
12            ptr->ePalavra = true
13    else//caso ja exista a letra
14        if i == tamanhoPalavra-1 then
15            ptr->ePalavra = true
16            ptr = ptr->filho[index];
17
```

Por ultimo é necessário pesquisar a frequência com que cada palavra do dicionário apareceu no texto e depois imprimir. A função que faz essa busca é bem similar a anteriormente apresentada, basicamente a unica diferença é que a ultima utiliza o próprio dicionário para fazer a busca. A impressão é feita percorrendo a estrutura dicionario e imprimindo a frequência de cada palavra da estrutura, como pode ser visto no algoritmo 4.

Algoritmo 3: BuscaPalavra(*TrieRaiz,palavra)

```
1 //Pesquisa palavra do texto na árvore
2 //Caso exista a palavra incrementa a frequencia
3 //Caso a palavra nao exista na arvore a busca é
4 //interrompida após o primeiro erro.
5 for i = 0 to k do
6     index = palavra[i] - 'a'
7 if ptr->filho[index] == NULL then // se a letra nao existe sai do loop
8     break
9     else
10         if ptr->filho[index]->ePalavra == true AND i == (tamanho-1) //Verifica se
            letra é folha
11             ptr->filho[index]->freq = ptr->filho[index]->freq + 1
12             ptr = ptr->filho[index]
13         else
14             ptr = ptr->filho[index]; //Caso não seja sufixo, aponta para o proximo
15
```

Algoritmo 4: imprimeFreq(*dicionario,n)

```
1 // Imprime todas as frequencias
2 for i = 0 to m do
3     print dicionario[i].freq
4
```

3 Análise de Complexidade

Nesta seção, será apresentada a análise do custo teórico de tempo e de espaço para as principais funções do algoritmo.

3.1 Análise Teórica do Custo Assintótico de Tempo

Sejam m e n e $k = 16$ o número de palavras no dicionário, no texto e o número máximo de caracteres em cada palavra respectivamente. Para esse problema o pior caso ocorre quando todas as palavras do dicionário e do texto não possuem os prefixos iguais e as palavras tem o número máximo de caracteres permitido, o que nos dá o caso em que teríamos que percorrer toda palavra para buscar ou inserir. Porém com a implementação por indexação temos o acesso ao ramo onde se encontra o prefixo da palavra com $O(1)$ para busca ou inserção, logo o custo máximo para localizar ou inserir uma palavra seria constante. Podemos analisar o custo assintótico do programa avaliando cada um dos algoritmos apresentados juntamente com entrada e função de *freeTrie*.

- O custo assintótico do algoritmo 2 é $O(m)$, para cada palavra inserida o algoritmo faz no k iterações, portanto a inserção de cada palavra tem custo constante, para m palavras inseridas temos a complexidade assintótica $O(k*m)$ logo, $O(m)$.
- O algoritmo 3 para cada palavra no texto temos uma chamada dessa função para o número máximo de iterações igual ao tamanho da palavra. Portanto assintoticamente a função seria $O(k*n)$, como K tem valor máximo igual a dezesseis, essa função tem complexidade assintótica igual a $O(n)$.
- O algoritmo 4 é $O(m)$, sendo que sempre será impresso m frequências na saída.
- A função leitura tem é feita em duas partes, primeiramente se lê o dicionário e em seguida o texto. A complexidade nesse caso será $O(n)$, pois o texto sempre será maior que o dicionário.
- Para desalojar a Trie a função *freeTrie* percorre todos os nós da raiz e para cada nó diferente de *null* a função percorre também esse nó, em chamadas recursivas na qual a entrada é o filho do nó anterior. Nesse caso a complexidade será dada por $O(m*TAM_ALFABETO)$ pois a função percorre todas as posições de todos os nós com filhos. Como $TAM_ALFABETO$ é uma constante a complexidade assintótica dessa função é $O(m)$.

Conclui-se assim que a complexidade assintótica de tempo do algoritmo é $\max(O(m), O(n))$.

3.2 Análise Teórica do Custo Assintótico de Espaço

Temos duas estruturas de dados principais *TrieNo* e *Dict* com custos de 116 e 18 bytes respectivamente. As duas estruturas são dependentes do número de palavras do dicionário m , logo, a complexidade de espaço de cada uma será $O(m*116)$ e $O(m*18)$, a complexidade assintótica de espaço do algoritmo será dada então por, $O(116*m)$ ou $O(m)$.

4 Análise Experimental

A análise experimental da implementação é mostrada pelas figuras (2) e (3). Para realizar os experimentos foi utilizado um gerador de casos teste além dos casos teste disponibilizados via moodle. O tempo de execução foi obtido pelo terminal ao fim da execução do programa. O espaço utilizado em cada execução foi obtido por meio da utilização do Valgrind. Os testes foram realizados em uma máquina com processador Core i7-4720HQ 2.6GHz e 8GB de memória ram.

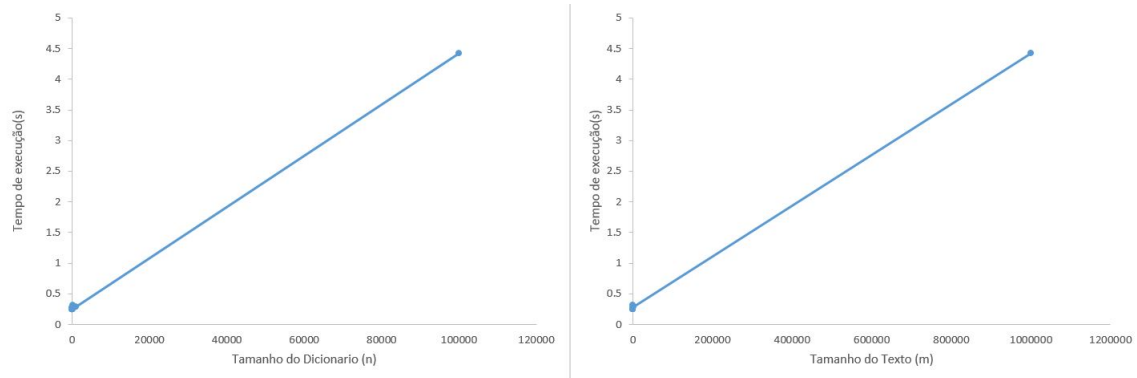


Figura 2: Grafico mostrando como varia o tempo de execução de acordo com o tamanho da entrada

A figura (2) fornece uma visão do tempo de execução do algoritmo para variações no tamanho do dicionário m e do texto n . Como esperado os graficos se comportam de forma linear conforme previsto na análise de complexidade assintótica de tempo. Ainda é possível observar que o tempo de execução para entradas grandes, realmente responde à maior entrada, m ou n .

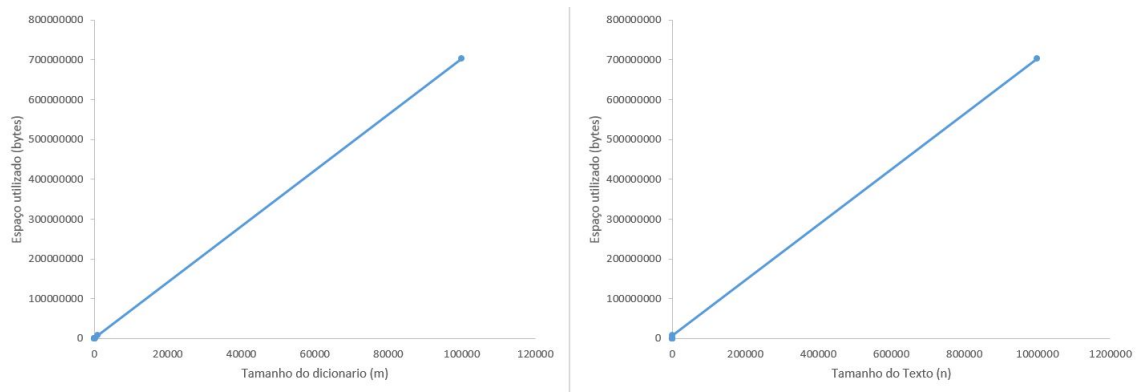


Figura 3: Estrutura de dados utilizada para armazenar o dicionário.

A figura (3) mostra a quantidade de memória, bytes, utilizada pelo programa a cada execução. É possível preceber que como previsto na análise assintótica do espaço utilizado pelo programa, os graficos só respondem á variações no tamanho do dicionário m . O que faz muito sentido já que só as palavras do dicionario são armazenadas.

5 Conclusões

Nesse trabalho, foi implementada a estrutura Trie, uma estrutura de dados mais eficiente que a tabela hash e árvores binárias, por não precisar percorrer toda a palavra. O problema a ser solucionado foi, dado um dicionário e um texto, calcular quantas vezes as palavras do dicionário aparecem no texto sendo que a frequência deveria estar presente no nó da árvore. O problema foi solucionado de forma a minimizar o tempo de processamento e o uso de memória. A análise de complexidade teórica de tempo e espaço foi comprovada. Através da análise experimental podemos comprovar que de fato o algoritmo desenvolvido tem complexidade de tempo e espaço $\max(O(m), O(n))$ e $O(m)$ respectivamente.

Referências Bibliográficas

- [1] Robert Sedgewick and Kevin Wayne. *Algorithms, 4th Edition*. Addison-Wesley, 2011.