

Relatório do Trabalho de Estrutura de Dados Avançada Dicionários

Guilherme Andrade da Silva¹

¹Universidade Federal do Ceará (UFC)
Av. José de Freitas Queiroz, 5003 – Quixadá – CE – Brazil

`guilherme.andrade@alu.ufc.br`

Abstract. *This document describes the implementation of a program that generates a word frequency list from any given book. The program utilizes dictionaries to store its elements and perform operations. For this program, four distinct dictionary structures were implemented: An AVL (Adelson-Velsky and Landis) Tree, a self-balancing binary search tree; A Red-Black Tree, a self-balancing binary search tree; A Hash Table with collision resolution via Separate Chaining; A Hash Table with collision resolution via Open Addressing. The implementation and functionality of each data structure are explained, along with their respective performance metrics.*

Resumo. *Documentação sobre a implementação de um programa que gera uma lista das frequências de palavras de um livro qualquer. O programa usa dicionários para guardar seus elementos e realizar as operações. Para esse programa foi implementado 4 dicionários: usando Árvore Binária Balanceada de Busca AVL (Adelson-Velsky e Landis), Árvore Binária Balanceada de Busca Rubro Negra, Tabela Hash Com Tratamento de Colisão Por Encadeamento Exterio e Tabela Hash Com Tratamento de COLisão Por Endereçamento Aberto. É explicado a implementação e a funcionalidade de cada estrutura, assim como suas métricas.*

1. Descrição do Projeto

O projeto consiste em implementar um programa que gera uma lista das palavras de um livro qualquer e quantas vezes ela aparece no livro, ou seja, sua frequência. O programa não diferencia maiúsculas de minúsculas, sendo assim, "House" e "house" serão interpretadas como a mesma palavra. Além disso, essa lista é gerada em ordem alfabética.

Como estrutura para guardar a frequência de cada palavra, é usado o Dicionário. O dicionário é uma estrutura que guarda um par de chave e valor e realiza as operações de inserção, busca e remoção em tempo médio de $O(1)$ ou $O(\lg n)$, a depender se foi implementado com tabelas de dispersão ou árvores balanceadas, respectivamente. Nesse trabalho, foi implementado 4 tipos de dicionários: com tabela hash encadeada, com árvore balanceada avl, com árvore balanceada rubro negra e tabela hash com endereçamento aberto

No programa de frequência em específico, as chaves dos pares serão as palavras e o valor a frequência com que aparecem. Mas os dicionários foram implementados usando tipos genéricos, podendo ser usado para qualquer outro fim.

Por fim, existem diversas aplicações da estrutura dicionário, uma dessas é essa que é tratada nessa documentação, que é gerar a frequência de certos elementos. Também pode ser usado para representar grafos, onde as chaves são nós e os valores são as listas de nós adjacentes, armazenamento de dados em cache, entre outras formas.

1.1. Árvore Binária de Busca Balanceada - AVL

A implementação da estrutura AVL foi feita usando a estrutura auxiliar Node, que guarda o par em si (chave e valor) do tipo genérico, a altura do nó e os ponteiros para seu nó esquerdo e direito.

Todos os seus algoritmos foram feitos usando recursão, sendo assim, foi feito o insert, que insere um par na árvore, caso já exista um par com essa chave, simplesmente não fazemos nada.

Já a função update, procura o nó com a chave passada no seu parâmetro e atualiza com o novo valor passado na função. Caso não exista, é retornado um erro.

As funções de busca e at (que retorna o valor relacionado a uma chave), funcionam de maneira bastante geral, a busca retorna true se a chave existir e false se não. O at retorna uma exceção caso não ache aquele par com a chave que foi passada no seu parâmetro.

O remove também busca o nó, se o achar o deleta da árvore, caso não, nada é feito. O clear remove todos os nós (com exceção da raiz) e também existe função que retorna seu tamanho e suas métricas. Há também uma função que verifica se a árvore é vazia e outra que a imprime em formato legível.

Quanto a sua criação, foi feito um construtor padrão para a classe. A implementação também tem funções privadas principalmente para balancear a árvore e a consertar depois de uma inserção e remoção e para facilitar a sobrecarga dos operadores de atribuição, igualdade e diferença. Também foi criado um construtor de cópia.

1.2. Árvore Binária de Busca Balanceada - Rubro Negra

A implementação da árvore rubro negra foi feito bastante parecido com a avl, no sentido de ser usado tipos genéricos, ponteiros, rotações e uma estrutura Node para representar o nó. A única diferença é que o nó guarda cor, em vez de altura, que pode ser preto ou branco. E também há um ponteiro há mais que é um ponteiro para o pai daquele nó. E de maneira geral, o nil usado nos algoritmos de rubro negra foi implementado de maneira que só há um nil na memória e os algoritmos foram adaptados para isso. Isso economiza memória e facilita o gerenciamento dos nós.

As funções que temos são basicamente as mesmas que temos na AVL. Foi implementado o insert, que insere um par na árvore, caso o par já exista, o valor relacionado a ele é atualizado, essa função retorna true se o valor for inserido e false se for atualizado. A função remove busca o par e o remove se ele existir, caso não, nada é feito. O contains verifica se aquele chave existe na árvore e retorna True se e somente se ele existir.

Há funções recursivas, mas foi priorizado implementações iterativas para melhor uso dos ponteiros para pai que existe nos nós. Continuando a listagem das funções, há a função at, que assim como na AVL, ela retorna o valor relacionado aquele nó ou uma exceção caso não exista na árvore. As funções empty, clear, as que retornam as métricas

e os size, assim como sobre as sobrecargas e construtor de cópia, funcionam de maneira semelhante a AVL.

Dentre as dificuldades, houve confusão no gerenciamento quanto ao nullptr e nil, sobre a diferenciação deles dois. Ainda mais no clear e destrutor que há um processo delicado de gerenciamento de memória. Assim como no remove, o qual teve que ser usado para melhor implementação o algoritmo de rb_transplant e o remove e fixed_up_remove baseados no rb_transplant conforme definido por Cormen et al. (2009).

1.3. Tabela Hash Com Tratamento de Colisão Por Encadeamento Exterior

Diferentemente das outras classes, a tabela hash com encadeamento exterior não precisou de uma estrutura auxiliar como o node. Mas precisou de um tipo genérico a mais, que foi o hash para as operações de hash codes.

Assim, esta tabela tem as mesmas funções que a AVL e Rubro-Negra, não tendo o update, pois o seu insert funciona como o da rubro negra, retornando True se tiver sido inserido e False se tiver sido atualizado. A sobrecarga de iterador foi feita de maneira que se não existe aquele par, a função já o cria e retorna o construtor padrão vazio para o objeto.

Fora as operações padrões, há as funções próprias de uma tabela hash: rehash, o qual recebe o novo do tamanho da tabela e ela é adaptada para aquele tamanho, com seus elementos realocados, essa função também é realizada toda vez que a tabela atinge seu fator de carga máximo. Resize, que recebe o número de elementos que a tabela tem que se adaptar para receber. Set_max_load_factor que muda o fator de carga máximo da tabela, fazendo uma operação de reserve, que realiza um rehash se necessário. E também uma função que retorna o fator de carga atual da tabela.

A implementação foi feita também usando tipos genéricos e funções auxiliares que ajuda no tempo de operação: hash_code e função que retorna o próximo número primo de um número. Além disso, há uma função que retorna quantos elementos existem na tabela, outra que retorna quantos buckets existem e qual o bucket em que um par está.

Como estruturas auxiliares, foi usado um vetor de listas encadeadas, para melhor desempenho da estrutura e duas operações. Além da sobrecarga de iterador, foi feita uma sobrecarga de impressão, para facilitar a manipulação e uso da estrutura fora de sua classe.

Para fim das métricas, consideramos a colisão, que na implementação é verificado se a lista em que vamos inserir o nó é empty ou não, se não for, há uma colisão e logo em seguida inserimos o nó.

1.4. Tabela Hash Com Tratamento de Colisão Por Endereçamento Aberto

A tabela hash com endereçamento aberto foi implementada com uma estrutura auxiliar Node, o qual guarda o par e seus status. A mesma estrutura contém um construtor vazio, o qual foi pensado para não ter que inicializar cada elemento da estrutura que vai armazenar a tabela como Empty, pois o construtor vazio já cria o Node com os construtores vazios dos tipos genéricos e Empty para seu status.

Há também outra estrutura auxiliar que é o vetor que armazena os nós da tabela, no qual todos os seus slots estão "ocupados", em sentido de memória. Mas só consideramos

preenchidos aqueles que são Active, Empty são os que estão vazios e nunca foram usados e Deleted os que já foram Active, mas o seu elemento foi removido.

Nas operações, considera-se essencialmente esses estados dos slots. Na inserção, ao acharmos o bucket que inserirmos aquele par, mudamos seu estado para Active, se antes não era. Na inserção usamos um algoritmo de encontrar o primeiro slot empty, caso ache um deleted, verificamos os outros para ver se o elemento já não está na tabela para o atualizarmos. Retorna true se foi inserido e false se foi atualizado.

Na remoção, consideramos os elementos ativos e ao acharmos para mudar seu status para deleted e por segurança, mudamos seu par para o construtor vazio de cada objeto, da chave e do valor. Em ambas as funções de remoção e inserção, diminuimos e incrementamos o size, respectivamente. Também há funções constantes que retorna o size, as métricas da tabela e o número de buckets da tabela.

Na busca, consideramos também apenas os elementos ativos e retornamos true se e somente se o elemento existir. As operações at e iterador funcionam semelhante-mente as da tabela com encadeamento. Também foi sobrecarregado um cout para melhor manipulação fora da classe.

Quanto ao rehash e hashcode, o rehash acontece se o número de elementos for igual ao size/capacity do vetor, pois o número de elementos é um atributo a parte do vetor. Também ocorre quando é atingido o seu fator máximo de carga, o que ocorrer primeiro. Para melhor uso e desempenho da tabela, foi usado hashing duplo, ou seja, foi implementado duas funções de hash code para que as funções e operações sejam realizadas no melhor tempo possível.

Por fim, não foi implementado uma função que devolva para onde uma slot foi alocado, pois não há sentido fazer uma função assim usando hashing duplo.

2. Métricas utilizadas

Em todas as estruturas, há a métrica de comparação de chaves onde em cada função que realize de fato uma comparação de chave, ela é incrementada.

```
1 Node *fixup_node(Node *node, Key k) {
2     int bal = balance(node);
3
4     if(bal < -1) {
5
6         //Caso 1(a)
7         this->m_counter_comparator++;
8         if(k < node->left->n_pair.first) {
9             return right_rotation(node);
10        }
11
12        //Caso 1(b)
13        this->m_counter_comparator++;
14        if(k > node->left->n_pair.first) {
15            node->left = left_rotation(node->left);
16            return right_rotation(node);
17        }
18    }
```

```

18     }
19
20     if(bal > 1){
21
22         //Caso 2(a)
23         this->m_counter_comparator++;
24         if(k > node->right->n_pair.first){
25             return left_rotation(node);
26         }
27
28         //Caso 2(b)
29         this->m_counter_comparator++;
30         if(k < node->right->n_pair.first){
31             node->right = right_rotation(node->right);
32             return left_rotation(node);
33         }
34     }
35
36     node->height = 1 + std::max(height(node->left), height(node
->right));
37     return node;
38 }

```

Listing 1. Função de correção do balanceamento (fix-up) da AVL. Aqui podemos ver a implementação da métrica de comparações de chaves, que acontece antes do if para a implementação correta

Também existe uma função em cada estrutura que retorna quantas comparações de chaves foram realizadas. Caso a estrutura seja limpa, seja executada a função clear, as métricas são zeradas, para melhor reaproveitamento da estrutura para qualquer outro fim. Ademais, fora essa métrica geral, também existem métricas específicas para cada tipo de estrutura.

2.1. Árvores balanceadas

Para as árvores balanceadas, a métrica específica escolhida foi o número de rotações que ela realiza durante suas operações, tanto na AVL, como na Rubro-Negra. Sendo assim, toda vez que a estrutura realiza uma rotação o contador é incrementado.

```

1 void leftRotate(Node* x){
2     Node* y = x->right;
3
4     x->right = y->left;
5
6     if(y->left != this->nil){
7         y->left->parent = x;
8     }
9
10    y->parent = x->parent;
11
12    if(x->parent == this->nil){
13        this->root = y;

```

```

14     }else if(x== x->parent->left){
15         x->parent->left = y;
16     }else{
17         x->parent->right = y;
18     }
19
20     y->left = x;
21     x->parent = y;
22
23     this->m_counter_rotation++;
24 }

```

Listing 2. Função de rotação à esquerda da Rubro Negra, a métrica de rotação é incrementada no final, sendo uma maneira simples de ser corretamente controlada

Também existe uma função em cada uma das árvores que retorna a métrica de rotações.

2.2. Tabelas hash

Para as tabelas, a métrica específica escolhida foi o número de colisões que acontecem durante suas operações. Para fins de implementação, foi considerado como colisão quando um par durante sua inserção é alocado para um bucket que já está ocupado, conforme Cormen et al. (2009). Sendo assim, toda vez que acontece uma colisão o contador é incrementado.

```

1 bool add(const Key& k, const Value& v){
2
3     if(load_factor() >= m_max_load_factor){
4         rehash(m_table_size * 2);
5     }
6
7     int slot = hash_code(k);
8
9     for(auto p : m_table[slot]){
10         this->m_counter_comparator++;
11         if(p.first == k){
12
13             return false;
14         }
15     }
16
17     if(!m_table[slot].empty()){
18         this->m_counter_collision++;
19     }
20
21     m_table[slot].push_back(std::make_pair(k,v));
22
23
24     m_number_of_elements++;
25

```

```
26     return true;
27
28 }
```

Listing 3. Função de adição em uma tabela hash com encadeamento aberto, a colisão acontece verificando antes de inserir o elemento se aquela lista é empty ou não, caso seja, a métrica é incrementada

Também existe uma função em cada uma das tabelas que retorna a métrica de colisões.

3. Dicionários (Maps)

Os dicionários foram implementados seguindo um escopo básico que definiu suas funções mínimas:

- **Construtor** : que constrói um dicionário vazio conforme a estrutura auxiliar que tiver sendo usada.

- **Add** : insere um valor no dicionário.

- **Remove** : remove um valor do dicionário.

- **Contains** : verifica se uma chave existe.

- **Clear** : limpa o dicionário, removendo os pares alocados.

- **Update** : atualiza um par chave-valor.

- **At** : retorna o valor associado a uma chave.

- **Size** : retorna quantos pares estão contidos no dicionário.

Também foi implementado outras funções auxiliares:

- **counter_collision/counter_rotation** : Retorna quantas rotações o dicionário realizou, se tiver usando um que tenha a árvore como estrutura base, ou retorna o número de colisões que aconteceram, caso esteja sendo usado um dicionário que utilize as tabelas hash como estrutura base.

- **counter_comparator** : Retorna quantas comparações de chaves foram realizadas durante as operações da estruturas.

3.1. Função Vectorize

Foi implementado uma função em cada tipo de dicionário que retorna um vetor com todos seus elementos que estão guardados. No caso dos dicionários que utilizam de árvores como estrutura base, o vetor já está previamente ordenado. Caso não, os elementos do vetor estão na ordem em que estão dispostos.

Essa função foi criada para facilitar a ordenação no caso dos dicionários que utilizam de tabelas e também para facilitar a geração da tabela de frequência.

3.2. Sobre carga de Indexação

Foi feito a sobre carga do operador de indexação em cada estrutura, a sobre carga funciona de maneira em que apenas em uma busca, verificamos se existe, se existir, retornamos o valor do par, caso não, o inserimos ali mesmo, onde ele deveria estar.

Isso é muito poderoso, pois diminui bastante o custo em termos práticos, o que irá agilizar a geração da tabela de frequência. Sendo mais custoso para os dicionários que utilizam de árvores na primeira inserção do elemento, por conta das rotações e colarações no caso do dicionário que usa árvore rubro negra. Mas, a busca será mais eficiente, pois nas tabelas, a busca tem bastante custo ($O(N)$) no pior caso.

4. Implementação da Frequência

Para a geração da tabela de frequência, foi criado uma classe para isso, para limpeza de código e organização, separamos a main da classe frequência e também para seguir as diretrizes da programação orientada a objeto, no caso a modularidade. Nesse contexto, a frequência é gerada através do construtor da classe frequência que utiliza de um dicionário base, o arquivo a ser lido e opcionalmente o nome do arquivo de saída.

Ademais, foi criado um atributo **ordena**, que verifica se o vetor que irá gerar a tabela de frequência em ordem alfabética precisa ser ordenado utilizando a biblioteca do `icu`, para ser corretamente ordenado casos de acentos, cê-cedilha, entre outras nuances da língua portuguesa. Pois, a ordenação por `icu` tem custo das suas operações de comparação, o que também agiliza, por exemplo, para casos em que usamos dicionário que usa tabela open hash e texto em inglês, a apenas ordenamos sem custo do `icu`, o que em termos práticos é mais rápido.

No caso dos dicionários que utilizam as árvores, não os ordenamos se o texto que foi lido está em inglês, pois a estrutura já está ordenada. No caso contrário, ordenamos seu vetor com custo do `icu`. Uma outra alternativa seria sobrecarregar o comparador dentro da implementação das árvores que o dicionário utiliza, mas foi pensado que isso teria mais custos ao longo da formação da estrutura, sendo preferido pagar o custo do `icu` apenas uma vez com a ordenação do vetor.

Tabela 1. Complexidade de Operações em Dicionários

Dicionário	Inserção		Busca		Remoção	
	Caso Médio	Pior Caso	Caso Médio	Pior Caso	Caso Médio	Pior Caso
Dicionário AVL	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$
Dicionário Rubro Negro	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)$
Dicionário Chained	$O(1)$	$O(n)$	$O(1)$	$O(n)$	$O(1)$	$O(n)$
Dicionário Open	$O(1)$	$O(n)$	$O(1)$	$O(n)$	$O(1)$	$O(n)$

5. Testes Realizados e Análises

Todos os testes foram realizados 100 vezes e feito a média para ser feita uma boa análise.

5.1. Criação da tabela de frequência

Tabela 2. Desempenho na Montagem da Tabela de Frequência Da Bíblia (sem custo do icu)

Estrutura de Dicionário	Tempo Médio (ns)	Tempo Médio (ms)	Tempo Médio (s)
Dicionário Open (End. Aberto)	800.637.928	800,64	0,80
Dicionário Rubro-Negro	909.492.913	909,49	0,91
Dicionário AVL	1.399.751.937	1399,75	1,40
Dicionário Chained (Encadeado)	1.431.239.329	1431,24	1,43

Tabela 3. Desempenho na Montagem da Tabela (Textos em Português - Dom Casmurro)

Estrutura de Dicionário	Tempo Médio (ns)	Tempo Médio (ms)	Tempo Médio (s)
Dicionário Open (End. Aberto)	121.578.867	121,58	0,12
Dicionário Rubro-Negro	161.089.753	161,09	0,16
Dicionário AVL	169.312.570	169,31	0,17
Dicionário Chained (Encadeado)	197.444.949	197,44	0,20

5.2. Busca nos dicionários

Tabela 4. Desempenho na Operação de Busca

Estrutura de Dicionário	Tempo Médio (ns)	Tempo Médio (ms)
Dicionário Rubro-Negro	8.674.028	8,67
Dicionário Chained (Encadeado)	9.191.439	9,19
Dicionário Open (End. Aberto)	9.380.098	9,38
Dicionário AVL	18.888.670	18,89

5.3. Teste de inserções

Tabela 5. Desempenho na Operação de Inserção

Estrutura de Dicionário	Tempo Médio (ns)	Tempo Médio (ms)
Dicionário Open (End. Aberto)	3.504.054	3,50
Dicionário Chained (Encadeado)	4.890.458	4,89
Dicionário Rubro-Negro	13.210.342	13,21
Dicionário AVL	30.877.215	30,88

5.4. Teste com entradas repetidas

Tabela 6. Desempenho com Palavras Repetidas (Busca + Atualização)

Estrutura de Dicionário	Tempo Médio (ns)	Tempo Médio (ms)
Dicionário Chained (Encadeado)	6.644.708	6,64
Dicionário Rubro-Negro	6.913.233	6,91
Dicionário Open (End. Aberto)	7.117.437	7,12
Dicionário AVL	9.149.442	9,15

5.5. Métricas e Suas Comparações

Tabela 7. Análise de Desempenho: Livro "Bíblia"

Estrutura (Dicionário)	Tipo	Comparações de Chaves	Colisões / Rotações
Árvore AVL	Árvore	11.955.920	9.648
Árvore Rubro-Negra (RB)	Árvore	13.801.555	8.123
Tabela Hash (Open)	Tabela Hash	809.311	10.906
Tabela Hash (Chained)	Tabela Hash	809.311	12.685

Tabela 8. Análise de Desempenho: Livro "Dom Casmurro"

Estrutura (Dicionário)	Tipo	Comparações de Chaves	Colisões / Rotações
Árvore AVL	Árvore	1.143.543	7.433
Árvore Rubro-Negra (RB)	Árvore	1.341.235	6.248
Tabela Hash (Open)	Tabela Hash	58.543	7.156
Tabela Hash (Chained)	Tabela Hash	58.543	10.113

Tabela 9. Análise de Desempenho: Livro "O Jardim Secreto"

Estrutura (Dicionário)	Tipo	Comparações de Chaves	Colisões / Rotações
Árvore AVL	Árvore	1.255.035	4.300
Árvore Rubro-Negra (RB)	Árvore	1.401.915	3.512
Tabela Hash (Open)	Tabela Hash	77.786	4.259
Tabela Hash (Chained)	Tabela Hash	77.786	5.817

6. Especificações

As estrutura e suas funções foram implementadas em C++ e testadas em um ambiente com as seguintes especificações:

- Processador AMD Ryzen 5 7520U with Radeon Graphics, com frequência máxima de: (CPU MHz máx.: 6734,7651)
- Memória RAM de 5.52 GB
- Memória SWAP de 2.00 GB
- Arquitetura x86-64
- 250,4 GB de armazenamento em SSD
- Compilador g++ versão 13.3.0
- Sistema Operacional Linux Mint versão 22.1
- Kernel Linux 6.8.0-63-generic

7. Conclusão

7.1. Das Estruturas Auxiliares

As estruturas foram implementadas conforme suas operações e para funcionalidades genéricas. Com implementações que tem o tempo médio destinado de cada estrutura, no sentido geral $O(\lg n)$ para as árvores e $O(1)$ para as tabelas hash. Foi tomado um cuidado para serem as mais genéricas possíveis de maneiras que não estão atreladas a esse projeto em específico.

7.2. Dos Dicionários E Frequência

Em relação aos dicionários, é bem claro que o dicionário que utiliza a árvore avl é muito mais custoso para as demais operações, apenas perdendo para o dicionário que usa a tabela hash encadeada, essa tendo o maior tempo médio para montar a tabela de frequência.

Ademais, observamos que as dicionários que usam as tabelas são muitos poderosos para inserção e busca, pois observamos que eles tem o menor tempo, com exceção do dicionário que usa a árvore rubro negra e na inserção, o dicionário que utiliza a árvore rubro negra levou 13 milissegundos aproximadamente, sendo que o que usou tabela encadeada não levou nem 5, sendo o que faz uso da avl o pior, com 30 milissegundos.

Em resumo, o dicionário que usa a tabela hash com encadeamento aberto é o melhor na média, pois seu tempo não é o pior e nem nas outras operações. Logo em seguida ficamos com o que usa a árvore rubro negra, por conta da sua ordenação prévia, o que usa tabela encadeada e por fim o que usa a avl, por ter ficado com o pior tempo médio nas operações.

Referências

Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). *Introduction to Algorithms*. The MIT Press, 3rd edition.