

Universidade Federal de Minas Gerais

Programação e Desenvolvimento de Software II - (PDSII)

Alunos: André Luiz Alves Costa e Bruno Henrique Evangelista Pereira

Link do repositório: https://github.com/andrelac963/TP_maquina_busca

Trabalho Prático: Máquina de Busca

Introdução

O trabalho prático teve o objetivo de desenvolver uma máquina de busca feita a partir da linguagem de programação C++, onde em um banco de dados de arquivos txt, o usuário irá dizer o que deseja consultar e o programa desenvolvido retornará em quais arquivos as palavras dadas pelo usuário aparecem.

Ao início do projeto nos reunimos para decidir a estrutura da aplicação, a modelagem das classes que seriam necessárias e como seria implementado cada funcionalidade. Para isso foi necessário fazer uma abstração das funcionalidades aplicando os conhecimentos aprendidos na disciplina de PDS 2 e isso auxiliou na organização geral. E com as ideias bem organizadas podemos partir para a implementação do projeto de fato.

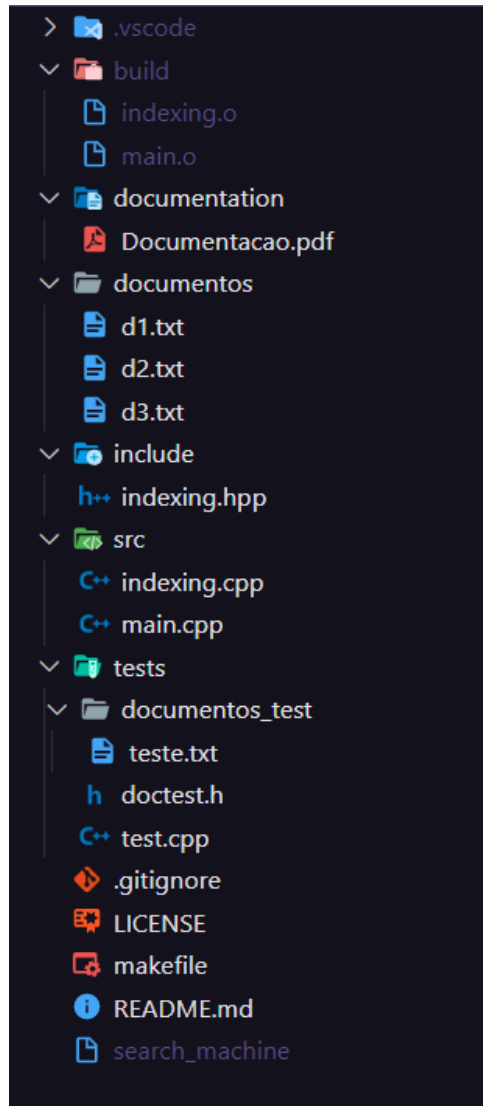
Implementação

Para versionamento do código foi utilizado o github, que funciona baseado em git, e também a utilização de pull requests para que os membros do projeto pudessem validar a implementação dos colegas de trabalho.

Para compilar o código foi utilizado o gcc versão 11.3.0 rodando no WLS com o Ubuntu 22.04.1 LTS. E para facilitar os testes foi feito um makefile com os comandos para compilar todos os arquivos do projeto.

A estrutura de pastas utilizadas foi da seguinte forma:

- build: para guardar os arquivos de compilação;
- documentos: para a base de dados lida pelo programa;
- include: para os arquivos de cabeçalho;
- src: para os arquivos de implementação.
- tests: para realização dos testes de unidade.



Decidimos que só seria necessário implementar uma classe indexing onde conteria as funcionalidades de indexação das palavras da base de dados e de recuperação delas.

```

4  #include <map>
5  #include <set>
6  #include <vector>
7  #include <string>
8  #include <iostream>
9  #include <bits/stdc++.h>
10 #include <dirent.h>
11
12 using namespace std;
13
14 class indexing
15 {
16 private:
17     map<string, set<string>> index;
18
19 public:
20     indexing();
21     ~indexing();
22
23     set<string> read_directory(const string &);
24     void read_files(const string &);
25
26     void insert(string, string);
27     map<string, set<string>> get_index();
28     string normalize(string);
29
30     vector<pair<string, int>> recovery(vector<string>);
31     void print_ordered_files(vector<string>);
32 };

```

Para guardar as informações do index decidimos por utilizar um map responsável por manter a string dada pelo base de dados e o set de strings que são os nomes dos arquivos txt onde ela ocorre. Como há necessidade de que outros arquivos acessem os métodos da classe somente o container responsável por guardar os dados fique privado. Partindo para implementação do indexing, além do construtor e destrutor de instâncias foi preciso desenvolver métodos para a leitura dos diretórios em que os .txt se encontravam e para ler os arquivos em si .

```

1  #include "indexing.hpp"
2
3  indexing::indexing()
4  {
5  }
6
7  indexing::~indexing()
8  {
9  }

```

```

11 set<string> indexing::read_directory(const string &name)
12 {
13     DIR *dirp = opendir(name.c_str());
14     if (dirp == NULL)
15     {
16         return set<string>();
17     }
18
19     struct dirent *dp;
20     set<string> files;
21
22     while ((dp = readdir(dirp)) != NULL)
23     {
24         if (dp->d_name[0] != '.')
25             files.insert(dp->d_name);
26     }
27     closedir(dirp);
28
29     return files;
30 }

```

Como o nome dos arquivos dentro do diretório ./Documentos não é fixo, foi necessário pesquisar uma forma de saber o nome de todos os arquivos dentro de uma pasta, que é o que foi implementado na imagem acima. Utilizando a biblioteca dirent.h, foi possível percorrer o nome de cada arquivo e guardá-los em um container set.

```

27 void indexing::read_files(const string &name)
28 {
29     set<string> files = read_directory(name);
30     fstream file;
31     string word, filename;
32
33     for (auto it = files.begin(); it != files.end(); it++)
34     {
35         filename = *it;
36         file.open(name + filename);
37         while (file >> word)
38         {
39             this->insert(word, filename);
40         }
41         file.close();
42     }
43 }

```

Agora utilizando a biblioteca bits/stdc++, foi possível abrir cada arquivo com o nome salvo no set e inserir cada palavra com seu respectivo nome de arquivo com a função insert.

```

50 void indexing::insert(string word, string filename)
51 {
52     string normalized_word = this->normalize(word);
53     if (normalized_word != "")
54     {
55         this->index[normalized_word].insert(filename);
56     }
57 }
58
59 map<string, set<string>> indexing::get_index()
60 {
61     return this->index;
62 }

```

A função insert, que recebe a palavra e o nome do arquivo onde ela foi encontrada, primeiro normaliza a palavra e depois a insere no index.


```

99  vector<pair<string, int>> indexing::recovery(vector<string> query)
100  {
101      set<string> normalized_query;
102      string aux;
103
104      for (int i = 0; i < query.size(); i++)
105      {
106          aux = this->normalize(query[i]);
107          if (aux != "")
108          {
109              normalized_query.insert(aux);
110          }
111      }
112
113      map<string, int> relevant_files;
114
115      for (auto it = this->index.begin(); it != this->index.end(); it++)
116      {
117          if (normalized_query.find(it->first) != normalized_query.end())
118          {
119              for (auto it2 = it->second.begin(); it2 != it->second.end(); it2++)
120              {
121                  relevant_files[*it2]++;
122              }
123          }
124      }

```

Agora para a função de recovery, responsável por retornar os arquivos com a pesquisa feita pelo usuário, precisamos primeiramente normalizar os parâmetros de pesquisa recebidos.

Feito isso, podemos então criar um map responsável por contar quantas vezes a pesquisa do usuário está presente no index. Dentro de um loop que percorre as chaves do index, o programa compara se a palavra do index está presente no container normalized_query, se sim então ele irá percorrer o container set do index, que guarda os arquivos onde a palavra ocorre, e colocar o arquivo no container relevant_files, onde cada vez que ele aparecer na pesquisa irá ser somado mais um ao seu valor.

```

126  vector<pair<string, int>> ordered_files;
127
128  for (auto it = relevant_files.begin(); it != relevant_files.end(); it++)
129  {
130      ordered_files.push_back(make_pair(it->first, it->second));
131  }
132
133  sort(ordered_files.begin(), ordered_files.end(), conditional_sort);
134
135  return ordered_files;
136  }

```

```

94  bool conditional_sort(pair<string, int> &a, pair<string, int> &b)
95  {
96      return a.second > b.second;
97  }
98

```

Para que fosse possível ordenar os arquivos por ordem de relevância, foi necessário transferi-los para um container vector, onde o valor de int, e utilizar a função sort() para ordenar os valores do container. A função sort recebe como parâmetro o início e fim do container que desejamos ordenar, e também a lógica de ordenação que foi implementada da função conditional_sort. Com isso o vector ordered_files é ordenado tanto pela quantidade de ocorrências quanto pela ordem lexicográfica, e finalmente retorna o vector. E por fim, para imprimir o resultado no terminal do usuário foi implementado a função print_ordered_files, que recebe os o vetor ordenado e o percorre para imprimir em tela.

```

138  void indexing::print_ordered_files(vector<string> query)
139  {
140      vector<pair<string, int>> ordered_files = this->recovery(query);
141
142      if (ordered_files.size() == 0)
143      {
144          cout << "\nNenhum arquivo encontrado!" << endl;
145      }
146      else
147      {
148          cout << "\nArquivos encontrados:" << endl;
149          for (int i = 0; i < ordered_files.size(); i++)
150          {
151              cout << ordered_files[i].first << " - " << ordered_files[i].second << " ocorrências" << endl;
152          }
153      }
154  }

```

```

O que deseja buscar?
casa quem apartamento quer

```

```

Arquivos encontrados:
d1.txt - 4 ocorrências
d2.txt - 3 ocorrências
d3.txt - 2 ocorrências

```

```

Deseja fazer outra pesquisa? (s/n)

```



```

5  int main()
6  {
7
8      string option, word, aux;
9      indexing search;
10     search.read_files("./documentos/");
11
12     do
13     {
14         vector<string> query;
15
16         cout << "O que deseja buscar?\n ";
17         getline(cin, word);
18
19         for (int i = 0; i <= word.length(); i++)
20         {
21             if (word[i] == ' ' || i == word.length())
22             {
23                 query.push_back(aux);
24                 aux = "";
25             }
26             else
27             {
28                 aux += word[i];
29             }
30         }
31
32         search.print_ordered_files(query);
33
34     do
35     {
36         cout << "\n\nDeseja fazer outra pesquisa? (s/n)" << endl;
37         cin >> option;
38
39         if (option != "s" && option != "n")
40         {
41             cout << "Opção inválida!" << endl;
42         }
43     } while (option != "s" && option != "n");
44
45     if (option == "s")
46     {
47         cout << "\n-----\n"
48             << endl;
49     }
50
51     cin.ignore();
52
53 } while (option == "s");
54
55 return 0;
56
57 }

```

Para que o usuário tenha uma interface onde utilizar o programa, no main foi implementado um menu com do/while. O usuário insere a string de pesquisa na variável através do getline, essa string é particionada palavra por palavra dentro de um vector para facilitar a manipulação, e é chamado a função de recovery passando este vector como parâmetro, retornando o resultado que foi mostrado mais acima. Além disso, o usuário tem a opção de sair do programa ou realizar uma nova consulta.

E com o objetivo de garantir o desempenho correto de todos os métodos do projeto, foi feito uma série de testes de unidade, via doctest.h, colocando diversas formas de entradas para corroboração da corretude do algoritmo desenvolvido.

```

1  #define DOCTEST_CONFIG_IMPLEMENT_WITH_MAIN
2  #include "doctest.h"
3
4  #include "../include/indexing.hpp"
5
6  TEST_CASE("Testa Indexing::Indexing()")
7  {
8      indexing test;
9      CHECK(test.get_index().size() == 0);
10 }
11
12 TEST_CASE("Testa Indexing::read_directory()")
13 {
14     indexing test;
15     set<string> files = test.read_directory("./tests");
16     CHECK(files.size() == 3);
17     CHECK(files.find("test.cpp") != files.end());
18     CHECK(files.find("doctest.h") != files.end());
19     CHECK(files.find("documentos_test") != files.end());
20 }
21
22 TEST_CASE("Testa Indexing::read_files()")
23 {
24     indexing test;
25     test.read_files("./tests/documentos_test/");
26     CHECK(test.get_index().size() == 3);
27 }
28
29 TEST_CASE("Testa Indexing::Insert()")
30 {
31     indexing test;
32     test.insert("teste", "teste.txt");
33     CHECK(test.get_index().size() == 1);
34 }

```

Conclusão

Ao término da aplicação foi possível desenvolver não apenas conhecimentos a respeito da matéria e da linguagem C++, como também o desenvolvimento prático de um projeto.

A priori os fundamentos discutidos em aula foram fundamentais para o desenvolvimento do trabalho como boas práticas e modularização que permitiram uma boa compreensão e organização do código, bem como a programação orientada a objeto e controle de versão que tornaram possível o desenvolvimento do mesmo projeto em paralelo entre os integrantes.

Além dos conceitos discutidos em aula já citados, houve também a busca por informações mais técnicas da ferramenta como por exemplo as funções `isalpha()`, que verifica se o caracter é uma letra, `tolower()`, que transforma letras maiúsculas em minúsculas, além do `getline()`, que foi necessário utilizar a função `cin.ignore()` para que funcionasse mais de uma vez no menu do usuário.

Vale destacar que graças a esse TP descobrimos coisas como uma letra acentuada vale por 2 espaços numa string, por exemplo, a string `comAcento` tem o dobro do tamanho da `semAcento`, presentes no método `normalize`, e como um espaço sozinho gera um caractere inválido, necessitando percorrer esse vetor de duas em duas casas.

Concluindo, foi um projeto para treinar nossos conhecimentos básicos, mas também foi desafiador em certos pontos, o que agregou bastante para o aprendizado na disciplina.