

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.

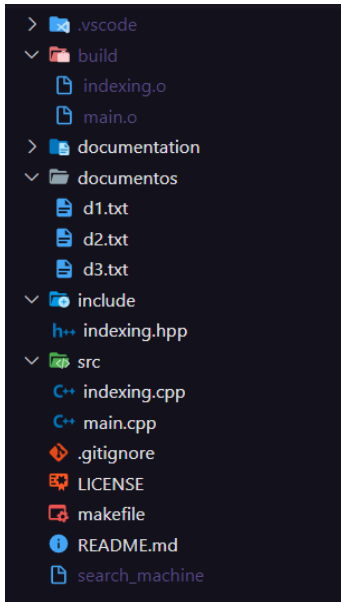
Implementação

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 version 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.



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.

```

1  #ifndef _INDEXING_H_
2  #define _INDEXING_H_
3
4  #include <map>
5  #include <set>
6  #include <vector>
7  #include <string>
8  #include <iostream>
9  #include <bits/stdc++.h>
10 #include <cstdio>
11 #include <cerrno>
12 #include <dirent.h>
13
14 using namespace std;
15
16 class indexing
17 {
18 private:
19     map<string, set<string>> index;
20
21     void insert(string, string);
22     string normalize(string);
23
24 public:
25     indexing();
26     ~indexing();
27
28     void recovery(vector<string>);
29 };
30
31 #endif

```

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. Além disso os métodos insert e normalize estão no como privados pois não há necessidade de outros arquivos os acessarem, garantindo mais segurança para eles, enquanto o construtor, destrutor e recuperação serão acessados pela main.cpp para a pesquisa do usuário.

Partindo para implementação do indexing, como só seria necessário ler uma vez a base de dados, pareceu conveniente a funcionalidade de leitura dos arquivos estarem no construtor, então assim que o usuário iniciar o programa essa base de dados será lida.

```

3 indexing::indexing()
4 {
5     // Listando os arquivos do diretório ./Documentos
6     vector<string> files;
7
8     struct dirent **namelist;
9     int n;
10
11     n = scandir("./documentos", &namelist, NULL, alphasort);
12
13     if (n < 0)
14         perror("scandir");
15     else
16     {
17         while (n > 0)
18         {
19             if (namelist[n - 1]->d_name[0] != '.')
20                 files.push_back(namelist[n - 1]->d_name);
21             delete namelist[n - 1];
22             n--;
23         }
24
25         delete namelist;
26     }
27

```

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 as bibliotecas cstdio, cerrno e dirent.h, foi possível percorrer o nome de cada arquivo e guardá-los em um container vector.

```

28 // Inserindo as palavras dos arquivos no índice
29 fstream file;
30 string word, filename;
31
32 for (int i = 0; i < files.size(); i++)
33 {
34     filename = files[i];
35     file.open("./documentos/" + filename);
36
37     while (file >> word)
38     {
39         this->insert(word, filename);
40     }
41
42     file.close();
43 }
44 }

```

Agora utilizando a biblioteca bits/stdc++, foi possível abrir cada arquivo com o nome salvo no vector 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     this->index[normalized_word].insert(filename);
54 }

```

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.


```

89 void indexing::recovery(vector<string> query)
90 {
91     set<string> normalized_query;
92     string aux;
93
94     // Normalizando a query
95     for (int i = 0; i < query.size(); i++)
96     {
97         aux = this->normalize(query[i]);
98         if (aux != "")
99         {
100             normalized_query.insert(aux);
101         }
102     }
103
104     // Contando a quantidade de vezes que cada arquivo aparece na query
105     map<string, int> relevant_files;
106
107     for (auto it = this->index.begin(); it != this->index.end(); it++)
108     {
109         if (normalized_query.find(it->first) != normalized_query.end())
110         {
111             for (auto it2 = it->second.begin(); it2 != it->second.end(); it2++)
112             {
113                 relevant_files[*it2]++;
114             }
115         }
116     }

```

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.

```

118 if (relevant_files.size() == 0)
119 {
120     cout << "\nNenhum documento encontrado!" << endl;
121 }
122 else
123 {
124
125     // Ordenando os arquivos por ocorrências na query
126     set<pair<int, string>> ordered_files;
127
128     for (auto it = relevant_files.begin(); it != relevant_files.end(); it++)
129     {
130         ordered_files.insert(make_pair(it->second, it->first));
131     }
132
133     cout << "\nDocumentos encontrados:" << endl;
134
135     for (auto it = ordered_files.rbegin(); it != ordered_files.rend(); it++)
136     {
137         cout << it->second << " - " << it->first << " ocorrências" << endl;
138     }
139 }

```

Para que fosse possível ordenar os arquivos por ordem de relevância, foi necessário transferi-los para um container set, onde o valor de int, que representa a quantidade de vezes que aquele arquivo foi contabilizado na pesquisa, seria o parâmetro de ordenação do set. E no fim, imprimindo o container ordered_files teremos esse resultado no terminal:

```
O que deseja buscar?
casa quem apartamento

Documentos encontrados:
d1.txt - 3 ocorrências
d3.txt - 2 ocorrências
d2.txt - 2 ocorrências
```

Obs: Não conseguimos implementar a ordenação lexicográfica dos documentos que têm a mesma quantidade de ocorrências, pois o set implementado com pair<int, string> só leva em conta uma das chaves para ordenação, que no caso foi a chave int que representa as ocorrências. Pesquisando encontramos algumas soluções envolvendo bubble sort que vão além dos nossos conhecimentos.

```
7  int main()
8  {
9
10     string option, word, aux;
11     indexing search;
12
13     do
14     {
15         vector<string> query;
16
17         cout << "O que deseja buscar?\n ";
18         getline(cin, word);
19
20         for (int i = 0; i <= word.length(); i++)
21         {
22             if (word[i] == ' ' || i == word.length())
23             {
24                 query.push_back(aux);
25                 aux = "";
26             }
27             else
28             {
29                 aux += word[i];
30             }
31         }
32
33         search.recovery(query);
34
35     }
36     do
37     {
38         cout << "\nDeseja fazer outra pesquisa? (s/n)" << endl;
39         cin >> option;
40
41         if (option != "s" && option != "n")
42         {
43             cout << "Opção inválida!" << endl;
44         }
45     } while (option != "s" && option != "n");
46
47     if (option == "s")
48     {
49         cout << "\n-----\n"
50             << endl;
51     }
52
53     cin.ignore();
54
55 } while (option == "s");
56
57 return 0;
58 }
```

Por fim, 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.

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.