

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №3 по курсу «Информационный поиск»

Студент: И. А. Мариничев
Преподаватель: А. А. Кухтичев
Группа: М8О-408Б-19
Дата:
Оценка:
Подпись:

Москва, 2023

Лабораторная работа №3 «Булев индекс»

Требуется построить поисковый индекс, пригодный для булева поиска, по подготовленному в ЛР1 корпусу документов. Требования к индексу:

- Самостоятельно разработанный, бинарный формат представления данных. Формат необходимо описать в отчёте, в побайтовом (или побитовом) представлении.
- Формат должен предполагать расширение, т.к. в следующих работах он будет меняться под требования новых лабораторных работ.
- Использование текстового представления или готовых баз данных не допускается.
- Кроме обратного индекса, должен быть создан «прямой» индекс, содержащий в себе как минимум заголовки документов и ссылки на них (понадобятся для выполнения ЛР4, при генерации страницы поисковой выдачи).
- Для термов должна быть как минимум понижена капитализация.

В отчёте должно быть отмечено как минимум:

- Выбранное внутренне представление документов после токенизации.
- Выбранный метод сортировки, его достоинства и недостатки для задачи индексации.

Среди результатов и выводов работы нужно указать:

- Количество термов.
- Средняя длина терма. Сравнить со средней длиной токена, вычисленной в ЛР1 по курсу ОТЕЯ. Объяснить причину отличий.
- Скорость индексации: общую, в расчёте на один документ, на килобайт текста.
- Оптимальна ли работа индексации? Что можно ускорить? Каким образом? Чем она ограничена? Что произойдёт, если объём входных данных увеличится в 10 раз, в 100 раз, в 1000 раз?

1 Описание

По выбранному корпусу документов строятся прямой и обратный индексы. Построение индекса занимает 379853 ms (6.33088333 min), сохранение занимает 1.14806e+06 ms (19.13433333 min).

Прямой индекс имеет следующее представление:

```
<Documents Total>
<Title 0> <URL 0> <Word Count 0>
<Title 1> <URL 1> <Word Count 1>
...
<Title N> <URL N> <Word Count N>
```

Обратный индекс имеет следующее представление:

```
<Documents Total>
<Word 0> <Postings Size 0> <Postings 0>
<Word 1> <Postings Size 1> <Postings 1>
...
<Word N> <Postings Size N> <Postings N>
```

2 Исходный код

Ниже приведен класс индекса, реализующий функционал построения и сохранения индекса в файл:

```
1 // index.h
2 #include <iostream>
3 #include <fstream>
4 #include <sstream>
5 #include <locale>
6 #include <vector>
7 #include <string>
8 #include <unordered_map>
9 #include <unordered_set>
10 #include <memory>
11 #include <chrono>
12
13 struct Document
14 {
15     std::wstring title;
16     std::wstring url;
17     uint32_t wordCount;
18 };
19
20 class Index
21 {
22 public:
23     void Build(std::string &inputFile);
24
25     void Save(std::string &outputFile);
26     void SaveIndex(std::string &outputFile);
27     void SaveInvertedIndex(std::string &outputFile);
28
29     void Load(std::string &inputFile);
30     void LoadIndex(std::string &inputFile);
31     void LoadInvertedIndex(std::string &inputFile);
32
33     friend class Query;
34
35 private:
36     uint32_t docTotal;
37     std::vector<Document> docIndex;
38     std::unordered_map<std::wstring, std::vector<uint32_t>> invertedIndex;
39
40     bool GetDocInfo(std::wifstream &input);
41     void GetDocContents(uint32_t id, std::wifstream &input);
42 };
43
44 1 // index.cpp
45 2 #include "index.h"
```

```

3
4 void Index::Build(std::string &inputFile)
5 {
6     std::cout << "Building index...\n";
7     auto start = std::chrono::high_resolution_clock::now();
8
9     std::wifstream input((inputFile).c_str());
10    for (size_t id = 1; input; ++id)
11    {
12        if (GetDocInfo(input))
13        {
14            GetDocContents(id, input);
15        }
16    }
17    docTotal = docIndex.size();
18    input.close();
19
20    auto end = std::chrono::high_resolution_clock::now();
21    std::cout << "Successfully built index in "
22              << std::chrono::duration<double, std::milli>(end - start).count() << " ms\
23              n";
24 }
25 bool Index::GetDocInfo(std::wifstream &input)
26 {
27     uint32_t length;
28     std::wstring line, url, title;
29     Document doc;
30     if (std::getline(input, line))
31     {
32         uint32_t url_start = line.find(L"url=");
33         if (url_start == std::string::npos)
34         {
35             return false;
36         }
37         uint32_t title_start = line.find(L"title=");
38
39         url_start += 4; // len("url=")
40         length = title_start - url_start - 1;
41         url = line.substr(url_start, length);
42
43         title_start += 6; // len("title=")
44         length = line.length() - title_start - 1;
45         title = line.substr(title_start, length);
46
47         doc.title = title;
48         doc.url = url;
49         docIndex.push_back(doc);
50

```

```

51         return true;
52     }
53     else
54     {
55         return false;
56     }
57 }
58
59 void Index::GetDocContents(uint32_t id, std::wifstream &input)
60 {
61     std::wstring line, word;
62     std::unordered_set<std::wstring> words;
63     uint32_t wordCount = 0;
64     while (std::getline(input, line))
65     {
66         if (line == L"</doc>")
67         {
68             break;
69         }
70
71         for (const auto &c : line)
72         {
73             if (isalnum(c, std::locale()))
74             {
75                 word += tolower(c);
76             }
77             else if (word.length())
78             {
79                 words.insert(word);
80                 wordCount++;
81                 word.clear();
82             }
83         }
84
85         if (word.length())
86         {
87             words.insert(word);
88             wordCount++;
89             word.clear();
90         }
91     }
92
93     for (const auto &w : words)
94     {
95         const auto [it, success] = invertedIndex.insert(
96             std::make_pair(w, std::vector<uint32_t>()));
97         it->second.push_back(id - 1);
98     }
99

```

```

100     docIndex[id - 1].wordCount = wordCount;
101 }
102
103 void Index::Save(std::string &outputFile)
104 {
105     std::cout << "Saving index...\n";
106     auto start = std::chrono::high_resolution_clock::now();
107
108     SaveIndex(outputFile);
109     SaveInvertedIndex(outputFile);
110
111     auto end = std::chrono::high_resolution_clock::now();
112     std::cout << "Successfully saved index in "
113         << std::chrono::duration<double, std::milli>(end - start).count() << " ms\n"
114         << "\n";
115 }
116
117 void Index::SaveIndex(std::string &outputFile)
118 {
119     std::wofstream wFileOut(outputFile.c_str());
120     if (wFileOut)
121     {
122         wFileOut << docIndex.size() << L"\n";
123         for (const auto &doc : docIndex)
124         {
125             wFileOut << doc.title << L" " << doc.url << L" \"" << doc.wordCount << L"
126             << "\n";
127         }
128     }
129     wFileOut.close();
130 }
131
132 void Index::SaveInvertedIndex(std::string &outputFile)
133 {
134     std::wofstream wFileOut((outputFile + "_inverted").c_str());
135     if (wFileOut)
136     {
137         wFileOut << docTotal << L"\n";
138         for (const auto &it : invertedIndex)
139         {
140             wFileOut << it.first << L" ";
141             wFileOut << it.second.size() << L" ";
142             for (size_t i = 0; i < it.second.size(); ++i)
143             {
144                 wFileOut << it.second[i] << L" ";
145             }
146             wFileOut << L"\n";
147         }
148     }
149 }

```

```

147     wFileOut.close();
148 }
149
150 void Index::Load(std::string &inputFile)
151 {
152     std::cout << "Loading index...\n";
153     auto start = std::chrono::high_resolution_clock::now();
154
155     LoadIndex(inputFile);
156     LoadInvertedIndex(inputFile);
157
158     auto end = std::chrono::high_resolution_clock::now();
159     std::cout << "Successfully loaded index in "
160         << std::chrono::duration<double, std::milli>(end - start).count() << " ms\
161         n";
162 }
163
164 void Index::LoadIndex(std::string &inputFile)
165 {
166     std::wifstream wFileIn(inputFile.c_str());
167     std::wstring docLine, size;
168     uint32_t reserveSize;
169
170     std::getline(wFileIn, size);
171     std::wstringstream wSizeIn(size);
172     wSizeIn >> reserveSize;
173
174     Document doc;
175     docIndex.reserve(reserveSize);
176
177     size_t pos = 0;
178     std::wstring delimiter = L"\" \";
179     while (std::getline(wFileIn, docLine))
180     {
181         pos = docLine.find(delimiter);
182         doc.title = docLine.substr(0, pos + 1);
183         docLine.erase(0, pos + delimiter.length() - 1);
184
185         pos = docLine.find(delimiter);
186         doc.url = docLine.substr(0, pos + 1);
187         docLine.erase(0, pos + delimiter.length());
188
189         std::wstringstream wWordCountIn(docLine.substr(0, docLine.length() - 1));
190         wWordCountIn >> doc.wordCount;
191
192         docIndex.push_back(doc);
193
194         docLine.clear();
195     }

```



```

195     wFileIn.close();
196 }
197
198 void Index::LoadInvertedIndex(std::string &inputFile)
199 {
200     std::wifstream wFileIn((inputFile + "_inverted").c_str());
201     uint32_t size;
202     std::wstring word;
203
204     wFileIn >> docTotal;
205     while (!wFileIn.eof())
206     {
207         std::vector<uint32_t> postings;
208         wFileIn >> word;
209         wFileIn >> size;
210
211         uint32_t posting;
212         postings.reserve(size);
213         for (size_t i = 0; i < size; ++i)
214         {
215             wFileIn >> posting;
216             postings.push_back(posting);
217         }
218
219         invertedIndex.insert(std::make_pair(word, postings));
220         word.clear();
221     }
222     wFileIn.close();
223 }

```

3 Выводы

Выполнив третью лабораторную работу по курсу «Информационный поиск», я познакомился с двумя основными структурами данных в этой области:

- **Прямой индекс** — список документов, в котором можно найти этот документ по его id. Иными словами, прямой индекс — это массив строк (вектор документов), где id документа — это его индекс.
- **Обратный индекс** — список слов, которые мы «выпотрошили» из всех документов. За каждым словом закреплён список отсортированных id документов (posting list), в которых это слово встретилось.

Список литературы

- [1] Маннинг, Рагхаван, Шютце *Введение в информационный поиск* — Издательский дом «Вильямс», 2011. Перевод с английского: доктор физ.-мат. наук Д. А. Ключина — 528 с. (ISBN 978-5-8459-1623-4 (рус.))