

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №2 по курсу «Информационный поиск»

Студент: П. А. Харьков
Преподаватель: А. А. Кухтичев
Группа: М8О-406Б
Дата:
Оценка:
Подпись:

Москва, 2023

Лабораторная работа №2

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

- Самостоятельно разработанный бинарный формат представления данных. Формат необходимо описать в отчете, в побайтовой 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.
- Для термов должна быть как минимум понижена капитализация.

В отчете должно быть отмечено как минимум:

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

1 Описание

Индексирование происходит следующим образом:

1. Выбирается обработанный в ЛР1 документ.
2. Из него считываются первые три строки: каждая отвечает за свою зону.
3. Затем для каждого слова в строке:
 - (а) В обратный индекс для слова в зоне добавляется номер документа.
 - (б) В прямой индекс для документа в зоне добавляется номер слова.

Для работы программы необходимо хранить словарь названий файлов, словарь слов, прямой индекс и обратный индекс. Каждый из этих объектов я храню в отдельном файле. Объекты сохраняются в бинарный файл наивным образом, в том виде, как они и хранятся в памяти.

2 Исходный код

```
1 word_id SearchEngine::get_word_id(const string& word) {
2     static word_id next_word_id = 0;
3     auto it = words_dict.find(word);
4     if (it == words_dict.end()) {
5         words_dict[word] = next_word_id;
6         return next_word_id++;
7     }
8     return it->second;
9 }
10 void SearchEngine::index_zone(istringstream& iss, doc_id doc_counter, Zone zone) {
11     string word;
12     while (iss >> word) {
13         word_id id = get_word_id(word);
14         inverted_index[word][zone].push_back(doc_counter);
15         forward_index[doc_counter][zone].push_back(id);
16         doc_zone_word_count[doc_counter][zone]++;
17     }
18 }
19 void SearchEngine::index_file(const string& file_path, int doc_counter) {
20     ifstream infile(file_path);
21     string title, first_paragraph, rest;
22
23     getline(infile, title);
24     getline(infile, first_paragraph);
25     getline(infile, rest);
26
27     infile.close();
28
29     istringstream iss_title(title), iss_first(first_paragraph), iss_rest(rest);
30     index_zone(iss_title, doc_counter, TITLE);
31     index_zone(iss_first, doc_counter, FIRST_PARAGRAPH);
32     index_zone(iss_rest, doc_counter, REST);
33 }
34 void SearchEngine::index_folder() {
35     doc_id doc_counter = 0;
36     for (const auto& entry : filesystem::directory_iterator(INPUT_DIR)) {
37         if (entry.is_regular_file()) {
38             string file_path = entry.path().string();
39             docs_dict[doc_counter] = file_path;
40             index_file(file_path, doc_counter);
41             doc_counter++;
42         }
43     }
44     compute_tf_idf();
45 }
```

3 Выводы

Выполнив вторую лабораторную работу по курсу «Информационный поиск», я научился реализовывать прямой и обратные индексы для того, чтобы в последствии использовать их для булева поиска. Наибольшей сложностью для меня в этой лабораторной работе было не написание правильной индексации файлов, а написание сохранения и загрузки проиндексированных данных в бинарные файлы.

Список литературы

- [1] Маннинг, Кристофер Д. Введение в информационный поиск [Текст] / Кристофер Д. Маннинг, Прабхакар Рагхаван, Хайнрих Шютце ; пер. с англ. М. Л. Суркова. - Москва : Вильямс, 2020. - 528 с.