

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №4 по курсу «Информационный поиск»

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

Москва, 2023

Лабораторная работа №4

Необходимо сделать ранжированный поиск на основании схемы ранжирования TF-IDF.

1 Описание

После индексации всех файлов, я начинаю индексировать данные для ранжирования с помощью TF-IDF. Для каждого слова в каждой зоне обратного индекса, я прохожусь по документу:

- Подсчитывается количество вхождений слов в зоне документа.
- Вычисляется TF по формуле

$$\text{TF}(t, d) = \frac{f_{t,d}}{\max f_{w,d} : w \in d}$$

где $f_{t,d}$ - количество раз, когда термин t появляется в документе d , а $\max f_{w,d} : w \in d$ - максимальное количество раз, когда любой термин появляется в документе d .

- Вычисляется IDF по формуле

$$\text{IDF}(t, D) = \log \frac{|D|}{|d \in D : t \in d|}$$

где $|D|$ - общее количество документов в коллекции, а $|d \in D : t \in d|$ - количество документов в коллекции, которые содержат термин t .

- Получается значение

$$\text{TF-IDF}(t, d, D) = \text{TF}(t, d) \times \text{IDF}(t, D)$$

2 Исходный код

```
1 double SearchEngine::tf(int term_count, int zone_word_count) {
2     return (double)term_count / zone_word_count;
3 }
4
5 double SearchEngine::idf(int doc_count, int total_docs) {
6     return log((double)total_docs / (1 + doc_count));
7 }
8
9 int SearchEngine::term_count(doc_id doc, Zone zone, const string& word) {
10     return count(forward_index[doc][zone].begin(), forward_index[doc][zone].end(),
11                 get_word_id(word));
12 }
13
14 void SearchEngine::compute_tf_idf() {
15     int total_docs = docs_dict.size();
16     for (const auto& word_entry : inverted_index) {
17         const string& word = word_entry.first;
18         for (const auto& zone_entry : word_entry.second) {
19             Zone zone = zone_entry.first;
20             const vector<doc_id>& doc_ids = zone_entry.second;
21             int doc_count = doc_ids.size();
22             for (doc_id doc : doc_ids) {
23                 int term_cnt = term_count(doc, zone, word);
24                 int zone_word_count = doc_zone_word_count[doc][zone];
25                 double tf_value = tf(term_cnt, zone_word_count);
26                 double idf_value = idf(doc_count, total_docs);
27                 tf_idf_index[doc][zone][get_word_id(word)] = abs(tf_value * idf_value);
28             }
29         }
30     }
```

3 Выводы

Выполнив четвертую лабораторную работу по курсу «Информационный поиск», я научился реализовывать ранжирование TF-IDF. В итоге, у меня получилось улучшить результаты при поисковой выдаче.

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.