

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

Москва, 2023

Лабораторная работа №4 «Булев поиск»

Нужно реализовать ввод поисковых запросов и их выполнение над индексом, получение поисковой выдачи. Синтаксис поисковых запросов:

- Размер «сырых» данных.
- Пробел или два амперсанда, «&&», соответствуют логической операции «И».
- Две вертикальных «палочки», «||» – логическая операция «ИЛИ»
- Восклицательный знак, «!» – логическая операция «НЕТ»
- Могут использоваться скобки.

Парсер поисковых запросов должен быть устойчив к переменному числу пробелов, максимально толерантен к введённому поисковому запросу.

Примеры запросов:

- [московский авиационный институт]
- [(красный || желтый) автомобиль]
- [руки !ноги]

Для демонстрации работы поисковой системы должен быть реализован веб-сервис, реализующий базовую функциональность поиска из двух страниц:

- Начальная страница с формой ввода поискового запроса.
- Страница поисковой выдачи, содержащая в себе форму ввода поискового запроса, 50 результатов поиска в виде текстов заголовков документов и ссылок на эти документы, а так же ссылку на получение следующих 50 результатов.

Так же должна быть реализована утилита командной строки, загружающая индекс и выполняющая поиск по нему для каждого запроса на отдельной строчке входного файла. В отчёте должно быть отмечено:

- Скорость выполнения поисковых запросов.
- Примеры сложных поисковых запросов, вызывающих длительную работу.
- Каким образом тестировалась корректность поисковой выдачи

1 Описание

Наши запросы представлены в стандартной булевой модели. Она использует алгебру логики и теорию множеств. Булевы выражения строятся над множеством, над элементами которого определены три операции – конъюнкция, дизъюнкция и отрицание:

1. Оператор «&&» указывает на то, что оба слова должны присутствовать в статье.
2. Оператор «||» указывает на то, что одно из слов или оба слова должны присутствовать в статье.
3. Оператор «!» указывает на то, что слово не должно присутствовать в статье.

Основные положения при вычислении:

1. Операции выражения в скобках выполняются первыми, т.е. имеют наивысший приоритет. Для выражения, находящегося в скобках, правила вычисления выражения те же, что и для обычного выражения.
2. Наивысший приоритет из трех возможных операций имеет операция отрицания. Затем операция пересечения, затем операция объединения. При отсутствии скобок операции с наибольшим приоритетом выполняются раньше, чем операции с меньшим приоритетом.
3. Операции с равными приоритетами выполняются в порядке их появления в выражении.

На каждом шаге у нас шесть возможных вариантов действий. Для выбора того или иного действия необходимо проверить сочетание состояния вершины стека операций и рассматриваемый символ. Рассмотрим варианты действий.

Для символов «||» и «&&» правила:

1. Если стек пуст, поместим операцию в стек и перейдем к следующему символу.
2. Пока на вершине стека операция с равным или большим приоритетом, выполнить ее, поместить рассматриваемую в данный момент операцию в стек и перейти к следующему символу.
3. Если на вершине стека операция с меньшим приоритетом, поместить рассматриваемую в данный момент операцию в стек.

Для символа «!» правила:

1. Если стек пуст, поместим операцию в стек и перейдем к следующему символу.
2. Пока на вершине стека операция с большим приоритетом, выполнить ее, поместить рассматриваемую в данный момент операцию в стек и перейти к следующему символу.
3. Если на вершине стека операция с меньшим приоритетом, поместить рассматриваемую в данный момент операцию в стек.

В случае левой скобки необходимо добавить символ в стек операций и перейти к следующему символу.

В случае правой скобки возможны следующие варианты:

1. Стек пуст. Следовательно, не хватает левой скобки, т.е. выражение ошибочно.
2. Выполняем операции из вершины стека, пока операция на вершине стека не левая скобка.

И наконец, рассмотрим действия, необходимые для обработки символа конца выражения:

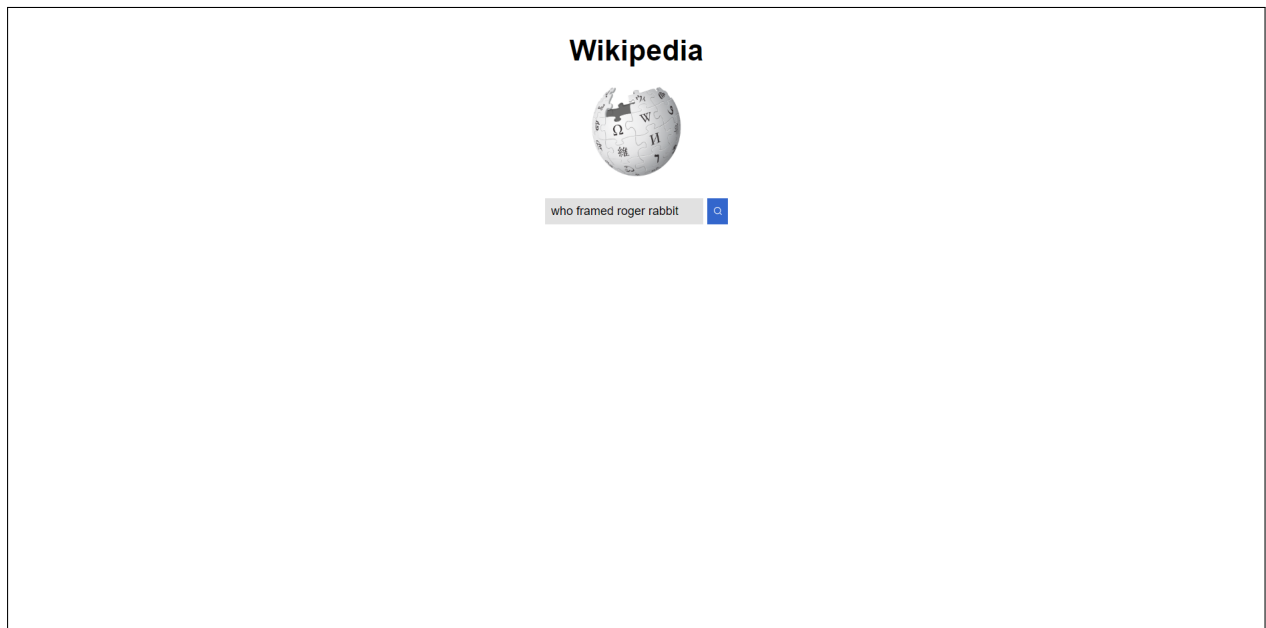
1. Если стек пуст, то выражение вычислено. Следовательно, нужно удалить из стека операндов результат и вывести его.
2. На вершине – левая скобка. Следовательно, левых скобок в выражении было больше, чем правых, т.е. выражение ошибочно.
3. На вершине стека – логическая операция. Необходимо удалить ее из стека и выполнить.

После обработки запроса выводим получившийся ответ.

Ниже приведены несколько запросов и результаты по ним.

Запрос	Время (в ms)	Количество найденных файлов
moscow && aviation && institute	12.1651	36
(red yellow) && automobile	12.0109	356
arms && !legs	8.3885	2997

Внешний вид веб-сервиса:



Here are your results for 'who framed roger rabbit':

- Who Framed Roger Rabbit
URL: <https://en.wikipedia.org/wiki?curid=76018>
- Bugs Bunny
URL: <https://en.wikipedia.org/wiki?curid=50286>
- Mickey Mouse universe
URL: <https://en.wikipedia.org/wiki?curid=77553>
- Song of the South
URL: <https://en.wikipedia.org/wiki?curid=61141>
- Steven Spielberg
URL: <https://en.wikipedia.org/wiki?curid=26940>
- The Walt Disney Company
URL: <https://en.wikipedia.org/wiki?curid=37398>

2 Исходный код

Ниже приведен класс булева поиска, реализующий функционал обработки пользовательских запросов и записи результатов в файл:

```
1 // query.h
2 #include "index.h"
3
4 #include <stack>
5 #include <set>
6 #include <algorithm>
7 #include <cmath>
8
9 #define AND L'&'
10 #define OR L'|'
11 #define NOT L'!'
12 #define LEFT_BRACKET L'('
13 #define RIGHT_BRACKET L')'
14 #define SPACE L' '
15
16
17 class Query
18 {
19 public:
20     void GetIndex(std::string &inputFile);
21     void ParseQueries(std::string &outputFile);
22     void ParseQueriesFromFile(std::string &inputFile, std::string &outputFile);
23
24 private:
25     Index index;
26     std::stack<std::shared_ptr<std::vector<uint32_t>>> operands;
27     std::stack<wchar_t> operations;
28
29     const std::vector<uint32_t> &GetDocIndices(std::wstring &word);
30
31     void ProcessingQuery(std::wstring &query);
32
33     bool IsOperation(wchar_t op);
34     uint32_t Priority(wchar_t op);
35
36     void ExecuteOperation(wchar_t op);
37     void Union();
38     void Intersection();
39     void Negation();
40 };
41
42 // query.cpp
43 #include "query.h"
44
45 void Query::GetIndex(std::string &inputFile)
```

```

5 | {
6 |     index.Load(inputFile);
7 | }
8 |
9 | void Query::ParseQueries(std::string &outputFile)
10 | {
11 |     std::wstring query;
12 |     bool isFuzzy = false;
13 |     while (std::getline(std::wcin, query))
14 |     {
15 |         if (!query.length())
16 |         {
17 |             break;
18 |         }
19 |
20 |         std::wofstream wFileOut(outputFile.c_str());
21 |
22 |         isFuzzy = IsFuzzy(query);
23 |         auto start = std::chrono::high_resolution_clock::now();
24 |         if (isFuzzy)
25 |         {
26 |             ProcessingFuzzyQuery(query);
27 |         }
28 |         else
29 |         {
30 |             ProcessingQuery(query);
31 |         }
32 |         std::shared_ptr<std::vector<uint32_t>> result_ptr = operands.top();
33 |         operands.pop();
34 |         if (isFuzzy)
35 |         {
36 |             Ranking(*result_ptr, query);
37 |         }
38 |         auto end = std::chrono::high_resolution_clock::now();
39 |
40 |         std::cout << "Found " << (*result_ptr).size() << " result(s) in "
41 |             << std::chrono::duration<double, std::milli>(end - start).count() << "
42 |                 ms\n";
43 |
44 |         wFileOut << (*result_ptr).size() << L'\n';
45 |         for (const auto &i : (*result_ptr))
46 |         {
47 |             wFileOut << index.docIndex[i].title << L' ' << index.docIndex[i].url << L'\n';
48 |         }
49 |         wFileOut.close();
50 |     }
51 | }

```

```

52 void Query::ParseQueriesFromFile(std::string &inputFile, std::string &outputFile)
53 {
54     std::wifstream wFileIn(inputFile.c_str());
55     std::wofstream wFileOut(outputFile.c_str());
56
57     std::wstring query;
58     bool isFuzzy = false;
59     while (std::getline(wFileIn, query))
60     {
61         if (!query.length())
62         {
63             break;
64         }
65
66         isFuzzy = IsFuzzy(query);
67         auto start = std::chrono::high_resolution_clock::now();
68         if (isFuzzy)
69         {
70             ProcessingFuzzyQuery(query);
71         }
72         else
73         {
74             ProcessingQuery(query);
75         }
76         std::shared_ptr<std::vector<uint32_t>> result_ptr = operands.top();
77         operands.pop();
78         if (isFuzzy)
79         {
80             Ranking(*result_ptr, query);
81         }
82         auto end = std::chrono::high_resolution_clock::now();
83
84         std::cout << "Found " << (*result_ptr).size() << " result(s) in "
85                 << std::chrono::duration<double, std::milli>(end - start).count() << "
86                 << "ms\n";
87
88         wFileOut << (*result_ptr).size() << L'\n';
89         for (const auto &i : (*result_ptr))
90         {
91             wFileOut << index.docIndex[i].title << L' ' << index.docIndex[i].url << L'\n';
92         }
93     }
94     wFileIn.close();
95     wFileOut.close();
96 }
97
98 const std::vector<uint32_t> &Query::GetDocIndices(std::wstring &word)
99 {

```



```

99     static const std::vector<uint32_t> empty(0);
100
101     auto it = index.invertedIndex.find(word);
102     if (it != index.invertedIndex.end())
103     {
104         return it->second;
105     }
106     else
107     {
108         return empty;
109     }
110 }
111
112 void Query::ProcessingQuery(std::wstring &query)
113 {
114     std::wstring word;
115     std::wstring quote;
116     for (size_t i = 0; i < query.length(); ++i)
117     {
118         wchar_t c = query[i];
119
120         if (query.substr(i, 4) == L" && ")
121         {
122             c = AND;
123             i += 3;
124         }
125         else if (query.substr(i, 4) == L" || ")
126         {
127             c = OR;
128             i += 3;
129         }
130         else if (c == SPACE)
131         {
132             c = AND;
133         }
134
135         if (c == QUOTE)
136         {
137             do
138             {
139                 i++;
140                 c = query[i];
141                 quote += tolower(c);
142             } while (query[i + 1] != QUOTE);
143             i++;
144             ProcessingQuote(quote);
145             quote.clear();
146             continue;
147         }

```

```

148
149     if ((IsOperation(c) || c == LEFT_BRACKET || c == RIGHT_BRACKET) && word.length
150         ())
151     {
152         auto postings = std::make_shared<std::vector<uint32_t>>(GetDocIndices(word)
153             );
154         operands.push(postings);
155         word.clear();
156     }
157
158     if (c == LEFT_BRACKET)
159     {
160         operations.push(LEFT_BRACKET);
161     }
162     else if (c == RIGHT_BRACKET)
163     {
164         while (operations.top() != LEFT_BRACKET)
165         {
166             ExecuteOperation(operations.top());
167             operations.pop();
168         }
169         operations.pop();
170     }
171     else if (IsOperation(c))
172     {
173         while (!operations.empty() && (((c != NOT) && (Priority(operations.top())
174             >= Priority(c))) ||
175             ((c == NOT) && (Priority(operations.top()) >
176                 Priority(c)))))
177         {
178             ExecuteOperation(operations.top());
179             operations.pop();
180         }
181         operations.push(c);
182     }
183     else
184     {
185         word += tolower(c);
186     }
187 }
188
189 if (word.length())
190 {
191     auto postings = std::make_shared<std::vector<uint32_t>>(GetDocIndices(word));
192     operands.push(postings);
193 }
194
195 while (!operations.empty())
196 {

```

```

193     wchar_t op = operations.top();
194     ExecuteOperation(op);
195     operations.pop();
196 }
197 }
198
199 bool Query::IsOperation(wchar_t op)
200 {
201     return (op == AND) || (op == OR) || (op == NOT);
202 }
203
204 uint32_t Query::Priority(wchar_t op)
205 {
206     switch (op)
207     {
208     case OR:
209         return 1;
210     case AND:
211         return 2;
212     case NOT:
213         return 3;
214     default:
215         return 0;
216     }
217 }
218
219 void Query::ExecuteOperation(wchar_t op)
220 {
221     switch (op)
222     {
223     case OR:
224         Union();
225         break;
226     case AND:
227         Intersection();
228         break;
229     case NOT:
230         Negation();
231         break;
232     default:
233         break;
234     }
235 }
236
237 void Query::Union()
238 {
239     std::vector<uint32_t> result;
240     std::shared_ptr<std::vector<uint32_t>> postings1_ptr = operands.top();
241     operands.pop();

```

```

242     std::shared_ptr<std::vector<uint32_t>> postings2_ptr = operands.top();
243     operands.pop();
244     std::set_union((*postings1_ptr).begin(), (*postings1_ptr).end(),
245                  (*postings2_ptr).begin(), (*postings2_ptr).end(),
246                  std::back_inserter(result));
247     auto result_ptr = std::make_shared<std::vector<uint32_t>>(result);
248     operands.push(result_ptr);
249 }
250
251 void Query::Intersection()
252 {
253     std::vector<uint32_t> result;
254     std::shared_ptr<std::vector<uint32_t>> postings1_ptr = operands.top();
255     operands.pop();
256     std::shared_ptr<std::vector<uint32_t>> postings2_ptr = operands.top();
257     operands.pop();
258     std::set_intersection((*postings1_ptr).begin(), (*postings1_ptr).end(),
259                          (*postings2_ptr).begin(), (*postings2_ptr).end(),
260                          std::back_inserter(result));
261     auto result_ptr = std::make_shared<std::vector<uint32_t>>(result);
262     operands.push(result_ptr);
263 }
264
265 void Query::Negation()
266 {
267     std::vector<uint32_t> result;
268     std::shared_ptr<std::vector<uint32_t>> postings_ptr = operands.top();
269     operands.pop();
270     size_t j = 0;
271     for (uint32_t i = 0; i < index.docTotal; ++i)
272     {
273         if (j == (*postings_ptr).size() || i < (*postings_ptr)[j])
274         {
275             result.push_back(i);
276         }
277         else if (i == (*postings_ptr)[j])
278         {
279             ++j;
280         }
281     }
282     auto result_ptr = std::make_shared<std::vector<uint32_t>>(result);
283     operands.push(result_ptr);
284 }

```

3 Выводы

Выполнив четвертую лабораторную работу по курсу «Информационный поиск», я посторил свой булев поиск, который обрабатывает выражения в скобках и реализует все три оператора: «&&», «||», «!». Кроме того был реализован веб-сервис, написанный при помощи HTML и Flask.

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.))