

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа по курсу «Информационный поиск»

Студент: А. К. Киреев
Преподаватель: А. А. Кухтичев
Группа: М8О-406Б
Дата:
Оценка:
Подпись:

Москва, 2022

Лабораторная работа №1 «Добыча корпуса документов»

Необходимо подготовить корпус документов, который будет использован при выполнении остальных лабораторных работ:

- Скачать его к себе на компьютер. В отчёте нужно указать источник данных.
- Ознакомиться с ним, изучить его характеристики. Из чего состоит текст? Есть ли дополнительная мета-информация? Если разметка текста, какая она?
- Разбить на документы.
- Выделить текст.
- Найти существующие поисковики, которые уже можно использовать для поиска по выбранному набору документов (встроенный поиск Википедии, поиск Google с использованием ограничений на URL или на сайт). Если такого поиска найти невозможно, то использовать корпус для выполнения лабораторных работ нельзя!
- Привести несколько примеров запросов к существующим поисковикам, указать недостатки в полученной поисковой выдаче.

В результатах работы должна быть указаны статистическая информация о корпусе:

- Размер «сырых» данных.
- Количество документов.
- Размер текста, выделенного из «сырых» данных.
- Средний размер документа, средний объём текста в документе.

Лабораторная работа №2 «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.
- 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.

Лабораторная работа №3 «Булев поиск»

Нужно реализовать ввод поисковых запросов и их выполнение над индексом, получение поисковой выдачи. 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. дополнительная мета-информация? Если разметка текста, какая она?

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.

Лабораторная работа №4 «Ранжирование TF-IDF»

Необходимо сделать ранжированный поиск на основании схемы ранжирования TF-IDF. Теперь, если запрос содержит в себе только термины через пробелы, то его надо трактовать как нечёткий запрос, т.е. допускать неполное соответствие документа терминам запроса и т.п.

В отчёте нужно привести несколько примеров выполнения запросов, как удачных, так и не удачных.

1 Описание

Лабораторная работа №1 «Добыча корпуса документов»

Для подготовки корпуса было принято решение написать собственного скрапера. В качестве сайта была выбрана английская википедия: <https://en.wikipedia.org>. Для настройки робота используем аргументы командной строки и данные из robots.txt. Начиная с главной страницы, скачиваем страницу в формате html. На каждой странице находим все ссылки, фильтруем эти ссылки (есть пути запрещенные сайтом и ссылки ведущие на другие сайты или посещенные ссылки). Одновременно с этим получаем видимый текст статьи и при помощи фреймворка nltk нормализуем его. Эта библиотека производит лексический анализ и выдает значимые слова в приведенной форме. Сохраняем полученный текст, ссылку и заголовок текста в mongodb.

Количество документов: 20000

Размер текста: 423 МВ

Средний размер документа: 36 КВ

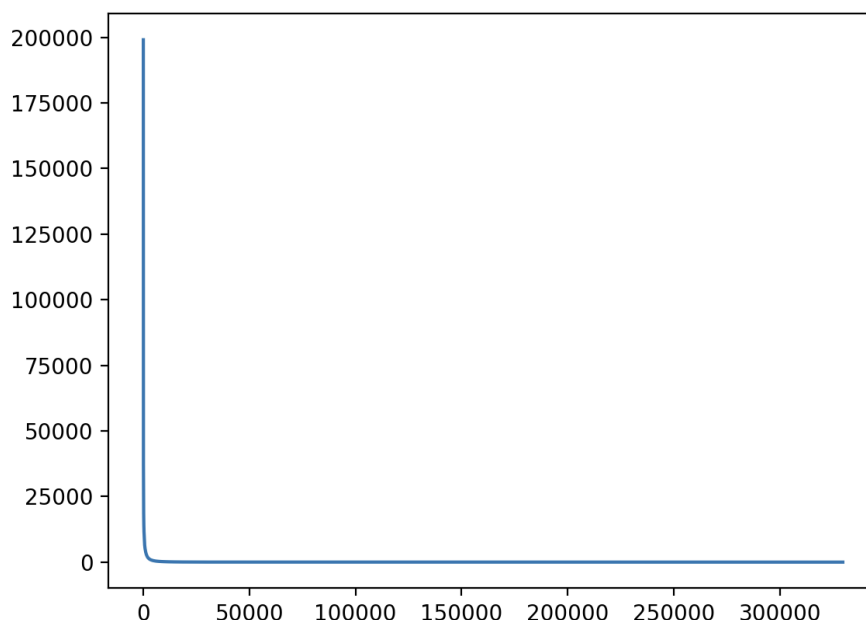


Рис. 1: Закон Ципфа для корпуса

2 Исходный код

```
scraper.py

1  import asyncio
2  import coloredlogs
3  import logging
4
5  from aiolimiter import AsyncLimiter
6  from httpx import AsyncClient, Response, codes
7  from bs4 import BeautifulSoup
8  from robots import RobotsParser
9  from typing import Tuple, Set, List
10 from functools import partial, wraps
11
12 from utils import Settings
13 from bucket_queue import BucketQueue
14 from db import save_documents, save_words, save_bigrams,\
15     load_visited_urls, dump_visited_urls, load_pending_urls, dump_pending_urls
16 from text_enrich import enrich_text, get_text
17
18
19 class Scraper:
20     logger = logging.getLogger(__name__)
21
22     def __init__(self, settings: Settings):
23         self.settings: Settings = settings
24         Scraper.logger.setLevel(level=self.settings.log_level)
25         coloredlogs.install(level=Scraper.logger.level)
26
27         self.robots = RobotsParser.from_uri(uri=f"{self.settings.url_base}/robots.txt")
28         if self.settings.rps > 50:
29             Scraper.logger.warning("Your rps(%s) too large", self.settings.rps)
30         self.throttler = AsyncLimiter(max_rate=self.settings.rps, time_period=1)
31
32     async def _load_scraper_state(self):
33         loaded_visited_urls: Set[str] = await load_visited_urls(Scraper.logger)
34         if loaded_visited_urls:
35             self.visited = loaded_visited_urls
36         else:
37             self.visited = {self.settings.start_url}
38             await dump_visited_urls({self.settings.start_url}, Scraper.logger)
39
40         self.queue = BucketQueue(settings=self.settings)
41         loaded_pending_urls: Set[str] = await load_pending_urls(Scraper.logger)
```

```

42     if loaded_pending_urls:
43         self.queue.extend(loaded_pending_urls)
44     else:
45         self.queue.extend(self.visited)
46
47     def _filter_urls(self, urls: Set[str]) -> Set[str]:
48         filtered: Set[str] = {f'{self.settings.url_base}{url}' for url in urls
49             if url.startswith('/wiki/')
50             and self.robots.can_fetch("*", url)
51             and ":" not in url
52             and f'{self.settings.url_base}{url}' not in self.visited}
53         Scraper.logger.info("Discarded %s refs from current batch of documents",
54             ↳ len(urls) - len(filtered))
55         return filtered
56
57     def _find_hrefs(self, sources: Tuple) -> Set[str]:
58         result: Set[str] = set()
59         for s in sources:
60             if not s:
61                 continue
62             soup = BeautifulSoup(s.text, features="html.parser")
63             result |= {a['href'] for a in soup.find_all('a', href=True)}
64             Scraper.logger.info("Got %s refs from current batch of documents", len(result))
65             return result
66
67     @staticmethod
68     def _handle_response(task):
69         @wraps(task)
70         async def wrapper(self, url: str, client: AsyncClient):
71             retries: int = 0
72             max_retries: int = 5
73             response: Response = await task(self, url, client)
74
75             while ((response.status_code == codes.TOO_MANY_REQUESTS or
76                 ↳ response.status_code == codes.SERVICE_UNAVAILABLE)
77                 and retries <= max_retries):
78                 Scraper.logger.warning("Too many request for url: %s, retry: %s", url,
79                     ↳ retries)
80                 await asyncio.sleep(1)
81                 response = await task(self, url, url, client)
82                 retries += 1
83
84             if retries > max_retries:
85                 Scraper.logger.error("Can't scrape url: %s, max retries was reached", url)

```

```

83         return None
84
85     if codes.is_redirect(response.status_code):
86         Scraper.logger.warning("Redirect on url: %s", url)
87         response = await task(self, response.headers['Location'], client)
88
89     if response.status_code != codes.OK:
90         Scraper.logger.error("Response error on url: %s, code %s", url,
91                               ↪ response.status_code)
92         return None
93
94     return response
95
96     return wrapper
97
98 @_handle_response
99 async def _scrape(self, url: str, client: AsyncClient) -> Response:
100     async with self.throttler:
101         try:
102             return await client.get(url)
103         except Exception as e:
104             return None
105
106 async def run(self):
107     await self._load_scraper_state()
108     async with AsyncClient() as client:
109         task = partial(self._scrape, client=client)
110         self.queue.set_task(task)
111         for tasks in self.queue:
112             results: Tuple = await asyncio.gather(*tasks)
113
114             if len(self.visited) > self.settings.max_scraped_count:
115                 break
116
117             enriched: List[List[str]] = list(map(lambda r: enrich_text(get_text(r)),
118                                                 ↪ results))
119             await save_documents(enriched, results, Scraper.logger)
120             await save_words(enriched, Scraper.logger)
121
122             processed: Set[str] = self._filter_urls({result.url.path for result in
123                                                 ↪ results})
124             await dump_visited_urls(processed, Scraper.logger)
125             self.visited |= processed

```



```

124         found_refs: Set[str] = self._filter_urls(self._find_hrefs(results))
125         await dump_pending_urls(found_refs, Scraper.logger)
126         self.queue.extend(found_refs)
127
128         Scraper.logger.info("Scraped: %s", len(self.visited))
129
130         await save_bigrams()

```

Лабораторная работа №2 «Булев индекс»

Для реализации программы построения прямого и обратного индексов, был выбран Rust. Прямой индекс в свою очередь разделен на файл-заголовок с мета-информацией об индексе, и на файл с текстом и значениями tf для слов текста.

Запись в файле с мета-информацией представляется в следующем виде:

```

< Doc_id 12 байт >< отступ в файле с текстом u64 >
...

```

Запись в файле с текстом представляется в следующем виде:

```

< Количество слов в документе u64 > [
< размер слова в байтах u64 >< слово >< значение tf f64 >,
...
]

```

Обратный индекс также, как и прямой, разбит на два файла.

Запись в файле с мета-информацией представляется в следующем виде:

```

< размер слова в байтах u64 >< слово >
< отступ в файле с документами u64 >< значение idf f64 >
...

```

Запись в файле с документами представляется в следующем виде:

```
< размер слова в байтах u64 >
< слово >
< Количество документов для данного слова u64 > [
< Doc_id 12 байт >,
...
]
```

Для построения обратного индекса во внешней памяти будем использовать SPIMI алгоритм.

3 Исходный код

```
inverted_index.rs
1 pub async fn build(&mut self) -> Result<(), Box<dyn std::error::Error>> {
2     let dbase = db::connect_to_docs_database().await?;
3     let mut cur = dbase.get_cursor().await?;
4
5     let mut inverted_idx_block = BTreeMap::new();
6
7     let mut tasks = Vec::new();
8
9     while let Some(result) = cur.next().await {
10         let doc: db::Doc = bson::from_document(result)?;
11         let mut used = HashSet::new();
12         for word in doc.words {
13             used.insert(word.to_owned());
14             if !inverted_idx_block.contains_key(&word) {
15                 inverted_idx_block.insert(word, vec![doc._id]);
16             } else {
17                 if let Some(posting_list) = inverted_idx_block.get_mut(&word) {
18                     posting_list.push(doc._id);
19                 }
20             }
21
22             self.cur_map_size += 1;
23
24             if self.cur_map_size >= self.max_map_size {
25                 print!("Block {:?} dumped\n", self.count_of_written_blocks);
```

```

26         tasks.push(InvertedIndex::write_block_to_disk(
27             inverted_idx_block,
28             self.count_of_written_blocks,
29             self.blocks_directory,
30         ));
31         self.count_of_written_blocks += 1;
32         inverted_idx_block = BTreeMap::new();
33         self.cur_map_size = 0;
34     }
35 }
36
37     for word in used {
38         self.head
39             .entry(word)
40             .and_modify(|pair| *pair = (pair.0, pair.1 + 1.0))
41             .or_insert((0, 1.0));
42     }
43 }
44
45 if inverted_idx_block.len() > 0 {
46     print!("EXTRA Block {:?} dumped\n", self.count_of_written_blocks);
47     tasks.push(InvertedIndex::write_block_to_disk(
48         inverted_idx_block,
49         self.count_of_written_blocks,
50         self.blocks_directory,
51     ));
52     self.count_of_written_blocks += 1;
53 }
54
55 let results = join_all(tasks).await;
56 for res in results {
57     res?;
58 }
59 print!("Joined\n");
60
61 self.merge_blocks().await?;
62
63 let mut resulting_block_file = self.blocks_directory.to_owned();
64 resulting_block_file.push_str(&format!(
65     "/block_{}.inv_idx.bin",
66     self.count_of_written_blocks - 1
67 ));
68 copy(resulting_block_file, self.index_content_file).await?;
69 print!("Inverted index's content file was built\n");

```

```

70
71     let count_of_documents = dbase.get_count_of_documents().await?;
72     for pair in self.head.values_mut() {
73         *pair = (pair.0, f64::log10(count_of_documents as f64 / pair.1));
74     }
75
76     self.set_offsets().await?;
77     self.write_head_to_disk().await?;
78     print!("Inverted index was built\n");
79
80     Ok(())
81 }

```

Лабораторная работа №3 «Булев поиск»

Написан web-сервис на Flask, реализовано общение с поисковым движком при помощи httpx. Также в рамках данной работы реализовано исправление ошибок при помощи биграмм и алгоритма Дамерау-Левенштейна.

Среднее время запроса: 10 сек.

4 Исходный код

```

1  from flask import Flask, render_template, request, url_for, flash, redirect
2
3  from search_helper import BigramIndex
4  from request_enrich import get_enriched_words
5  from db import get_documents
6
7  from utils import Settings
8  from typing import List, Dict
9
10 import coloredlogs
11 import logging
12 import asyncio
13 import httpx
14 import time
15
16

```

```

17 loop = asyncio.get_event_loop()
18
19 logger = logging.getLogger(__name__)
20 logger.setLevel(level=logging.DEBUG)
21 coloredlogs.install(level=logging.DEBUG)
22
23 app = Flask(__name__)
24 app.config['SECRET_KEY'] = Settings().secret
25
26
27 @app.route('/', methods=('GET', 'POST'))
28 def search():
29     if request.method == 'POST':
30         search_request = request.form['request']
31         enriched_request: List[str] = get_enriched_words(search_request)
32
33         bi: BigramIndex = BigramIndex(enriched_request)
34         loop.run_until_complete(bi.build(logger))
35
36         search_dict: Dict = bi.get_search_dict()
37         logger.debug("Supposed request structure: %s", search_dict)
38
39         search_engine_request: List = [supposed[0] for supposed in
40             ↪ search_dict.values()]
41         logger.debug("Request for search engine: %s", search_engine_request)
42
43         start = time.time()
44         search_engine_response = httpx.post('http://localhost:8080/search',
45             json={'words': search_engine_request},
46             timeout=None)
47         end = time.time()
48         logger.debug("Search engine request time [sec]: %s", end - start)
49
50         doc_ids: List = search_engine_response.json()["doc_ids"]
51         logger.debug("Got %s documents in search engine", len(doc_ids))
52         results = loop.run_until_complete(get_documents(doc_ids))
53         logger.debug("Response: %s", results)
54
55         front_request: List[Dict] = []
56         for given, changed in zip(enriched_request, search_engine_request):
57             if given != changed:
58                 front_request.append({"word": changed, "color": "red"})
59             else:
60                 front_request.append({"word": changed, "color": "black"})

```

```
60
61     if not search_request:
62         flash('Request is required!')
63     else:
64         return render_template('result.html', results=results,
65                                ↪ changed_request=front_request)
66
67 return render_template('search.html')
68
69 if __name__ == '__main__':
70     app.run(debug=False, use_reloader=False)
```

Лабораторная работа №4 «Ранжирование TF-IDF»

Во время построения прямого индекса считается tf , idf считается при построении обратного индекса. Для запроса также считаются значения tf , потом по полученным значениям строятся вектор запроса и векторы найденных документов. Полученные векторы нормируем и считаем косинусы. Для получения ранжированной выборки сортируем по значению косинуса (по убыванию).

5 Исходный код

```
engine.rs
1  impl Engine {
2      pub async fn search(
3          &self,
4          query: Vec<String>,
5      ) -> Result<Vec<ObjectId>, Box<dyn std::error::Error>> {
6          let qvec = self.query_into_vec(&query);
7
8          let mut result_of_bool_search: Option<Vec<Vec<u8>>> = None;
9
10         for word in &query {
11             let exists = self.inverted_index.head.get(word);
12             match exists {
13                 Some(_) => {}
14                 None => {
15                     continue;
16                 }
17             }
18             match result_of_bool_search {
19                 Some(intersec) => {
20                 ↪ word).await?);
21                     result_of_bool_search = Some(self.intersec(&intersec,
22                                     None => {
23                     result_of_bool_search = Some(self.get_posting_list(word).await?);
24                                     }
25                                 }
26             }
27
28             match result_of_bool_search {
29                 Some(documents) => {
30                     return Ok(self.get_best_documents(&qvec, documents).await?);
31                 }
32                 None => return Ok(vec![]),
33             }
34         }
35
36         async fn get_best_documents(
37             &self,
38             qvec: &HashMap<String, f64>,
39             documents: Vec<Vec<u8>>,
40         ) -> Result<Vec<ObjectId>, Box<dyn std::error::Error>> {
```

```

41     let mut result = Vec::default();
42     let mut heap = BinaryHeap::new();
43
44     for doc in documents {
45         let dvec = self.doc_into_vec(&doc).await?;
46         let mut cosine = 0.0;
47
48         for word in qvec.keys() {
49             if dvec.contains_key(word) {
50                 cosine += qvec[word] * dvec[word];
51             }
52         }
53
54         heap.push((NotNan::new(cosine).unwrap(), doc));
55     }
56
57     for _ in 0..self.top_n {
58         if let Some( (_, doc_id)) = heap.pop() {
59             result.push(ObjectId::from_bytes(doc_id.try_into().unwrap()));
60         } else {
61             break;
62         }
63     }
64
65     Ok(result)
66 }
67
68 async fn intersec(
69     &self,
70     curr: &Vec<Vec<u8>>,
71     word: &String,
72 ) -> Result<Vec<Vec<u8>>, Box<dyn std::error::Error>> {
73     let mut result = Vec::default();
74
75     let other = self.get_posting_list(word).await?;
76
77     let (mut i, mut j) = (0, 0);
78
79     while i < curr.len() && j < other.len() {
80         if curr[i] == other[j] {
81             result.push(curr[i].clone());
82             i += 1;
83             j += 1;
84         } else if curr[i] < other[j] {

```



```

85         i += 1;
86     } else {
87         j += 1;
88     }
89 }
90
91     Ok(result)
92 }
93
94 async fn get_posting_list(
95     &self,
96     word: &String,
97 ) -> Result<Vec<Vec<u8>>, Box<dyn std::error::Error>> {
98     let mut result = Vec::default();
99
100     let offset = self.inverted_index.head[word].0;
101
102     let content_file = File::open(self.inverted_index.index_content_file).await?;
103     let mut reader = BufReader::new(content_file);
104     reader.seek(SeekFrom::Start(offset)).await?;
105
106     let posting_list_len = reader.read_u64().await?;
107
108     for _ in 0..posting_list_len {
109         let mut bytes = vec![0u8; 12];
110         reader.read_exact(&mut bytes).await?;
111         result.push(bytes);
112     }
113
114     Ok(result)
115 }
116
117 fn query_into_vec(&self, query: &Vec<String>) -> HashMap<String, f64> {
118     let mut tf = HashMap::new();
119     for word in query {
120         tf.entry(word.to_owned())
121             .and_modify(|counter| *counter += 1.0)
122             .or_insert(1.0);
123     }
124     for val in tf.values_mut() {
125         *val = 1.0 + f64::log10(*val);
126     }
127
128     let keys = tf.keys().cloned().collect::<Vec<_>>();

```

```

129
130     for word in keys {
131         let default = (0, 0.0);
132         let idf = self.inverted_index.head.get(&word).unwrap_or(&default).1;
133         tf.entry(word).and_modify(|tf| *tf *= idf);
134     }
135
136     let len = f64::sqrt(tf.values().map(|&val| val * val).sum());
137     for val in tf.values_mut() {
138         *val /= len;
139     }
140
141     tf
142 }
143
144 async fn doc_into_vec(
145     &self,
146     doc_id: &Vec<u8>,
147 ) -> Result<HashMap<String, f64>, Box<dyn std::error::Error>> {
148     let mut tf = HashMap::new();
149
150     let offset = self.forward_index.head[doc_id];
151
152     let content_file = File::open(self.forward_index.index_content_file).await?;
153     let mut reader = BufReader::new(content_file);
154     reader.seek(SeekFrom::Start(offset)).await?;
155
156     let words_count = reader.read_u64().await?;
157
158     for _ in 0..words_count {
159         let len = reader.read_u64().await?;
160         let mut bytes = vec![0u8; len as usize];
161         reader.read_exact(&mut bytes).await?;
162         let tf_val = reader.read_f64().await?;
163         tf.insert(String::from_utf8(bytes).unwrap(), tf_val);
164     }
165
166     let keys = tf.keys().cloned().collect::<Vec<_>>();
167
168     for word in keys {
169         let default = (0, 0.0);
170         let idf = self.inverted_index.head.get(&word).unwrap_or(&default).1;
171         tf.entry(word).and_modify(|tf| *tf *= idf);
172     }

```

```

173
174         let len = f64::sqrt(tf.values().map(|&val| val * val).sum());
175         for val in tf.values_mut() {
176             *val /= len;
177         }
178
179         Ok(tf)
180     }
181 }

```

main.rs

```

1  #[macro_use]
2  extern crate lazy_static;
3  extern crate futures;
4
5  lazy_static! {
6      static ref ENGINE: Mutex<Option<engine::Engine>> = Mutex::new(None);
7  }
8
9  pub mod db;
10 pub mod engine;
11 pub mod index;
12 pub mod inverted_index;
13
14 use actix_web::{post, App, HttpResponse, HttpServer, Responder};
15 use serde::{Deserialize, Serialize};
16 use std::sync::Mutex;
17
18 #[derive(Debug, Serialize, Deserialize)]
19 struct Request {
20     words: Vec<String>,
21 }
22
23 #[derive(Debug, Serialize, Deserialize)]
24 struct Response {
25     doc_ids: Vec<[u8; 12]>,
26 }
27
28 #[post("/search")]
29 async fn search(req_body: String) -> impl Responder {
30     println!("Request: {req_body}");
31     let data: Request = serde_json::from_str(&req_body).unwrap();
32     let mut result =
↪ ENGINE.lock().unwrap().as_mut().unwrap().search(data.words).await.unwrap();

```

```

33     result.dedup();
34
35     print!("Found doc_ids:\n");
36     for oid in &result {
37         print!("{:?}", oid.to_hex());
38     }
39
40     let response = Response {
41         doc_ids: result.iter().map(|oid| oid.bytes()).collect(),
42     };
43
44     return HttpResponse::Ok().body(serde_json::to_string(&response).unwrap());
45 }
46
47 #[tokio::main]
48 async fn main() -> std::io::Result<()> {
49     *ENGINE.lock().unwrap() = Some(engine::init_engine().await.unwrap());
50     HttpServer::new(|| App::new().service(search))
51         .bind(("localhost", 8080))?
52         .run()
53         .await
54 }
55

```

6 Выводы

Выполнив первую лабораторную работу по курсу «Информационный поиск», я научился писать веб-скрапер на Python, обрабатывать документы и познакомился с MongoDB. Во второй лабораторной работе научился строить прямой и обратный индексы на Rust для поиска. При написании третьей лабораторной работы я изучил фреймворки Flask и httpx и написал с их помощью веб-сервис для поисковой выдачи и общения с поисковым движком. В четвертой лабораторной я изучил алгоритмы ранжирования, а также нечеткий поиск.

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.))