

Московский авиационный институт  
(национальный исследовательский университет)

Факультет информационных технологий и прикладной  
математики

Кафедра вычислительной математики и программирования

Лабораторная работа по курсу «Информационный поиск»

Студент: И. Д. Недосеков  
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 робота для обхода сайта [neolurk.org](http://neolurk.org). Начиная с главной страницы, скачиваем страницу в формате html. На каждой страничке находим все ссылки, фильтруем (есть пути запрещенные сайтом и ссылки ведущие на другие сайты). После чего вытаскиваем только видимый текст статьи и передаем его в [mystem](#). Эта библиотека производит лексический анализ и выдает значимые слова в приведенной форме. Сохраняем полученный текст и исходный документ в mongodb.

Количество документов: 20000

Размер текста: 408 MB

Средний размер документа: 21 KB

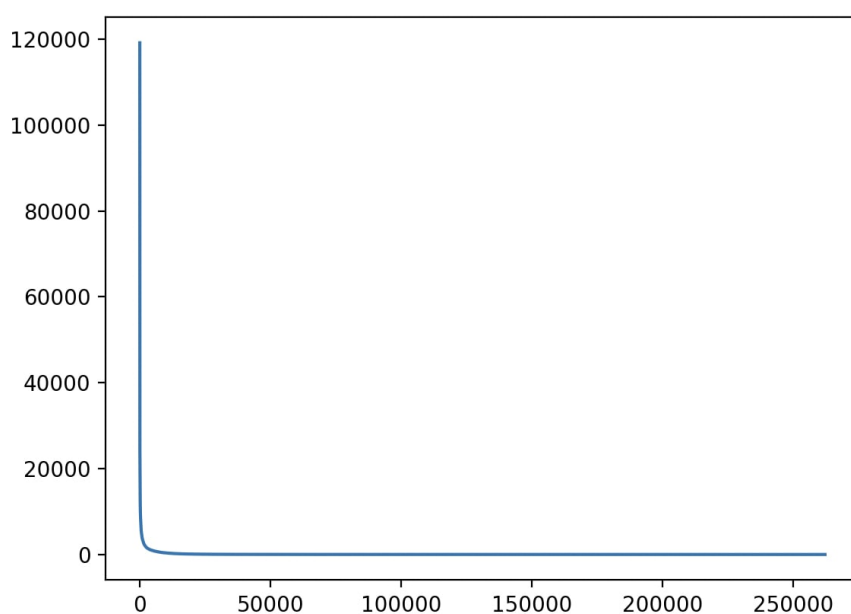


Рис. 1: Закон Ципфа для корпуса

## Лабораторная работа №2 «Булев индекс»

Для реализации программы построения индекса, был выбран Rust. Индекс разделяется на прямой (doc\_id: слова) и обратный (слово: doc\_ids). Прямой индекс в свою очередь разделен на файл мета информации об индексе, и сам индекс. Запись в нем представляется в следующем виде:

```
< Doc_id 12 байт >< Количество слов в документе uint64 > [  
  < размер слова в байтах uint64 >< слово в байтах >,  
  ...  
]
```

В файле мета информации хранится ID документа, сдвиг от начала и размер записи в байтах.

Обратный индекс хранится аналогично прямому.

Для построения будем использовать SPIMInvert алгоритм.

## Лабораторная работа №3 «Булев поиск»

Написан web интерфейс на FastApi.

На этой стороне проверяется корректность запроса и заменяются опечатки через биграммы.

Среднее время запроса: 6.2 сек.

## Лабораторная работа №4 «Ранжирование TF-IDF»

Во время построения прямого индекса считается tf, idf считается после нахождения подходящих документов. Строится вектор запроса и найденных документов. Сортируем по величине косинуса и выдаем результат.

## 2 Исходный код

```
engine/index_api/forward_index.rs
1  use std::{
2      collections::HashMap,
3      io::{Read, Seek, Write},
4      mem::size_of,
5      ops::AddAssign,
6  };
7
8  lazy_static! {
9      static ref FORWARD_INDEX_RECORD_METADATA_SIZE: usize =
10      bincode::serialized_size(&ForwardIndexRecordMetadata::default()).unwrap() as
11      ↪ usize;
12  }
13
14  use futures::TryStreamExt;
15  use mongodb::bson::Document;
16
17  use std::fs::OpenOptions;
18
19  use serde::{de::IntoDeserializer, Deserialize, Serialize};
20
21  #[derive(Serialize, Deserialize, PartialEq, Debug)]
22  pub struct ForwardIndexRecord {
23      pub id: [u8; 12],
24      pub lemmas: Vec<Lemma>,
25  }
26
27  #[derive(Serialize, Deserialize, PartialEq, Debug)]
28  pub struct Lemma {
29      pub frequency: f64,
30      pub text: String,
31  }
32
33  impl ForwardIndexRecord {
34      fn from_document(document: Document) -> Self {
35          let lemmas: Vec<String> = document
36              .get_array("lemmas")
37              .unwrap()
38              .into_iter()
39              .map(|lemma_info| {
40                  lemma_info
41                      .as_document()
```

```

41         .unwrap()
42         .get_str("text")
43         .unwrap()
44         .to_string()
45     })
46     .collect();
47
48     let mut lemmas_count = HashMap::new();
49     let count_words = lemmas.len() as f64;
50     for lemma in lemmas {
51         lemmas_count.entry(lemma).or_insert(0).add_assign(1);
52     }
53
54     return Self {
55         id: document.get_object_id("_id").unwrap().bytes(),
56         lemmas: lemmas_count
57             .into_iter()
58             .map(|(lemma, count)| Lemma {
59                 frequency: count as f64 / count_words,
60                 text: lemma,
61             })
62             .collect(),
63     };
64 }
65 }
66
67 #[derive(Serialize, Deserialize, Debug, Default)]
68 pub struct ForwardIndexRecordMetadata {
69     id: [u8; 12],
70     offset_in_index_file: u64,
71     size_of_record: u64,
72 }
73
74 pub async fn build_forward_index(cursor: &mut mongodb::Cursor<Document>) {
75     let mut index_file = OpenOptions::new()
76         .truncate(true)
77         .create(true)
78         .write(true)
79         .open("./forward_index.bin")
80         .unwrap();
81
82     let mut index_metadata_file = OpenOptions::new()
83         .truncate(true)
84         .create(true)

```



```

85     .write(true)
86     .open("./forward_index_metadata.bin")
87     .unwrap();
88
89     while let Some(document) = cursor.try_next().await.unwrap() {
90         let record = ForwardIndexRecord::from_document(document);
91         let record_as_bytes = bincode::serialize(&record).unwrap();
92
93         let current_metadata = ForwardIndexRecordMetadata {
94             id: record.id,
95             offset_in_index_file: index_file.metadata().unwrap().len(),
96             size_of_record: record_as_bytes.len() as u64,
97         };
98
99         index_file.write(&record_as_bytes[..]).unwrap();
100
101         let record_as_bytes = bincode::serialize(&current_metadata).unwrap();
102         index_metadata_file.write(&record_as_bytes[..]).unwrap();
103     }
104 }
105
106 use std::fs::File;
107 use std::io::SeekFrom;
108
109 type DocId = [u8; 12];
110 pub struct ForwardIndexApi {
111     index_file: File,
112     index_metadata: File,
113
114     loaded_metadata: Option<HashMap<DocId, ForwardIndexMetaRecord>>,
115 }
116
117 struct ForwardIndexMetaRecord {
118     offset_in_index_file: u64,
119     size_of_record: u64,
120 }
121
122 impl ForwardIndexApi {
123     pub fn new(index_path: &str, index_metadata_path: &str) -> Self {
124         Self {
125             index_file: File::open(&index_path).unwrap(),
126             index_metadata: File::open(&index_metadata_path).unwrap(),
127
128             loaded_metadata: None,

```

```

129     }
130 }
131
132 pub fn next(&mut self) -> Result<ForwardIndexRecord, Box<dyn std::error::Error>> {
133     let mut metadata_buf = vec![0; *FORWARD_INDEX_RECORD_METADATA_SIZE];
134     self.index_metadata.read_exact(&mut metadata_buf)?;
135
136     let next_record_info: ForwardIndexRecordMetadata =
137         bincode::deserialize(&metadata_buf).unwrap();
138
139     let mut record_buf = vec![0u8; next_record_info.size_of_record as usize];
140     self.index_file.read_exact(&mut record_buf.as_mut_slice())?;
141
142     let next_record: ForwardIndexRecord =
143     ↪ bincode::deserialize(&record_buf.as_slice()).unwrap();
144
145     return Ok(next_record);
146 }
147
148 fn _load_metadata_record(
149     &mut self,
150 ) -> Result<ForwardIndexRecordMetadata, Box<dyn std::error::Error>> {
151     let mut metadata_buf = vec![0; *FORWARD_INDEX_RECORD_METADATA_SIZE];
152     self.index_metadata.read_exact(&mut metadata_buf)?;
153     let next_record_info: ForwardIndexRecordMetadata =
154     ↪ bincode::deserialize(&metadata_buf)?;
155
156     return Ok(next_record_info);
157 }
158
159 pub fn load_metadata(&mut self) {
160     if let Some(_) = self.loaded_metadata {
161         return;
162     }
163
164     let mut metadata = HashMap::new();
165
166     while let Ok(next_record) = self._load_metadata_record() {
167         metadata.insert(
168             next_record.id,
169             ForwardIndexMetaRecord {
170                 offset_in_index_file: next_record.offset_in_index_file,
171                 size_of_record: next_record.size_of_record,
172             },
173         );
174     }
175 }

```

```

171     );
172 }
173 self.loaded_metadata = Some(metadata);
174 }
175
176 pub fn load_record_by(&mut self, id: &DocId) -> ForwardIndexRecord {
177     self.load_metadata();
178     let loaded_meta = self.loaded_metadata.as_ref().unwrap();
179     let meta_info = loaded_meta.get(id).unwrap();
180     self.index_file
181         .seek(SeekFrom::Start(meta_info.offset_in_index_file))
182         .unwrap();
183     let mut buf = vec![0u8; meta_info.size_of_record as usize];
184     self.index_file.read_exact(&mut buf).unwrap();
185     let record: ForwardIndexRecord = bincode::deserialize(&buf).unwrap();
186
187     return record;
188 }
189 }
190

```

engine/index\_api/reversed\_index.rs

```

1  #[macro_use]
2  extern crate lazy_static;
3
4  pub mod forward_index;
5  pub mod reversed_index;
6

```

engine/index\_api/lib.rs

```

1  use std::{
2      collections::{BTreeMap, HashMap},
3      io::{Read, Write},
4  };
5
6  use crate::forward_index;
7
8  pub async fn spimi_invert(
9      stream: &mut forward_index::ForwardIndexApi,
10 ) -> Result<(), Box<dyn std::error::Error>> {
11     let mut dictionary = BTreeMap::new();
12     let mut count_documents = 0;
13     while let Ok(docs) = stream.next() {

```

```

14     println!("{count_documents}");
15     count_documents += 1;
16     for lemma in docs.lemmas {
17         dictionary
18             .entry(lemma.text)
19             .or_insert(Vec::default())
20             .push(docs.id);
21     }
22
23 }
24
25 ReversedIndexAPI::dump_to_file(dictionary, count_documents)?;
26
27 Ok(())
28 }
29
30 #[derive(serde::Serialize, serde::Deserialize)]
31 struct ReversedIndexRecord {
32     word: String,
33     doc_ids: Vec<[u8; 12]>,
34     idf: f64,
35 }
36
37 #[derive(serde::Serialize, serde::Deserialize)]
38 struct ReversedIndex {
39     index: Vec<ReversedIndexRecord>,
40 }
41
42 pub struct ReversedIndexAPI {
43     file: std::fs::File,
44 }
45
46 impl ReversedIndexAPI {
47     fn new(path: &str) -> Self {
48         Self {
49             file: std::fs::File::open(&path).unwrap(),
50         }
51     }
52
53     fn dump_to_file(
54         dictionary: BTreeMap<String, Vec<[u8; 12]>>,
55         count_documents: u64,
56     ) -> Result<(), Box<dyn std::error::Error>> {
57         let mut file = std::fs::OpenOptions::new()

```

```

58 .truncate(true)
59 .create_new(true)
60 .write(true)
61 .open("./reversed_index.bin")?;
62
63 let index = dictionary
64 .into_iter()
65 .map(|(word, doc_ids)| ReversedIndexRecord {
66     word: word,
67     idf: f64::ln(count_documents as f64 / doc_ids.len() as f64),
68     doc_ids: doc_ids,
69 })
70 .collect();
71
72 let buf = bincode::serialize(&ReversedIndex { index: index })?;
73
74 file.write(&buf.as_slice())?;
75
76 Ok(())
77 }
78
79 pub fn load_from_file(path: &str) -> (BTreeMap<String, Vec<[u8; 12]>>,
80 ↪ HashMap<String, f64>) {
81     let mut file = std::fs::File::open(path).unwrap();
82     let buf_size = file.metadata().unwrap().len();
83
84     let mut buf = vec![0u8; buf_size as usize];
85
86     file.read_exact(&mut buf.as_mut_slice()).unwrap();
87     let index: ReversedIndex = bincode::deserialize(&buf).unwrap();
88
89     let words_idf = HashMap::from_iter(
90     index
91     .index
92     .iter()
93     .map(|record| (record.word.clone(), record.idf)),
94     );
95
96     return (BTreeMap::from_iter(
97     index
98     .index
99     .into_iter()
100     .map(|record| (record.word, record.doc_ids)),
101     ), words_idf);

```

```
101 }
102 }
```

```
engine/indexer/main.rs
1  use serde::{Deserialize, Serialize};
2
3  use mongodb::bson::doc;
4  use mongodb::bson::oid::ObjectId;
5
6  #[derive(Debug, Serialize, Deserialize)]
7  struct Page {
8      _id: ObjectId,
9      path: String,
10     lemmas: Vec<Lemma>,
11 }
12
13 use index_api::forward_index::*;
14 use index_api::reversed_index::*;
15
16 #[derive(Debug, Serialize, Deserialize)]
17 struct Lemma {
18     text: String,
19 }
20
21 #[tokio::main]
22 async fn main() -> Result<(), Box<dyn std::error::Error>> {
23     let client = mongodb::Client::with_uri_str("mongodb://localhost").await?;
24
25     let db = client.database("test");
26
27     let test_collection = db.collection::<Page>("docs");
28
29     let aggregate_pipeline = vec![doc! {
30         "$sort": {
31             "_id": 1
32         }
33     }];
34     let mut cursor = test_collection.aggregate(aggregate_pipeline, None).await?;
35
36     build_forward_index(&mut cursor).await;
37
38     let mut fapi = ForwardIndexApi::new("./forward_index.bin",
    ↪  "./forward_index_metadata.bin");
```

```

39
40     spimi_invert(&mut fapi).await?;
41
42     Ok(())
43 }
44

```

```

                                     engine/server/boolean_search.rs
1  use std::collections::BTreeMap;
2
3  // use futures::stream::Iter;
4
5  use std::cmp::Ordering;
6
7  #[derive(Debug)]
8  struct SearchState<'a> {
9      ids: &'a Vec<u8; 12>,
10     position: usize,
11 }
12
13 pub fn boolean_search(
14     request: &Vec<String>,
15     reversed_index: &BTreeMap<String, Vec<u8; 12>>,
16 ) -> Vec<u8; 12> {
17     let mut docs_containing_words = Vec::with_capacity(request.len());
18     let mut maximum_doc_id = [u8::MIN; 12];
19
20     for word in request {
21         match reversed_index.get(word) {
22             Some(docs) => {
23                 docs_containing_words.push(SearchState {
24                     ids: docs,
25                     position: 0,
26                 });
27                 maximum_doc_id = std::cmp::max(maximum_doc_id, docs[0])
28             }
29             None => continue,
30         }
31     }
32     let mut result = Vec::new();
33
34     'main_loop: loop {
35         let mut some_checked = false;

```

```

36     let mut all_previous_equal = true;
37
38     'word_cycle: for state in &mut docs_containing_words {
39         'doc_id_cycle: while let Some(next_doc_id) = state.ids.get(state.position) {
40             some_checked = some_checked || true;
41
42             match next_doc_id.cmp(&maximum_doc_id) {
43                 Ordering::Less => {
44                     state.position += 1;
45                     continue 'doc_id_cycle;
46                 }
47                 Ordering::Equal => {
48                     all_previous_equal = all_previous_equal && true;
49                     break 'doc_id_cycle;
50                 }
51                 Ordering::Greater => {
52                     all_previous_equal = false;
53                     maximum_doc_id = next_doc_id.clone();
54                     continue 'main_loop;
55                 }
56             }
57         }
58     }
59     if all_previous_equal {
60         result.push(maximum_doc_id);
61         docs_containing_words
62             .iter_mut()
63             .for_each(|state| state.position += 1)
64     }
65
66     if !some_checked {
67         break 'main_loop;
68     }
69 }
70
71 return result;
72 }
73

```

---

engine/server/main.rs

---

```

1     #[macro_use]
2     extern crate lazy_static;
3

```



```

4 lazy_static! {
5     static ref REVERSED_INDEX: BTreeMap<String, Vec<[u8; 12]>> =
6         ReversedIndexAPI::load_from_file("../indexer/reversed_index.bin").0;
7     static ref WORDS_IDF: HashMap<String, f64> =
8         ReversedIndexAPI::load_from_file("../indexer/reversed_index.bin").1;
9 }
10
11 use std::collections::{BTreeMap, HashMap};
12
13 use index_api::{forward_index, reversed_index::ReversedIndexAPI};
14 use serde::{Deserialize, Serialize};
15
16 pub mod boolean_search;
17
18 #[derive(Serialize, Deserialize, Debug)]
19 struct RequestSchema {
20     words: Vec<String>,
21 }
22
23 #[derive(Serialize, Deserialize, Debug)]
24 struct ResponseSchema {
25     doc_ids: Vec<mongodb::bson::oid::ObjectId>,
26 }
27
28 use actix_web::{get, post, App, HttpResponse, HttpServer, Responder};
29
30 #[get("/")]
31 async fn hello() -> impl Responder {
32     HttpResponse::Ok().body("Hello world!")
33 }
34
35 use itertools::Itertools;
36 use std::cmp::Ordering;
37
38
39 fn calculate_cos(v1: &Vec<f64>, v2: &Vec<f64>) -> f64 {
40     let numerator = v1.iter().zip(v2.iter()).map(|(a, b)| a * b).sum::<f64>();
41     let denominator =
42         f64::sqrt(v1.iter().map(|a| a * a).sum()) * f64::sqrt(v2.iter().map(|b| b *
↪ b).sum());
43
44     return numerator / denominator;
45 }
46

```

```

47 fn rank_docs(request: &Vec<String>, docs: &Vec<[u8; 12]>) -> Vec<[u8; 12]> {
48     let mut forward_api = forward_index::ForwardIndexApi::new(
49         "../indexer/forward_index.bin",
50         "../indexer/forward_index_metadata.bin",
51     );
52
53     let request_words_count = request.len() as f64;
54     let request_vector = request
55         .iter()
56         .map(|word| WORDS_IDF.get(word).unwrap_or(&0f64) * 1f64 / request_words_count)
57         .collect_vec();
58
59     let result = docs.iter().cloned().sorted_by_key(|doc| {
60         let forward_record = forward_api.load_record_by(doc);
61         let searching_words: Vec<forward_index::Lemma> = forward_record
62             .lemmas
63             .into_iter()
64             .filter(|lemma| request.contains(&lemma.text))
65             .sorted_by_key(|lemma| lemma.text.clone())
66             .collect_vec();
67
68         let document_vector = searching_words
69             .iter()
70             .map(|lemma| lemma.frequency * WORDS_IDF.get(&lemma.text).unwrap())
71             .collect();
72
73         return (calculate_cos(&request_vector, &document_vector) * 1000f64) as i64;
74     })
75     .rev().collect();
76
77     return result;
78 }
79
80 #[post("/search")]
81 async fn search(req_body: String) -> impl Responder {
82     println!("{req_body}");
83     let data: RequestSchema = serde_json::from_str(&req_body).unwrap();
84     let result_ids = boolean_search::boolean_search(&data.words, &REVERSED_INDEX);
85     let ranked_results = rank_docs(&data.words, &result_ids);
86     let result = ResponseSchema {
87         doc_ids: ranked_results.into_iter().map(|doc_id|
88     ↪ mongodb::bson::oid::ObjectId::from_bytes(doc_id)).collect(),
89     };

```

```

90     return HttpResponse::Ok().body(serde_json::to_string(&result).unwrap());
91 }
92
93 #[tokio::main]
94 async fn main() -> std::io::Result<()> {
95     HttpServer::new(|| App::new().service(hello).service(search))
96     .bind(("localhost", 8080))?
97     .run()
98     .await
99 }
100

```

scraper/scrapy.py

```

1  import asyncio
2  from collections import deque
3  from itertools import chain, islice
4  from typing import Iterable, TypeVar
5
6  from aiolimiter import AsyncLimiter
7  from bs4 import BeautifulSoup
8  from httpx import AsyncClient, Response, codes
9  from loguru import logger
10
11  T = TypeVar('T')
12
13
14  class ChankedQueue(deque[tuple[T, ...]]):
15  def __init__(self, chunk_size: int = 50, *arg, **kwargs) -> None:
16  self.chunk_size = chunk_size
17  super(ChankedQueue, self).__init__(*arg, **kwargs)
18
19  def extend(self, iterable: Iterable[T]) -> None: # type: ignore
20  real_dequeue = super(ChankedQueue, self)
21  start = 0
22  if real_dequeue.__bool__():
23  last_chunk = real_dequeue.pop()
24  stop = self.chunk_size - len(last_chunk)
25  else:
26  last_chunk = tuple()
27  stop = self.chunk_size
28
29  while True:
30  next_chunk = tuple(chain(last_chunk, islice(iterable, start, stop)))

```

```

31     if not next_chunk:
32         break
33
34     real_dequeue.append(next_chunk)
35     last_chunk = tuple()
36     start, stop = stop, stop + self.chunk_size
37
38
39     def get_all_url_from_html(source: str) -> set[str]:
40         soup = BeautifulSoup(source, features="html.parser")
41         return {a['href'] for a in soup.find_all('a', href=True)}
42
43
44     async def make_request(url: str, session: AsyncClient, throttler: AsyncLimiter) ->
45         ↳ Response:
46     async with throttler:
47         logger.debug('scrapy {}'.format(url))
48         return await session.get(url)
49
50     async def scrape(url: str, session: AsyncClient, throttler: AsyncLimiter) ->
51         ↳ Response | None:
52     response = await make_request(url, session, throttler)
53
54     if response.status_code == codes.OK:
55         return response
56
57     if response.status_code == codes.TOO_MANY_REQUESTS or response.status_code ==
58         ↳ codes.SERVICE_UNAVAILABLE:
59     logger.warning('retry later by {} on url {}'.format(response.status_code, url))
60     await asyncio.sleep(1.5)
61     response = await make_request(url, session, throttler)
62
63     if codes.is_redirect(response.status_code):
64     logger.warning('redirect on {} {}'.format(url, response.status_code))
65     response = await make_request(response.headers['Location'], session, throttler)
66
67     if response.status_code != codes.OK:
68     logger.error('cannot scrape {} code {}'.format(url, response.status_code))
69     return None
70
71     return response

```

---

 scrapper/scrapy.py
 

---

```

1  import asyncio
2  from itertools import chain
3
4  from httpx import AsyncClient, Timeout
5  from loguru import logger
6
7  from documents import upload_documents
8  from scrapy import ChankedQueue, get_all_url_from_html, scrape
9  from config import Config
10 from postprocess import build_bigram_words_dictionary
11
12 def get_new_urls(paths: set[str], visited: set[str]) -> set[str]:
13     return {
14         f'{Config.url}/{path}'
15         for path in paths
16         if path.startswith('/wiki/') and path not in visited
17     }
18
19
20 async def empty_coroutine() -> None:
21     return
22
23
24 async def run():
25     visited_paths = {f'{Config.url}/{Config.start_path}'}
26     urls_to_process = ChankedQueue()
27     urls_to_process.append((f'{Config.url}/{Config.start_path}',))
28
29     upload_documents_task = empty_coroutine()
30
31     session = AsyncClient(timeout=Timeout(100, connect=60))
32
33     while len(visited_paths) < Config.count_documents:
34
35         url_chunk = urls_to_process.popleft()
36         scrape_tasks = [scrape(url, session=session, throttler=Config.throttler) for url in
37             ↪ url_chunk]
38         _, *results = await asyncio.gather(upload_documents_task, *scrape_tasks)
39         visited_paths.update(url_chunk)

```

```

40 paths_in_results = set(
41     chain.from_iterable(
42         get_all_url_from_html(result.text)
43     for result in results
44     if result
45     )
46 )
47 urls_to_process.extend(get_new_urls(paths_in_results, visited_paths))
48
49 upload_documents_task = upload_documents(results)
50
51 await upload_documents_task
52 await session.aclose()
53 logger.debug('finished downloading, start build dictionary and bigrams')
54
55 await build_bigram_words_dictionary()
56
57
58 logger.debug("finished scraping")
59
60
61
62
63 if __name__ == "__main__":
64     asyncio.run(run())
65

```

```

web-api/main.py
1  import json
2
3  import httpx
4  from fastapi import FastAPI, Request, Form
5  from fastapi.templating import Jinja2Templates
6  from bson import ObjectId
7
8  from helpers.mongodb_connector import mongo_client
9
10 from loguru import logger
11
12 from enrich import get_enrich_query
13
14 app = FastAPI()
15 templates = Jinja2Templates(directory="templates/")

```

```

16
17 EXAMPLE_DOCS = ['http://neolurk.org/wiki/Заглавная_страница',
18 ↪ 'http://neolurk.org/wiki/Заглавная_страница']
19
20 @app.get('/')
21 def read_form():
22     return 'hello world'
23
24
25 @app.get("/form")
26 def form_post(request: Request):
27     result = ""
28     return templates.TemplateResponse('form.html', context={'request': request,
29 ↪ 'result': result})
30
31
32 async def get_search_results(query: list[str]) -> list[str]:
33     ids = await get_search_results_ids(query)
34     ids = [ObjectId(id) for id in ids]
35
36     results = mongo_client.test.docs.aggregate(
37     [
38     {'$match': {'_id': {'$in': ids}}},
39     {'$project': {'path': 1}}
40     ]
41     )
42     res_url = []
43     async for result in results:
44         res_url.append(f'http://neolurk.org{result["path"]}')
45
46     return res_url
47
48
49 async def get_search_results_ids(query: list[str]) -> list[str]:
50     async with httpx.AsyncClient() as client:
51         response = await client.post(
52         'http://localhost:8080/search',
53         headers={'Content-Type': 'application/x-www-form-urlencoded'},
54         content=f'{{"words": { [word for word in sorted(query)] }}}'.replace("'", '"')
55         )
56
57     if response.status_code != httpx.codes.OK:
58         logger.error('something went wrong', response)

```

```

58     raise Exception('asd')
59
60     return [doc['$oid'] for doc in response.json()['doc_ids']]
61
62
63     @app.post("/form")
64     async def form_post(request: Request, statement: str = Form(default='')):
65         result = EXAMPLE_DOCS if statement else []
66         enrich_query = await get_enrich_query(statement)
67         query = [
68             enrich_query.fuzzy_mapping[lemma][0]
69             if lemma in enrich_query.fuzzy_mapping
70             else lemma
71             for lemma in enrich_query.query
72         ]
73         search_result = await get_search_results(query)
74         logger.debug('Enrich query {} mapping {}'.format(enrich_query.query,
75                                                         ↪ enrich_query.fuzzy_mapping))
76         return templates.TemplateResponse('form.html', context={'request': request,
77                                                         ↪ 'results': search_result})

```

```

web-api/fuzzy.py
1  from helpers.mongodb_connector import mongo_client
2  from helpers.bigrams import get_bigrams_of_word
3  from config import Config
4
5
6
7  class Dimension:
8
9  def __init__(self, word: str) -> None:
10     self.word = word
11     self.words_bigrams = set(get_bigrams_of_word(word))
12
13  def distance_to(self, candidate: str) -> float:
14     if abs(len(self.word) - len(candidate)) > 3:
15         return 0
16     words_bigrams = self.words_bigrams
17     candidates_bigrams = set(get_bigrams_of_word(candidate))
18     distance = abs(len(words_bigrams) - len(words_bigrams & candidates_bigrams))
19     return abs(1 - distance / len(words_bigrams | candidates_bigrams))
20

```



```

21
22 async def get_fuzzy_words(word: str) -> list[str]:
23     bigrams = get_bigrams_of_word(word)
24     word_dimension = Dimension(word)
25     fuzzy_words_rating: set[tuple[str, float]] = set()
26     async for bigram_with_words in mongo_client.test.bigram_words.find({'bigram': {'$in':
27         ↪ bigrams}}):
28         fuzzy_words_rating.update(
29             (candidate, word_dimension.distance_to(candidate))
30         for candidate in bigram_with_words['words'])
31         if word_dimension.distance_to(candidate) >= Config.minimum_distance_for_fuzzy_words
32         )
33     return [
34         word
35         for word, _ in reversed(sorted(fuzzy_words_rating, key=lambda row: row[1]))
36     ]
37
38 async def is_dictionary_word(word: str) -> bool:
39     return bool(
40         await mongo_client.test.word_bigrams.count_documents({'word': word})
41     )

```

```

web-api/enrich.py
1 import asyncio
2
3 from helpers.mysystem import system
4 from helpers.stemmer import stemmer
5
6 from typing import NamedTuple
7
8 from fuzzy import is_dictionary_word, get_fuzzy_words
9
10
11 def get_lexemes(text: str) -> list[str]:
12     return [
13         stemmer.stem(lexem)
14         for lexem in system.lemmatize(text)
15         if lexem.strip()
16     ]
17
18
19 class EnrichQuery(NamedTuple):

```

```

20 query: list[str]
21 fuzzy_mapping: dict[str, str]
22
23
24 async def process_lexeme(lexeme: str) -> tuple[str, str]:
25     if await is_dictionary_word(lexeme):
26         return lexeme, lexeme
27     else:
28         return lexeme, await get_fuzzy_words(lexeme) or lexeme
29
30
31 async def get_enrich_query(text: str) -> EnrichQuery:
32     lexemes = get_lexemes(text)
33     query: list[str] = []
34     fuzzy_mapping: dict[str, str] = {}
35     process_lexeme_tasks = (process_lexeme(lexeme) for lexeme in lexemes)
36
37     for lexeme, processed_lexeme in await asyncio.gather(*process_lexeme_tasks):
38         query.append(lexeme)
39         if lexeme != processed_lexeme:
40             fuzzy_mapping[lexeme] = processed_lexeme
41
42     return EnrichQuery(query, fuzzy_mapping)
43
44

```

### 3 Выводы

Выполнив лабораторную работу по курсу «Информационный поиск», я написал веб паука, научился строить индекс для поиска. Изучил алгоритмы ранжирования, лемматизации и основы работы поисковых движков.

## 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.))