

Московский авиационный институт  
(национальный исследовательский университет)

Факультет информационных технологий и прикладной  
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №1 по курсу «Информационный поиск»

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

Москва, 2023

# Лабораторная работа №1

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

- Скачать его к себе на компьютер. В отчете нужно указать источник данных.
- Ознакомиться с ним, изучить его характеристики. Из чего состоит текст? Есть ли дополнительная мета-информаци? Если разметка текста, какая она?
- Разбить на документы.
- Выделить текст.
- Найти существующие поисковики, которые уже можно использовать для поиска по выбранному набору документов.
- Привести несколько примеров запросов к существующим поисковикам, указать недостатки в полученной поисковой выдаче.

В результате работы должна быть указана статистическая информация о корпусе:

- Размер "сырых" данных.
- Количество документов.
- Размер текста, выделенного из "сырых" данных.
- Средний размер документа, средний объем текста в документе.

# 1 Описание

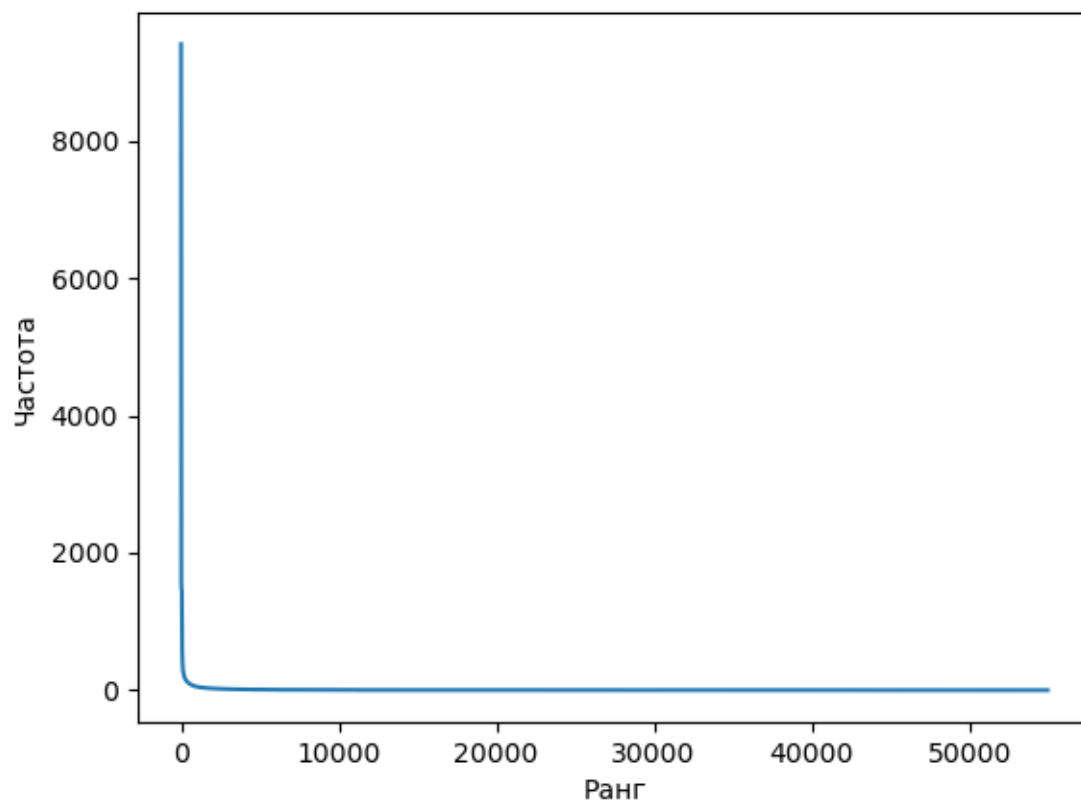
В качестве набора данных я решил выбрать английские статьи, начинающиеся на букву «а», википедии. Для того, чтобы загрузить ее html страницы, я решил воспользоваться общедоступным архивом википедии 2008 года. Открывая каждую страницу, я выделяю в ней заголовок, первый абзац и остальной текст, затем с помощью фреймворка nltk нормализую его. После нормализации я сохраняю текст из каждой зоны в отдельный файл.

Количество документов: 153512

Количество слов в документах: 51543484

Размер полученных документов: 365 MB

Также я построил закон Ципфа для своего набора данных:



## 2 Исходный код

```
1 from bs4 import BeautifulSoup
2 import os
3 from nltk.tokenize import word_tokenize, sent_tokenize
4 from nltk.corpus import stopwords
5 from nltk.stem import PorterStemmer, WordNetLemmatizer
6 import string
7 from tqdm import tqdm
8 import time
9
10 from collections import Counter
11 import matplotlib.pyplot as plt
12
13 lemmatizer = WordNetLemmatizer()
14 all_tokens = []
15
16 stop_words = set(stopwords.words("english"))
17
18 cnt_files = 0
19 cnt_words = 0
20
21 def get_visible_text(soup):
22     #
23     title = soup.find("h1", {"class": "firstHeading"}).text
24
25     #
26     content = soup.find("div", {"class": "mw-parser-output"}) or soup.body
27     first_paragraph = None
28     other_text = ""
29     for el in content.find_all(["p", "table"]):
30         if el.name == "p" and el.text.strip() != "":
31             first_paragraph = el.text
32             break
33
34     for p in content.find_all("p")[1:]:
35         other_text += p.text
36     for p in content.find_all(['ul', 'ol']):
37         other_text += p.text
38     return (title, first_paragraph, other_text)
39
40 def normalize_text(text):
41     words = word_tokenize(text)
42
43     punct = string.punctuation
44     words = [word for word in words if word not in punct and word != '‘’' and word != '
45         \'\' and word != '-']
46     words_lower = [word.lower() for word in words]
```

```

47 words_no_stop = [word for word in words_lower if word not in stop_words]
48 lemmatized_words = [lemmatizer.lemmatize(word) for word in words_no_stop]
49
50 global cnt_words
51 cnt_words += len(lemmatized_words)
52 # all_tokens.extend(lemmatized_words)
53 return lemmatized_words
54
55 def process_file(file_path):
56     try:
57         #
58         with open(file_path, 'r', encoding='utf-8') as f:
59             html_content = f.read()
60
61             soup = BeautifulSoup(html_content, "html.parser")
62
63             # ,
64             redirect = soup.find("meta", {"http-equiv": "Refresh"})
65             if redirect:
66                 return
67
68             (title, first_paragraph, other_text) = get_visible_text(soup)
69
70             #
71             title = normalize_text(title)
72             first_paragraph = normalize_text(first_paragraph)
73             other_text = normalize_text(other_text)
74
75             with open(os.path.join(res_folder, os.path.basename(file_path)), mode='w',
76                       encoding="utf-8") as f:
77                 f.write(' '.join(title) + '\n')
78                 f.write(' '.join(first_paragraph) + '\n')
79                 f.write(' '.join(other_text))
80
81             global cnt_files
82             cnt_files += 1
83
84             if cnt_files % 1000 == 0:
85                 print(cnt_files)
86     except Exception as e:
87         return
88
89 def process_folder(folder_path, res_folder):
90     global cnt_files
91     for root, dirs, files in os.walk(folder_path):
92         for file in files:
93             if not str(os.path.basename(file)).startswith('A') and not str(os.path.
94                                     basename(file)).startswith('a'):
95                 file_path = os.path.join(root, file)

```

```

94         os.remove(file_path)
95         continue
96
97         file_path = os.path.join(root, file)
98         process_file(file_path)
99
100     folder_path = os.path.join('data', 'en', 'a')
101     res_folder = os.path.join('texts', 'en')
102     start = time.time()
103     process_folder(folder_path, res_folder)
104     end = time.time()
105     print('Time: ', int(end - start) // 60, 'm ', int(end - start) % 60, 's')
106     print('Words: ', cnt_words)
107     print('Files: ', cnt_files)
108
109     #
110     # freq = Counter(all_tokens)
111     # most_common_words = freq.most_common()
112     # rank = range(1, len(most_common_words)+1)
113     # freqs = [f for (w, f) in most_common_words]
114     # plt.plot(rank, freqs)
115     # plt.xlabel('Rank')
116     # plt.ylabel('Freq')
117     # plt.show()

```

### 3 Выводы

Выполнив первую лабораторную работу по курсу «Информационный поиск», я научился получать видимый текст с html страниц, а затем обрабатывать его. Именно по полученному обработанному набору данных впоследствии я смогу реализовать поиск по статьям.

## 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.