

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

Москва, 2023

## 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.

# 1 Описание

В качестве корпуса документов был выбран дамп английской википедии. Из этого дампа было выбрано 50000 документов, которые будут использоваться в дальнейшем. Каждый документ имеет следующее строение:

```
<doc id=«<Id>» url=«https://en.wikipedia.org/wiki?curid=<Id>» title=«<Title>»>  
<Title>  
<Contents>  
</doc>
```

Где <Id> - это идентификатор документа в википедии, <Title> - это название документа, <Contents>- это содержимое документа.

## 2 Исходный код

Для того, чтобы получить корпус нужного размера, был использован следующий скрипт на языке Python:

```
1 FILE_NAME = 'corpus/wikipedia'
2 PREVIEW_SIZE = 1557199
3 BUFFER = list()
4 TOTAL = 0
5
6 with open(FILE_NAME) as f:
7     for _ in range(PREVIEW_SIZE):
8         BUFFER.append(f.readline())
9
10 with open(f'{FILE_NAME}_preview', 'w') as f:
11     for i in range(PREVIEW_SIZE):
12         if (BUFFER[i] == "</doc>\n"):
13             TOTAL += 1
14             f.write(BUFFER[i])
15
16 print(TOTAL)
```

### 3 Выводы

Выполнив первую лабораторную работу по курсу «Информационный поиск», я осознал понятие документа и корпуса документов, лучше понял строение википедии и узнал о том, какие метаданные содержатся на ее страницах.

## 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.))