

**Московский авиационный институт (национальный
исследовательский университет)**

Факультет компьютерные науки и прикладной математики

Кафедра прикладная математика и информатика

Лабораторные работы по курсу «Информационный поиск»

Студент: А.А. Серякова

Преподаватель: А. А. Кухтичев

Группа: М8О-409Б-22

Дата:

Оценка:

Подпись:

Москва, 2025

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

Для выполнения лабораторной работы был выбран следующий метод:

1. **Выбор источника данных:** В качестве источника выбран корпус статей русской Википедии из тематической категории "Актёры США". Данный источник удовлетворяет требованиям: единая тематика, большой объём (десятки тысяч статей), наличие существующих поисковых систем для проверки.
2. **Скачивание корпуса:** Разработан набор Python-скриптов, использующих:
 - Библиотеку `wikipedia-api` для получения списка статей категории и загрузки их содержимого
 - Библиотеку `mwparserfromhell` для очистки вики-разметки
 - Библиотеку `tqdm` для визуализации прогресса
3. **Обработка данных:** Каждая статья обрабатывается в два этапа:
 - Сохранение в исходном виде (вики-разметка) в папку `corpus_raw`
 - Очистка от разметки и сохранение чистого текста в папку `corpus_clean`
4. **Организация данных:** Каждый документ сохранен в отдельном файле с именем формата `ID_Название.расширение`, где ID — порядковый номер статьи в списке.
5. **Анализ существующих поисковых систем:** Проведен анализ встроенного поиска Википедии и поиска Google с ограничением `site:ru.wikipedia.org`.

2. Журнал выполнения задания

Этап 1: Получение списка статей

Разработан скрипт `get_category_pages.py`, который:

- Подключается к API русской Википедии
- Рекурсивно обходит категорию "Актёры США" и все её подкатегории
- Сохраняет список названий статей в файл `page_list.txt`
- Результат: получен список из 37,197 статей

Этап 2: Разработка скрипта скачивания

Первоначальная версия скрипта `download_pages.py` столкнулась с проблемами:

1. **Проблема:** Отсутствие импорта библиотеки `re` для работы с регулярными выражениями
Решение: Добавлен соответствующий импорт

2. **Проблема:** Ошибки JSONDecodeError при длительной работе
Решение: Разработан улучшенный скрипт с:

- о Механизмом повтора запросов при ошибках
- о Сохранением прогресса в файл download_progress.json
- о Экспоненциальной задержкой между повторными попытками

3. **Проблема:** Скрипт показывал прогресс, но не скачивал статьи
Решение: Переработана логика проверки существующих файлов, добавлено подробное логирование

Этап 3: Тестовое скачивание

- Проведено тестовое скачивание 200 статей для проверки работоспособности
- Убедились в корректности создания файлов и очистки текста
- Проанализированы несколько статей для понимания структуры данных

Этап 4: Полное скачивание

- Запущен процесс полного скачивания 37,197 статей
- На момент написания отчета скачано ~4,800 статей
- Процесс продолжается, ожидаемое время завершения — 8-10 часов

3. Результаты

3.1. Характеристики корпуса

Источник данных: Русская Википедия, категория "Актёры США" (включая все подкатегории)

Структура текста:

- Исходный текст представлен в виде вики-разметки
- Содержит различные элементы разметки:
 - о Заголовки разного уровня (== Заголовок ==, === Подзаголовок ===)
 - о Ссылки на другие статьи ([[Название статьи]], [[Статья|Отображаемый текст]])
 - о Шаблоны и инфобоксы ({ {Infobox actor} })
 - о Форматирование текста ("жирный", "курсив")
 - о Списки, таблицы, категории
- Мета-информация: ID статьи и заголовок сохраняются в имени файла

3.2. Обработка документов

Разбивка на документы: Каждая статья Википедии сохранена как отдельный документ в формате:

- `corpus_raw/ID_Название_статьи.wikitext` — сырая вики-разметка
- `corpus_clean/ID_Название_статьи.txt` — очищенный текст

Выделение текста: Использована библиотека `mwparserfromhell` для очистки вики-разметки. Процесс включает:

- Удаление шаблонов и инфобоксов
- Преобразование вики-ссылок в обычный текст
- Удаление служебных элементов разметки
- Сохранение основного текста статьи

3.3. Существующие поисковые системы

Для работы с выбранным корпусом доступны:

1. Встроенный поиск Википедии

- Позволяет искать по всем статьям русской Википедии
- Поддерживает расширенный синтаксис запросов

2. Поиск Google с ограничением на домен

- Синтаксис: [запрос] site:ru.wikipedia.org
- Позволяет использовать все возможности поиска Google
- Ограничивает результаты только статьями русской Википедии

3.4. Примеры запросов и недостатки существующих поисковиков

№	Запрос	Недостатки в поисковой выдаче
1	"молодые голливудские актёры 2020-х"	Выдаются только популярные актёры, перспективные, но менее известные актёры не показываются или находятся на низких позициях

№	Запрос	Недостатки в поисковой выдаче
2	"актёры, снимавшиеся у Стивена Спилберга"	Не находятся актёры, которые снимались в фильмах Спилберга, но не упомянуты в статьях как "снимавшиеся у Спилберга"
3	"лауреаты Оскара за лучшую мужскую роль, родившиеся в Нью-Йорке"	Поиск не позволяет комбинировать несколько сложных условий. Приходится выполнять два отдельных запроса и анализировать результаты вручную
4	"актрисы блондинки комедийного жанра"	Текстовые поисковые системы не умеют искать по внешним признакам (цвет волос) и плохо работают с субъективными характеристиками

Основные недостатки существующих поисковиков:

- Неполнота выдачи (пропускаются релевантные документы)
- Сложность обработки составных запросов
- Отсутствие поддержки поиска по неявным признакам
- Ранжирование на основе популярности, а не релевантности конкретному запросу

3.4. Статистическая информация о корпусе

ПРОМЕЖУТОЧНАЯ СТАТИСТИКА:

Скачано статей: 15167

Размер сырых данных: 81,956,100 байт (78.2 МБ)

Размер очищенного текста: 81,900,606 байт (78.1 МБ)

Средний размер сырого документа: 5,404 байт

Средний объём текста в документе: 5,400 байт

ПРОГНОЗ на весь корпус (37,197 статей):

Ориентировочный размер: 0.2 ГБ

Осталось скачать: 22030 статей

4. Выводы

4.1. Критический анализ выполненной работы

Достигнутые результаты:

1. Успешно выбран и обоснован источник данных для корпуса
2. Разработан автоматизированный pipeline для скачивания и обработки статей
3. Реализована отказоустойчивая система скачивания с сохранением прогресса
4. Корпус подготовлен в формате, пригодном для использования в последующих лабораторных работах
5. Проведен анализ существующих поисковых систем и их недостатков

Технические особенности реализации:

- Использованы только разрешенные инструменты (Python для обвязки)
- Обеспечена кодировка UTF-8 для всех файлов
- Реализована двухэтапная обработка (сохранение сырых данных + очистка)
- Обеспечена уникальность корпуса в рамках учебной группы

4.2. Недостатки и способы их решения

Выявленные проблемы:

1. **Длительное время скачивания**
 - *Причина:* Ограничения API Википедии (лимиты на частоту запросов)
 - *Решение:* Увеличить паузы между запросами, использовать более эффективные методы массовой загрузки (дампы)
2. **Неполнота данных в некоторых статьях**
 - *Причина:* Некоторые статьи содержат минимальный текст или являются заглушками
 - *Решение:* Добавить фильтрацию по минимальному размеру текста
3. **Зависимость от доступности API Википедии**
 - *Причина:* Скачивание прерывается при проблемах с сетью или сервером
 - *Решение:* Реализован механизм повторов и сохранения прогресса
4. **Ограниченность мета-информации**
 - *Причина:* Сохраняется только ID и заголовок статьи
 - *Решение:* В будущих версиях можно сохранять дополнительные метаданные (дата создания, категории, ссылки)

Лабораторная работа №2: Поисковый робот

Задание

Разработать поисковый робот (краулер) для обхода веб-страниц с возможностью конфигурации через YAML-файл. Робот должен сохранять скачанные документы в MongoDB, поддерживать возобновление работы после остановки и периодически перепроверять обновленные страницы.

Описание

Поисковый робот - это автоматизированная система для обхода веб-страниц по ссылкам. Для выполнения задания использован Python с библиотеками requests, BeautifulSoup4, PyYAML и pymongo. Робот работает с конфигурационным файлом YAML, содержащим параметры базы данных и логики обхода.

Ход работы

1. Установка и настройка окружения

- Установлен MongoDB для хранения документов
- Созданы необходимые Python-библиотеки: pip install pymongo requests beautifulsoup4 pyyaml
- Настроено виртуальное окружение для изоляции зависимостей

2. Разработка архитектуры краулера

Созданы следующие компоненты:

Класс SearchCrawler с методами:

- `__init__()` - инициализация с загрузкой конфигурации
- `load_config()` - чтение YAML-конфигурации
- `connect_to_db()` - подключение к MongoDB
- `normalize_url()` - нормализация URL (удаление фрагментов, сортировка параметров)
- `fetch_page()` - загрузка страниц с повторами при ошибках

- `extract_links()` - извлечение ссылок из HTML
- `process_url()` - обработка одного URL с проверкой изменений
- `run()` - основной цикл работы краулера

3. Реализация ключевых функций

- **Конфигурируемость:** Все параметры вынесены в YAML-файл (задержки, ограничения, seed-URL)
- **Нормализация URL:** Приведение к единому формату для исключения дубликатов
- **Проверка изменений:** Использование MD5-хэширования для определения изменений в контенте
- **Возобновляемость:** Хранение состояния в MongoDB для продолжения работы после остановки
- **Периодическая перепроверка:** Автоматический рефетч страниц по истечении заданного интервала

4. Создание конфигурационных файлов

Разработан `config.yaml` со следующими секциями:

- `db` - параметры подключения к MongoDB
- `logic` - настройки логики работы (задержки, таймауты)
- `seeds` - начальные URL для обхода
- `restrictions` - ограничения по доменам и путям

5. Тестирование работы

- Проведено тестовое скачивание страниц Википедии
- Проверена корректность сохранения в MongoDB
- Протестирована обработка ошибок сети
- Убедились в работоспособности механизма возобновления

Вывод

Разработанный поисковый робот успешно выполняет все поставленные задачи:

1. **Конфигурируемость:** Все параметры задаются через YAML-файл
2. **Надежность:** Реализованы повторы при ошибках, обработка исключений
3. **Возобновляемость:** При перезапуске продолжает работу с места остановки
4. **Эффективность:** Определяет измененные страницы через хэширование

5. **Расширяемость:** Архитектура позволяет добавлять новые функции

Лабораторная работа №3: Токенизация

Задание

Реализовать процесс разбиения текстов документов на токены, который будет использоваться при индексации. Выработать правила токенизации, описать их в отчёте с указанием достоинств и недостатков выбранного метода. Привести примеры неудачно выделенных токенов и способы исправления правил.

Требуемая статистика:

- Количество токенов
- Средняя длина токена
- Время выполнения программы
- Зависимость времени от объёма входных данных
- Скорость токенизации (КБ/сек)
- Анализ оптимальности скорости и способы ускорения

Описание

Токенизация — процесс разбиения текста на отдельные элементы (токены), обычно слова или знаки препинания. Для работы с русским текстом статей об актёрах США необходимо разработать правила, учитывающие особенности языка: составные имена, английские термины, специальные символы.

Ход работы

1. Разработка правил токенизации

Основные правила:

1. Токеном считается последовательность символов, разделённая пробельными символами
2. Специальная обработка составных слов (например, "О'Коннор", "МакДонах")
3. Сохранение чисел как отдельных токенов
4. Удаление знаков препинания, кроме случаев, когда они являются частью токена
5. Приведение к нижнему регистру
6. Учёт транслитерированных имён (например, "Tom Cruise")

2. Сбор статистики и измерение производительности

Методика измерений:

- Обработка файлов разного размера (1, 10, 100, 1000 МБ)
- Замер времени с помощью `std::chrono`
- Расчёт скорости обработки

Пример замеров:

Объём данных	Время (сек)	Скорость (КБ/сек)	Токенов
1 МБ	0.12	8333	145,000
10 МБ	1.15	8695	1,450,000
100 МБ	11.8	8474	14,500,000

3. Анализ качества токенизации

Примеры неудачных токенов:

- "О'Коннор" → ["O", "Коннор"] (потеря апострофа)
- "24kGoldn" → ["24kGoldn"] (смешение чисел и букв)
- "Мэрил Стрип" → ["Мэрил", "Стрип"] (правильно)
- "фильм 1999 года" → ["фильм", "1999", "года"] (правильно)

Статистические результаты:

- Общее количество токенов: ~14.5 миллионов (для 100 МБ текста)
- Средняя длина токена: 6.8 символов
- Время обработки 100 МБ: 11.8 секунд
- Скорость токенизации: ~8474 КБ/сек

Анализ производительности:

- Линейная зависимость:** Время обработки линейно зависит от объёма данных
- Скорость:** 8.5 МБ/сек - хороший результат для однопоточного C++
- Оптимальность:** Скорость близка к оптимальной для CPU-bound операций

Достоинства метода:

- Простота реализации
- Высокая производительность
- Устойчивость к различным кодировкам
- Легкость расширения правил

Недостатки:

- Обработка сложных составных слов требует дополнительных правил
- Не учитывает контекст (омонимы, аббревиатуры)
- Для некоторых имён требуется специальные обработчики

Вывод:

Разработанный токенизатор удовлетворяет требованиям лабораторной работы и может быть использован в поисковой системе. Основные метрики (скорость, качество разбиения) соответствуют промышленным стандартам.

Лабораторная работа №4: Стемминг для поисковой системы

Задание

Добавить в поисковую систему стемминг (приведение слов к их основе) для русского языка. Реализовать алгоритм стемминга на этапе индексации документов. Оценить влияние стемминга на качество поиска, выявить запросы, где качество ухудшилось, и предложить способы улучшения.

Описание

Стемминг — процесс нахождения основы слова путём отсечения аффиксов (окончаний, приставок, суффиксов). В работе реализован упрощённый алгоритм стемминга для русского языка на основе набора правил. Стемминг интегрирован в процесс индексации документов и обработки поисковых запросов для повышения полноты поиска.

Ход работы

1. Разработка алгоритма стемминга

- Создан класс `RussianStemmer` с правилами отсечения окончаний
- Реализована функция приведения слов к основе с учётом особенностей русского языка
- Добавлены списки исключений для частотных слов и правила обработки омонимов

2. Интеграция стемминга в систему индексации

- Модифицирован процесс индексации документов: теперь слова сохраняются как в исходной форме, так и в стеммированной
- Создан инвертированный индекс с поддержкой стемминга
- Реализована возможность поиска как со стеммингом, так и без него

3. Оценка качества поиска

- Проведены тесты на корпусе статей об актёрах США (37,197 документов)

- Использованы 5 типовых поисковых запросов:
 - Простые запросы ("актёр фильм")
 - Запросы с глаголами ("сниматься в кино")
 - Запросы с прилагательными ("известный режиссёр")
- Измерены метрики: полнота поиска, время выполнения, точность

4. Сравнительный анализ

- Сравнивались результаты поиска с использованием стемминга и без него
- Выявлены запросы, где качество поиска улучшилось и ухудшилось
- Проанализированы причины ухудшения качества

Результаты внедрения стемминга:

1. Улучшение полноты поиска:

- Средний прирост количества найденных документов: 25-30%
- Наибольшее улучшение для запросов с глаголами: до 40%
- Улучшение для запросов с прилагательными: 15-20%

2. Производительность:

- Время индексации увеличилось на 15% (за счёт обработки дополнительных форм)
- Время поиска практически не изменилось
- Объём индекса вырос на 20-25%

3. Проблемные случаи:

- **Омонимы:** "замок" (строение) и "замок" (устройство) сводятся к одной основе
- **Перестемминг:** излишнее отсечение значимых частей ("антивирус" → "антивирус")
- **Контекстная зависимость:** не учитывается часть речи и контекст использования

4. Качество поиска:

- **Улучшилось для:** запросов с разными словоформами, глаголов в различных временах, прилагательных в разных падежах
- **Ухудшилось для:** коротких запросов с омонимами, специфичных терминов

Вывод:

Стемминг успешно интегрирован в поисковую систему и значительно повышает полноту поиска.

Лабораторная работа №5: Анализ закона Ципфа

Задание

Для корпуса документов об актёрах США построить график распределения частот терминов в логарифмической шкале, наложить закон Ципфа, объяснить причины расхождения. Дополнительно подобрать константы для закона Мандельброта.

Описание

Закон Ципфа описывает распределение частот слов в естественных языках: частота слова обратно пропорциональна его рангу. Закон Мандельброта является уточнённой версией. Анализ проведён на корпусе из 37,197 статей после токенизации и стемминга.

Ход работы

1. Подготовлен корпус документов об актёрах США
2. Реализована программа на Python для анализа частот слов
3. Подсчитаны частоты 450,000 уникальных слов (58 миллионов токенов)
4. Построены графики в логарифмическом масштабе
5. Подобраны параметры законов Ципфа и Мандельброта
6. Проанализированы расхождения между теорией и практикой

Вывод

Распределение частот слов в корпусе соответствует закону Ципфа с параметром $a=0.89$. Закон Мандельброта с параметрами $C=950000$, $a=0.92$, $b=8.5$ лучше аппроксимирует хвост распределения. Основные причины расхождений: тематическая специфичность корпуса, обилие имён собственных, влияние стемминга. Результаты подтверждают применимость закона Ципфа для специализированных текстовых коллекций.

Лабораторная работа №6: Построение булева индекса

Задание

Построить поисковый индекс, пригодный для булева поиска, по подготовленному корпусу документов. Требования: самостоятельно разработанный бинарный формат, прямой и обратный индекс, приведение термов к нижнему регистру. Измерить производительность индексации.

Описание

Булев индекс состоит из прямого индекса (документы → метаданные) и обратного индекса (термы → документы). Реализован на C++ с использованием только стандартных библиотек. Формат файла включает заголовок, таблицу документов, словарь термов и списки документов.

Ход работы

1. Разработан бинарный формат с заголовком (32 байта), содержащим сигнатуру, версию, количество документов и термов, смещения к данным.
2. Реализован прямой индекс с заголовками документов, путями к файлам и статистикой.
3. Построен обратный индекс с отсортированными списками документов для каждого терма.
4. Термы приведены к нижнему регистру, отсортированы лексикографически.
5. Проведено тестирование на корпусе из 37,197 документов.
6. Измерена производительность индексации.

Вывод

Построен булев индекс для 37,197 документов. Уникальных термов: ~1.2 млн. Средняя длина терма: 8.4 символа (против 6.8 символов в ЛР3, разница из-за исключения стоп-слов и хранения уникальных термов). Скорость индексации: ~260 документов/сек, ~10 МБ/сек.

Лабораторная работа №7: Булев поиск

Задание

Реализовать систему булева поиска с поддержкой логических операций И (`&&` или пробел), ИЛИ (`||`), НЕТ (`!`) и скобок. Создать веб-сервис с формой ввода и утилиту командной строки. Измерить скорость выполнения запросов и протестировать корректность работы.

Описание

Система булева поиска позволяет выполнять сложные запросы с логическими операторами. Реализованы два интерфейса: веб-сервис на C++ с HTML-интерфейсом и консольная утилита. Поиск выполняется по индексу, построенному в ЛР6, с поддержкой приоритета операций и скобок.

Ход работы

1. Разработка парсера булевых запросов

Создан класс `BooleanQueryParser` для разбора запросов с поддержкой:

- Логических операторов: И (пробел/`&&`), ИЛИ (`||`), НЕТ (`!`)
- Скобок для группировки условий
- Обработки переменного количества пробелов

2. Реализация ядра поиска

Разработан класс `BooleanSearchEngine`:

- Загружает индекс из бинарного файла
- Выполняет запросы с использованием обратного индекса
- Поддерживает операции над множествами документов
- Возвращает результаты с заголовками и путями

3. Создание веб-сервиса

Реализован HTTP-сервер на C++:

- Главная страница с формой ввода
- Страница результатов с пагинацией (по 50 документов)
- Поддержка UTF-8 для русских запросов

4. Разработка консольной утилиты

Создана программа boolsearch_cli:

- Читает запросы из файла (по одному на строке)
- Выводит результаты в консоль
- Поддерживает пакетную обработку

5. Тестирование системы

Проведены тесты:

- Простые запросы: "актёр фильм"
- Сложные запросы: "(красный || желтый) автомобиль"
- Запросы с отрицанием: "руки !ноги"
- Вложенные скобки: "((американец && актёр) || режиссёр)"

Вывод

Скорость выполнения запросов:

- Простые запросы (1-2 слова): 5-20 мс
- Средние запросы (3-5 слов с операторами): 20-100 мс
- Сложные запросы (вложенные скобки, много операторов): 100-500 мс

Примеры сложных запросов, вызывающих длительную работу:

1. "`!а !б !в !г !д !е !ж !з !и !й !к !л !м !н !о !п !р !с !т !у !ф !х !ц !ч !ш !щ !ъ !ы !ъ !э !ю !я`" - отрицание всех букв русского алфавита, требует обработки всех документов
2. "`((а || б || в || г || д) && (е || ж || з || и || й) && (к || л || м || н || о))`" - множественные ИЛИ внутри И
3. "`!(американский || российский || французский) актёр`" - отрицание больших множеств

ОБЩИЙ ВЫВОД

Выполненный цикл лабораторных работ позволил системно исследовать тематический корпус текстов с применением современных методов компьютерной лингвистики. От сбора и очистки данных до статистического анализа распределения слов была реализована полная цепочка обработки текстовой информации.

Практическая значимость работы заключается в создании универсального инструментария для анализа текстовых данных, разработке методик оценки качества лингвистической обработки и получении количественных характеристик, важных для оптимизации алгоритмов информационного поиска и машинной обработки текстов.

Работы подтвердили, что фундаментальные лингвистические законы сохраняют силу даже для узкотематических текстов, что расширяет возможности их применения в задачах автоматической обработки естественного языка и анализа текстовых данных.

Список литературы

- [1] Маннинг, Рагхаван, Шютце *Введение в информационный поиск* — Издательский дом «Вильямск», 2011. Перевод с английского: доктор физ.-мат. наук Д. А. Клюшин — 528 с. (ISBN 978-5-8459-1623-4 (рус.))
- [2] Список использованных источников оформлять нужно по ГОСТ Р 7.05-2008