

МИНОБРНАУКИ РОССИИ
Федеральное государственное бюджетное образовательное учреждение
высшего образования

**«САРАТОВСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ
ИМЕНИ Н. Г. ЧЕРНЫШЕВСКОГО»**

Кафедра математической кибернетики и компьютерных наук

**АВТОМАТИЧЕСКАЯ ТЕМАТИЧЕСКАЯ КЛАССИФИКАЦИЯ
НОВОСТНОГО МАССИВА**

БАКАЛАВРСКАЯ РАБОТА

студента 4 курса 451 группы
направления 09.03.04 — Программная инженерия
факультета КНиИТ
Кондрашова Даниила Владиславовича

Научный руководитель
доцент, к. ф.-м. н.

С. В. Папшев

Заведующий кафедрой
к. ф.-м. н.

С. В. Миронов

Саратов 2025

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	3
1 Теоретические и методологические основы	4
1.1 Получение текстовых данных	4
1.1.1 Выбор инструмента	4
1.1.2 Подбор информационной платформы	4
1.2 Подготовка текстовых данных	5
1.2.1 Выбор инструментов	5
2 Практико-технологические основы	7
2.1 Получение новостного массива путём вебскраппинга	7
2.2 Подготовка новостного массива	11
2.2.1 Удаление лишних пробелов и переносов строк	11
2.2.2 Разделение строк на русские и английские фрагменты	12
2.2.3 Обработка двоеточий и временных меток	13
2.2.4 Токенизация, лемматизация и удаление стоп-слов по словарю	14
2.2.5 Удаление стоп-слов с помощью метрики tfidf	15
2.3 Количественные характеристики обработанного и необработанного датасета	18
2.4 Вычисление тематической модели	18
ЗАКЛЮЧЕНИЕ	18
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	18
Приложение А Листинг вебскраппера	18
Приложение Б Листинг обработчика новостного массива	22
Приложение В Количественные характеристики подготовленного и неподготовленного новостного массива	25

ВВЕДЕНИЕ

В настоящее время обработка больших объёмов текстовых данных, включая новостные потоки, становится критически важной задачей. Как в научной среде, так и в бизнесе требуется оперативно анализировать информацию, отслеживать тенденции и принимать решения. Однако анализ всего массива данных невозможен из-за его масштабов, необходимо фильтровать информацию, оставляя только нужную.

Помочь в решении этой проблемы может тематическая классификация. Хотя многие сайты и порталы предлагают рубрикацию контента, её точность часто оказывается низкой: теги присваиваются некорректно или поверхностно. Это приводит к ошибкам в поиске и анализе информации.

В таком случае необходим механизм позволяющий получать правильную тематическую классификацию данных, который смог бы присваивать темы тем же новостям автоматически. Одни из возможных инструментов, которые позволяют реализовать подобие такого механизма — это тематические модели и алгоритмы машинного и глубокого обучения. Первый из них позволяет косвенно выявить темы текстового набора данных и разметить данные для обучения второго инструмента, который сможет тематически классифицировать последующий текст.

Таким образом, целью данной работы является реализация механизма автоматической тематической классификации новостей с помощью методов тематического моделирования и глубокого и машинного обучения.

Для достижения этой цели необходимо решить следующие задачи:

1. Реализовать механизм получения новостных массивов данных;
2. Реализовать механизм подготовки текстовых данных;
3. Вычислить тематические модели;
4. Путём сравнительного анализа выявить наиболее удачную тематическую модель;
5. Разметить данные для обучения на них моделей машинного и глубокого обучения;
6. Обучить несколько моделей машинного и глубокого обучения и выявить наиболее удачную путём сравнительного анализа;
7. Провести анализ получившихся результатов.

1 Теоретические и методологические основы

1.1 Получение текстовых данных

1.1.1 Выбор инструмента

Для получения каких-либо данных с сайта существует три основных метода:

- Ручной метод — выписывание необходимой информации с помощью человека;
- Получение данных путём предоставления их запроса у владельца, с их последующим скачиванием;
- Получение данных программным путём.

Первый метод из-за своей неэффективности можно сразу отбросить. Второй метод далеко не всегда можно применить, кроме того вряд ли владельцы информационных платформ будут оперативно отсылать все данные по первой просьбе. Таким образом, остаётся только третий метод.

Оперативно и достаточно эффективно в большинстве случаев можно получить данные применяя инструменты вебскраппинга. Дальше будет использоваться этот вариант получения новостного массива.

Различные библиотеки для вебскраппинга доступны на разных языках, однако исходя из того, что наиболее популярным языком для обработки данных и работы с машинным и глубоким обучением является python, выберем библиотеки доступные на нём. Такими библиотеками являются requests, beautifulsoup4 и selenium. Первая библиотека позволяет отсылать http запросы. Вторая библиотека позволяет преобразовывать html код в подобие классов для удобного получения информации. Последняя библиотека позволяет обрабатывать сайты, которые по http запросу не выдают html код наблюдаемой пользователем страницы. Данная библиотека позволяет эмулировать работу браузера и получать html код страницы прямо из него.

Такого набора хватит для обработки подавляющего большинства сайтов.

1.1.2 Подбор информационной платформы

В рамках данной работы среди всех типов текстовых данных будут рассматриваться новостные. Теперь нужно подобрать сайт.

Если для получения информации есть несколько возможных веб-источников, то стоит выбирать сайт по следующим критериям:

1. Сайт имеет единую структуру документов;
2. Сайт не блокирует http запросы отправляемые вебскраппером;
3. Сайт не является реактивным, то есть в момент просмотра страницы html код страницы полностью сформирован и доступен по запросу клиенту.

Будет идеально, если все пункты соблюдаются, одако, даже в случае отсутствия пунктов 2 и 3, ограничения в большинстве случаев можно достаточно просто обойти. В случае несоответствия пункта 1 могут возникнуть серьёзные трудности, которые, в худшем случае, решить только методами веб скраппинга не получится.

В рамках данной работы будет использоваться новостной сайт ВШЭ. Данный сайт соответствует всем описанным выше критериям.

1.2 Подготовка текстовых данных

Полученные данные требуют предварительной обработки для устранения шума и повышения качества анализа. Основные этапы включают:

1. Очистка от технического шума:
 - Удаление лишних пробелов, переносов строк;
 - Очистка от спецсимволов (скобки, HTML-теги, эмодзи);
 - Нормализация регистра (приведение всего текста к нижнему регистру).
2. Токенизация: Разделение текста на слова или предложения;
3. Лемматизация: Приведение слов к начальной форме (например, «бежал» ⇒ «бежать»);
4. Удаление стоп-слов: Исключение частых слов без смысловой нагрузки (предлоги, частицы, местоимения);

Обоснование выбора лемматизации вместо стемминга: Стемминг (например, алгоритм Snowball) «обрубает» окончания по шаблонам («бежал» ⇒ «беж»), что искажает смысл. Лемматизация сохраняет семантику, что критично для тематического моделирования.

1.2.1 Выбор инструментов

Чтобы не повышать количество используемых языков, будем рассматривать только инструменты, доступные на Python. Среди них выделяются: NLTK, Rymorphy3, SpaCy и Gensim.

Сделаем выбор между связкой NLTK + Rymorphy3 и SpaCy. Обе группы

библиотек позволяют проводить лемматизацию и удаление стоп-слов, но реализуют это по-разному. NLTK и Rymorphy3 приводят слова к начальной форме без учёта контекста, тогда как SpaCy — нейросетевой инструмент, анализирующий окружение терминов. Определение стоп-слов в обоих случаях происходит по заранее заданным словарям, поэтому разницы здесь нет. Однако SpaCy обеспечивает не только более точную лемматизацию, но и лаконичный интерфейс, что упрощает интеграцию в проект.

Как упоминалось ранее библиотека SpaCy определяет стоп-слова только по предопределённому списку, который не является исчерпывающим. Это связано с тем, что набор стоп-слов зависит от тематики текста, и универсального решения не существует. Для дополнительной фильтрации применим метрику TF-IDF, которая оценивает значимость слов. Формула расчёта:

$$tfidf(w, d) = \frac{n_{wd}}{n_d} \cdot \log \left(\frac{|D|}{|\{d \in D : w \in d\}|} \right), \quad (1)$$

где:

- w — термин;
- d — документ;
- n_{wd} — частота встречаемости w в d ;
- n_d — число терминов в d ;
- $|D|$ — число документов в коллекции;
- $|\{d \in D : w \in d\}|$ — количество документов, содержащих w .

Данная метрика будет тем выше для термина w в документе d , чем чаще будет встречаться термин w в документе d и реже во всех остальных документах коллекции. Таким образом, данную метрику можно интерпретировать как метрику значимости слова w для документа d . Её расчёт будет производиться с помощью библиотеки Gensim.

Таким образом, для обработки текста выбраны SpaCy (токенизация, лемматизация, базовые стоп-слов?) и Gensim (расширенная фильтрация через TF-IDF).

2 Практико-технологические основы

2.1 Получение новостного массива путём вебскраппинга

Для обработки такого простого новостного сайта как у ВШЭ достаточно использования `requests` и `beautifulsoap4`, без `selenium`.

Чтобы наиболее просто и эффективно получить данные необходимо разобрать структуру сайта и разработать соответствующие функции под каждую из частей. Сам портал представляет собой многостраничный сайт, на каждой странице которого расположено по 10 новостей с краткой информацией по каждой: ссылка, дата, заголовок, краткое содержание. На каждую новость можно перейти по ссылке для получения полного её содержания.

Теперь последовательно реализуем функции-обработчики под соответствующие части сайта.

Чтобы получать html код страницы, необходимо воспользоваться библиотекой `requests` и методом `get`. Данный метод отправляет запрос на сайт и получает соответствующий код в качестве ответа, который можно сохранить в файл для последующей выгрузки и обработки. Соответствующая функция расположена в листинге 1.

```
1 def __getPage__(url: str, file_name: str) -> None:
2     # получение html кода страницы с помощью библиотеки requests
3     r = requests.get(url=url)
4     # сохранение полученного кода в текстовый файл
5     with open(file_name, "w", encoding="utf-8") as file:
6         file.write(r.text)
```

Листинг 1: Функция получения html кода страницы

Далее нужно реализовать получение краткой информации о новости: ссылка, дата, краткое содержание. Для этого нужно загрузить код страницы из файла и преобразовать его к классам с помощью библиотеки `beautifulsoap4`. Далее можно будет воспользоваться поиском по тегам и классам с помощью метода `find` и получить текстовое содержимое с помощью методов `text` и `get`. Пример получения ссылки и краткого содержания новости можно увидеть в данном листинге 2.

```
1 # получение html кода страницы из файла
2 with open(page_file_name, encoding="utf-8") as file:
3     src = file.read()
4 # преобразование html кода в классы
```

```

5 soup = BeautifulSoup(src, "lxml")
6 # переход к содержимому новости, которое находится
7 # в теге div с классом post
8 news = soup.find("div", class_="post")
9 try:
10     # получение текста ссылки из соответствующего тега
11     link = news.find("h2",
12                     class_="first_child").find("a").get("href")
13     # не все ссылки в теге сохранены полностью, данный
14     # код добавляет обрезанную часть
15     if not link.startswith("https://"):
16         link = 'https://www.hse.ru' + link
17 except:
18     link = ""
19 try:
20     # получение краткого описания новости из соответствующего тега
21     a
22     news_short_content = news.find("p",
23                                     class_="first_child").find_next_sibling("p").text.strip()
24 except:
25     news_short_content = ""

```

Листинг 2: Получение ссылки и краткого содержания

Теперь нужно реализовать функцию получения полного содержания новости. Для этого нужно воспользоваться реализованной функцией `get_page` (получить код страницы по полученной ранее ссылке на новость), преобразовать его в классы с помощью `beautifulsoar4` и получить текстовое содержимое с помощью методов `find` и `text`. Реализацию соответствующей функции можно увидеть в листинге 3.

```

1 def __parse_news__(url: str) -> str:
2     # получаем html код страницы по ссылке на новость
3     news_file_name = "news.html"
4     __getPage__(url, news_file_name)
5     # и сразу загружаем его из файла
6     with open(news_file_name, encoding="utf-8") as file:
7         src = file.read()
8     # преобразуем html код к классам и сразу получаем всё текстовое
9     # содержание
10    # новости. Это возможно так как весь контент новости содержит
11    # ся

```



```

10     # в теге post__text
11     content = BeautifulSoup(src, "lxml").find("div",
12         class_="main").find(
13         "div", class_="post__text"
14     ).text.strip()
15     # возвращаем полученное содержание новости в виде строки
16     return content

```

Листинг 3: Функция получения полного текстового содержания новости

Следующим шагом нужно вспомнить, что на странице располагается 10 новостей, каждая новость располагается в теге div с классом post. Таким образом, нужно 10 раз проитерироваться по данным тегам и получить 10 новостей. Сделать это можно с помощью метода `find_next_sibling` (он ищет следующий тег, который идентичен по типу и классу предыдущему) и обычного цикла. Хранить полученное содержимое удобно в pandas DataFrame, так как с помощью него удобно обрабатывать полученные массивы данных и вычислять их количественные характеристики. Ключевые части соответствующей функции представлены в следующем листинге 4.

```

1  def __parse_page__(page_file_name: str, news_container:
2      pd.DataFrame) -> None:
3      # скрытый фрагмент получения html кода страницы
4      for i in range(10):
5          # скрытый фрагмент получения краткой информации о новости
6          try: # получение полного содержания новости
7              if link.startswith("https://www.hse.ru/news/"):
8                  news_content = __parse_news__(link)
9          except:
10             news_content = ""
11             # сохранение содержимого новости, если она не пустое
12             if len(
13                 news_day + news_month + news_year + news_name +
14                 news_short_content +
15                 news_content
16             ) > 0:
17                 news_container.loc[len(news_container.index)] = [
18                     link, news_date, news_name, news_short_content,
19                     news_content]
20             # переход к следующей новости

```

```
18 news = news.find_next_sibling("div", class_="post")
```

Листинг 4: Функция обработки одной страницы новостей

Далее необходимо реализовать функцию обрабатывающую все страницы с новостями. Сделать это можно путём многократного применения описанной выше функции обработки одной новостной страницы к изменяемой ссылке страницы. Благодаря простому устройству сайта ВШЭ менять эту ссылку можно достаточно просто с помощью обычного цикла путём изменения индекса в одной части. Соответствующий код представлен в следующем листинге 5.

```
1 def __crawling_pages__(start: int, end: int, news_container:
    pd.DataFrame, num_of_thread: int) -> pd.DataFrame:
2     page_file_name = "page.html"
3     for i in range(start, end + 1):
4         try:
5             __getPage__( "https://www.hse.ru/news/page{0}.html".format(i),
                           page_file_name)
6             __parse_page__(page_file_name, news_container)
7         except:
8             continue
```

Листинг 5: Функция обработки всех страниц новостей

Осталось только для ускорения получения данных с файла добавить многопоточность. Сделать это можно с помощью стандартных средств языка python, только стоит учесть, что под каждый отдельный поток нужно будет создать свой отдельный контейнер pandas DataFrame, чтобы избежать проблем с записью. Соответствующий код представлен в следующем листинге 6.

```
1 def crawling_pages(off_pc: bool, pages: int) -> None:
2     columns = ["url", "date", "title", "summary", "content"]
3     # создание контейнеров под каждый из потоков
4     news_container1 = pd.DataFrame(columns=columns)
5     news_container2 = pd.DataFrame(columns=columns)
6     # создание потоков
7     thread1 = threading.Thread(target=__crawling_pages__,
                                args=(0, pages // 2, news_container1, 1))
8     thread2 = threading.Thread(target=__crawling_pages__,
                                args=(pages // 2, pages, news_container2, 2))
9     # запуск потоков
10    thread1.start()
11    thread2.start()
```

```

12     # ожидание завершения работы потоков
13     thread1.join()
14     thread2.join()
15     # объединение содержимого контейнеров потоков в один
16     try:
17         news = pd.concat([news_container1, news_container2],
18                           ignore_index=True)
19         news.to_excel("./news.xlsx")
20     except:
21         print("Не получилось!")

```

Листинг 6: Многопоточное получение новостей

Полный код вебскраппера можно увидеть в соответствующем приложении **A**.

2.2 Подготовка новостного массива

2.2.1 Удаление лишних пробелов и переносов строк

Для корректной токенизации и просто для удобства анализа текстовых данных важно удалить из них лишние пробелы и переносы строк, сделать это можно с помощью стандартных средств языка python.

Функция будет иметь следующий алгоритм:

1. Записываем в результирующую строку символы из исходной строки, пока не будет встречен символ пробела или переноса строки;
2. Добавляем к результирующей строке 1 символ пробела и прекращаем добавление символов, пока не встретим символ, отличный от пробела или переноса строки;
3. В случае, когда встретится символ не являющийся пробелом или переносом строки переходим к пункту 1. Повторяем описанные выше действия пока не будет пройдена вся исходная строка.

Реализация соответствующей функции представлена в следующем листинге **7**.

```

1 def __remove_extra_spaces_and_line_breaks__(self, text: str) ->
2     str:
3     processed = ""
4     if type(text) != str or len(text) == 0:
5         return ""
6     flag = True

```

```

6     for symb in text:
7         if flag and (symb == " " or symb == "\n"):
8             processed += " "
9             flag = False
10        if symb != " " and symb != "\n":
11            flag = True
12        if flag:
13            processed += symb
14    return processed.strip()

```

Листинг 7: Функция удаления лишних пробелов и переносов строк

2.2.2 Разделение строк на русские и английские фрагменты

Библиотека SpaCy обрабатывает текст с помощью различных предобученных нейронных сетей, такие сети обучаются работе только на одном языке, например, только на русском или английском языке.

Текст новостей с новостного сайта ВШЭ имеет вставки на английском языке, что делает некорректным использование только одной предобученной нейронной сети. Поэтому, чтобы применять сразу два типа нейронных сетей необходимо разбивать строки на русские и английские фрагменты. Решить данную задачу можно с помощью стандартных средств языка python.

Функция будет иметь следующий алгоритм:

1. Определяем к какому алфавиту принадлежит первый буквенный символ строки и устанавливаем идентификатор в состояние соответствующее типу алфавита;
2. Записываем последовательно символы строки во временную подстроку, пока не встретим букву другого алфавита;
3. После встречи символа противоположного алфавита записываем в список кортеж вида (идентификатор алфавита, временная подстрока);
4. Очищаем временную подстроку и изменяем состояние идентификатора на противоположное. После этого повторяем описанные выше действия, пока не будет пройдена вся исходная строка.

Реализация соответствующей функции представлена в следующем листинге 8.

```

1 def __first_is_en__(self, cell: str) -> bool:
2     index_first_en = re.search(r"[a-zA-Z]", cell)
3     index_first_ru = re.search(r"[a-яА-Я]", cell)

```

```

4         return True if index_first_en and (not(index_first_ru)
        or index_first_en.start() <
        index_first_ru.start()) else False
5 def __split_into_en_and_ru__(self, cell: str) -> list[(bool,
    str)]:
6     parts = []
7     is_en = self.__first_is_en__(cell)
8     part = ""
9     for symb in cell:
10        if is_en == (symb in string.ascii_letters) or not
            (symb.isalpha()):
11            part += symb
12        else:
13            parts.append((is_en, part))
14            part = symb
15            is_en = not (is_en)
16    if part:
17        parts.append((is_en, part))
18    return parts

```

Листинг 8: Функция разбиения строки на русские и английские фрагменты

2.2.3 Обработка двоеточий и временных меток

При вычислении тематической модели BigARTM использует символ двоеточия как служебный, поэтому наличие его в текстовых данных приведёт к возникновению ошибок.

Само двоеточие, чаще всего, используется при написании времени, данные случаи можно обработать. Другие случаи применения предусмотреть проблематично, поэтому работать функция будет следующим образом: если двоеточие располагается в шаблоне временной метки, то будем заменять её на строку time, в противном случае будем просто удалять двоеточие.

Реализация соответствующей функции представлена в следующем листинге 9.

```

1 def __time_processing__(self, text: str) -> str:
2     if re.match(r"\d{2}:\d{2}", text):
3         return "time"
4     else:
5         return text.replace(":", "")
6

```

```

7 def __processing_token__(self, token_lemma: str) -> str:
8     return self.__time_processing__(
9         self.__remove_extra_spaces_and_line_breaks__(token_lemma)
10    )

```

Листинг 9: Функция обработки двоеточий и временных меток

2.2.4 Токенизация, лемматизация и удаление стоп-слов по словарю

Библиотека SpaCy имеет простой и удобный интерфейс. Для проведения токенизации, лемматизации и обнаружении стоп слов достаточно просто передать ей на вход строку. На выходе будет получен список объектов, в каждом из которых содержится по одному из токенов, их принадлежность к стоп-словам из словаря, начальная и исходная формы. С помощью этих объектов удобно записать в результирующую строку начальные формы токенов, которые не являются стоп-словами.

Пример применения библиотеки SpaCy к одной строке русского языка имеет следующий вид 10.

```

1 result = " ".join(
2     [
3         token.lemma_
4         for token in
5             self.nlp_en(self.__processing_token__(russian_str))
6             if
7                 not (token.is_stop) and not (token.is_punct) and
8                 len(token.lemma_) > 1
9     ]
10 )

```

Листинг 10: Пример применения библиотеки SpaCy для обработки одной строки русского языка

Реализация полного алгоритма, соержащего описанные выше функции представлена в следующем листинге 11.

```

1 def __processing_cell__(self, cell: str) -> str:
2     parts = self.__split_into_en_and_ru__(cell)
3     tokens = []
4     for part in parts:
5         if part[0]:
6             tokens += [

```

```

7         token.lemma_
8         for token in
9             self.nlp_en(self.__processing_token__(part[1]))
10            if not (token.is_stop) and not (token.is_punct)
11                and
12                len(token.lemma_) > 1
13            ]
14        else:
15            tokens += [
16                token.lemma_
17                for token in
18                    self.nlp_ru(self.__processing_token__(part[1]))
19                    if not (token.is_stop) and not (token.is_punct)
20                        and
21                        len(token.lemma_) > 1
22                ]
23        return " ".join(tokens)

```

Листинг 11: Функция удаления лишних пробелов и переносов строк, токенизации, лемматизации и удаления стоп-слов по словарю

2.2.5 Удаление стоп-слов с помощью метрики tfidf

Как говорилось ранее удаление стоп-слов только по словарю не может быть исчерпывающим, поэтому можно применить метрику tfidf для расчёта значимости слов и удалять слова с малой значимостью.

Расчёт этой метрики удобно с помощью библиотеки Gensim. Для этого нужно вычислить по коллекции документов словарь, затем по словарю сформировать частотный словарь — corpus, а уже по нему вычислить tfidf метрики для слов.

Реализация соответствующей функции представлена в следующем листинге 12.

```

1 def
2     calc_tfidf_corpus_without_zero_score_tokens_and_tfidf_dictionary(s
3     -> None:
4     texts = []
5     self.original_tokens = []
6     for row in range(self.p_data.shape[0]):
7         words = []
8         for column in self.processing_columns:

```

```

7         for word in self.p_data.loc[row, column].split(" "):
8             words.append(word)
9         self.original_tokens.append(words)
10        texts.append(words)
11        dictionary = gensim.corpora.Dictionary(texts)
12        corpus = [dictionary.doc2bow(text) for text in texts]
13        tfidf = gensim.models.TfidfModel(corpus)
14        self.tfidf_corpus = tfidf[corpus]
15        self.tfidf_dictionary = dictionary

```

Листинг 12: Функция вычисления tfidf метрики для слов документов

Однако данное вычисление не является полным, так как библиотека Gensim не добавляет в словарь слова, значение tfidf которых точно будет равняться нулю. В таком случае необходимо добавить недостающие слова. Реализация соответствующей функции представлена в следующем листинге 13.

```

1 def add_in_tfidf_corpus_zero_score_tokens(self) -> None:
2     full_corpus = []
3     for doc_idx, doc in enumerate(self.tfidf_corpus):
4         original_words = self.original_tokens[doc_idx]
5         term_weights = {self.tfidf_dictionary.get(term_id):
6                         weight for term_id, weight in doc}
7         full_doc = []
8         for word in original_words:
9             if word in term_weights:
10                weight = term_weights[word]
11            else:
12                weight = 0.0
13            full_doc.append((word, weight))
14        full_corpus.append(full_doc)
15    self.tfidf_corpus = full_corpus

```

Листинг 13: Функция добавление недостающих tfidf слов

Последним шагом перед удалением стоп-слов является вычисление границы, по которой будет определяться принадлежность к стоп-словам. Сделать это можно следующим образом 14.

```

1 def add_in_tfidf_corpus_zero_score_tokens(self) -> None:
2     full_corpus = []
3     for doc_idx, doc in enumerate(self.tfidf_corpus):
4         original_words = self.original_tokens[doc_idx]

```



```

5         term_weights = {self.tfidf_dictionary.get(term_id):
                           weight for term_id, weight in doc}
6     full_doc = []
7     for word in original_words:
8         if word in term_weights:
9             weight = term_weights[word]
10        else:
11            weight = 0.0
12        full_doc.append((word, weight))
13    full_corpus.append(full_doc)
14    self.tfidf_corpus = full_corpus

```

Листинг 14: Функция вычисления tfidf границы

В данном случае к стоп-словам будут относиться слова, значение tfidf метрики которых будет относиться к n минимальным процентам значений.

Теперь осталось только пройти по датасету и удалить соответствующие стоп-слова. Реализация соответствующей функции представлена в следующем листинге 15.

```

1 def del_tfidf_stop_words(self, tfidf_percent_threshold) -> None:
2     self.calc_tfidf_corpus_without_zero_score_tokens_and_tfidf_dictionary
3     self.add_in_tfidf_corpus_zero_score_tokens()
4     self.calc_threshold_for_tfidf_stop_words(tfidf_percent_threshold)
5     for row, doc in zip(range(self.p_data.shape[0]),
6                          self.tfidf_corpus):
7         tfidf_stop_words = [word for word, tfidf_value in doc if
8                             tfidf_value < self.threshold_for_tfidf_stop_words]
9         for column in self.processing_columns:
10            words_without_tfidf_stop_words = []
11            for word in self.p_data.loc[row, column].split(" "):
12                if word in tfidf_stop_words:
13                    continue
14                words_without_tfidf_stop_words.append(word)
15            self.p_data.loc[row, column] = "
16            ".join(words_without_tfidf_stop_words)

```

Листинг 15: Функция удаление вычисленных по метрике tfidf стоп-слов

Также стоит сказать, что также дополнительно стоит добавить удаление низкочастотных слов, так как это может положительно повлиять на результаты тематического моделирования.

Полный код обработчика новостного массива можно увидеть в соответствующем приложении **Б**.

2.3 Количественные характеристики обработанного и необработанного датасета

В рамках данной работы была выполнена обработка новостного массива с различными параметрами (имеется ввиду разные пороги для tfidf метрик, а также некоторые другие). Количественные характеристики представлены в соответствующих таблицах **В**.

2.4 Вычисление тематической модели

ЗАКЛЮЧЕНИЕ

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

ПРИЛОЖЕНИЕ А

Листинг вебскраппера

```
1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4 import os
5 import time
6 import threading
7
8 def __loading_bar_and_info__(
9     start: bool, number_of_steps: int, total_steps: int,
10     number_of_thread: int
11 ) -> None:
12     '''Вывод информации о прогрессе выполнения программы.
13     start - нужно ли вывести начальную строку;
14     number_page - количество спаршенных страниц;
15     total_pages - всего страниц, которые нужно спарсить;
16     miss_count - число новостей, которые не удалось спарсить;
17     whitour_whole_content - число новостей, у которых не получило
18         сь полностью спарсить контент.'''
19     done = int(number_of_steps / total_steps * 100) if int(
20         number_of_steps / total_steps * 100
21     ) < 100 or number_of_steps == total_steps else 99
22     stars = int(
```

```

21         40 / 100 * done
22     ) if int(20 / 100 * done) < 20 or number_of_steps ==
        total_steps else 39
23     tires = 40 - stars
24
25     if start:
26         stars = 0
27         tires = 40
28         done = 0
29
30     print("thread{0} <".format(number_of_thread), end="")
31     for i in range(stars):
32         print("*", end="")
33
34     for i in range(tires):
35         print("-", end="")
36     print("> {0}% ||| {1} / {2}".format(done, number_of_steps,
        total_steps))
37
38 def __getPage__(url: str, file_name: str) -> None:
39     '''Получение html файла страницы.
40     url - ссылка на страницу;
41     file_name - имя файла, в который будет сохранена страница.'''
42     r = requests.get(url=url)
43
44     with open(file_name, "w", encoding="utf-8") as file:
45         file.write(r.text)
46
47 def __parse_news__(url: str) -> str:
48     '''Получение полного контента новости.
49     url - ссылка на новость.
50     Функция возвращает полный текст новости.'''
51     news_file_name = "news.html"
52     __getPage__(url, news_file_name)
53
54     with open(news_file_name, encoding="utf-8") as file:
55         src = file.read()
56
57     content = BeautifulSoup(src, "lxml").find("div",
        class_="main").find(
58         "div", class_="post__text"

```

```

59     ).text.strip()
60
61     return content
62
63 def __parse_page__(page_file_name: str, news_container:
64     pd.DataFrame) -> None:
65     '''Парсинг информации с новостной страницы: ссылка на новость
66         + короткая информация о ней.
67     page_file_name - имя файла, в который сохранён код страницы;
68     news_container - таблица, в которую заносится информация о но
69         вости.
70     Функция также возвращает количество новостей, которые не удал
71         ось спарсить
72     и количество новостей, полный контент которых спарсить не уда
73         лось.'''
74     with open(page_file_name, encoding="utf-8") as file:
75         src = file.read()
76
77     soup = BeautifulSoup(src, "lxml")
78
79     news = soup.find("div", class_="post")
80     for i in range(10):
81         try:
82             news_day = news.find("div",
83                                     class_="post-meta__day").text.strip()
84         except:
85             news_day = ""
86
87         try:
88             news_month = news.find("div",
89                                     class_="post-meta__month").text.strip()
90         except:
91             news_month = ""
92
93         try:
94             news_year = news.find("div",
95                                     class_="post-meta__year").text.strip()
96         except:
97             news_year = ""
98
99     news_date = news_day + "." + news_month + "." + news_year

```

```

93
94     try:
95         news_name = news.find("h2",
96                                class_="first_child").find("a").text.strip()
97     except:
98         news_name = ""
99
100    try:
101        news_short_content = news.find("p",
102                                         class_="first_child"
103                                         ).find_next_sibling("p").text.strip()
104    except:
105        news_short_content = ""
106
107    try:
108        link = news.find("h2",
109                          class_="first_child").find("a").get("href")
110        if not link.startswith("https://"):
111            link = 'https://www.hse.ru' + link
112    except:
113        link = ""
114
115    try:
116        if link.startswith("https://www.hse.ru/news/"):
117            news_content = __parse_news__(link)
118    except:
119        news_content = ""
120
121    if len(
122        news_day + news_month + news_year + news_name +
123        news_short_content +
124        news_content
125    ) > 0:
126        news_container.loc[len(news_container.index)] = [
127            link, news_date, news_name, news_short_content,

```

Листинг 16: Полный код вебскрапера

ПРИЛОЖЕНИЕ Б

Листинг обработчика новостного массива

```
1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4 import os
5 import time
6 import threading
7
8 def __loading_bar_and_info__(
9     start: bool, number_of_steps: int, total_steps: int,
10     number_of_thread: int
11 ) -> None:
12     '''Вывод информации о прогрессе выполнения программы.
13     start - нужно ли вывести начальную строку;
14     number_page - количество спаршенных страниц;
15     total_pages - всего страниц, которые нужно спарсить;
16     miss_count - число новостей, которые не удалось спарсить;
17     whitour_whole_content - число новостей, у которых не получило
18         сь полностью спарсить контент.'''
19     done = int(number_of_steps / total_steps * 100) if int(
20         number_of_steps / total_steps * 100
21     ) < 100 or number_of_steps == total_steps else 99
22     stars = int(
23         40 / 100 * done
24     ) if int(20 / 100 * done) < 20 or number_of_steps ==
25         total_steps else 39
26     tises = 40 - stars
27
28     if start:
29         stars = 0
30         tises = 40
31         done = 0
32
33     print("thread{0} <".format(number_of_thread), end="")
34     for i in range(stars):
35         print("*", end=" ")
36
37     for i in range(tises):
38         print("-", end=" ")
39     print("> {0}% ||| {1} / {2}".format(done, number_of_steps,
```

```

        total_steps))

37
38 def __getPage__(url: str, file_name: str) -> None:
39     '''Получение html файла страницы.
40     url - ссылка на страницу;
41     file_name - имя файла, в который будет сохранена страница.'''
42     r = requests.get(url=url)
43
44     with open(file_name, "w", encoding="utf-8") as file:
45         file.write(r.text)
46
47 def __parse_news__(url: str) -> str:
48     '''Получение полного контента новости.
49     url - ссылка на новость.
50     Функция возвращает полный текст новости.'''
51     news_file_name = "news.html"
52     __getPage__(url, news_file_name)
53
54     with open(news_file_name, encoding="utf-8") as file:
55         src = file.read()
56
57     content = BeautifulSoup(src, "lxml").find("div",
58         class_="main").find(
59         "div", class_="post__text"
60     ).text.strip()
61
62     return content
63
64 def __parse_page__(page_file_name: str, news_container:
65     pd.DataFrame) -> None:
66     '''Парсинг информации с новостной страницы: ссылка на новость
67     + короткая информация о ней.
68     page_file_name - имя файла, в который сохранён код страницы;
69     news_container - таблица, в которую заносится информация о но
70     вости.
71     Функция также возвращает количество новостей, которые не удал
72     ось спарсить
73     и количество новостей, полный контент которых спарсить не уда
74     лось.'''
75     with open(page_file_name, encoding="utf-8") as file:
76         src = file.read()

```

```

71
72 soup = BeautifulSoup(src , "lxml")
73
74 news = soup.find("div" , class_="post")
75 for i in range(10):
76     try:
77         news_day = news.find("div" ,
78                               class_="post-meta__day").text.strip()
79     except:
80         news_day = ""
81
82     try:
83         news_month = news.find("div" ,
84                                class_="post-meta__month").text.strip()
85     except:
86         news_month = ""
87
88     try:
89         news_year = news.find("div" ,
90                                class_="post-meta__year").text.strip()
91     except:
92         news_year = ""
93
94     news_date = news_day + "." + news_month + "." + news_year
95
96     try:
97         news_name = news.find("h2" ,
98                                class_="first_child").find("a").text.strip()
99     except:
100         news_name = ""
101
102     try:
103         news_short_content = news.find("p" ,
104                                         class_="first_child"
105                                         ).find_next_sibling("p").text.strip()
106     except:
107         news_short_content = ""
108
109     try:
110         link = news.find("h2" ,
111                           class_="first_child").find("a").get("href")

```



```

108         if not link.startswith("https://"):
109             link = 'https://www.hse.ru' + link
110     except:
111         link = ""
112
113     try:
114         if link.startswith("https://www.hse.ru/news/"):
115             news_content = __parse_news__(link)
116     except:
117         news_content = ""
118
119     if len(
120         news_day + news_month + news_year + news_name +
121         news_short_content +
122         news_content
123     ) > 0:
124         news_container.loc[len(news_container.index)] = [
125             link, news_date, news_name, news_short_content,
126             news_content
127
128         ]
129
130     news = news.find_next_sibling("div", class_="post")

```

Листинг 17: Полный код подготовки новостного массива

ПРИЛОЖЕНИЕ В

Количественные характеристики подготовленного и неподготовленного новостного массива

Характеристика	Неподгот.	Стоп-слова	+Низкочаст.	TF-IDF 1%	TF-IDF 2%	TF-IDF 3%
Кол. док.	17340	17340	17340	17340	17340	17340
Кол. токенов	1213111	16545045	-	6479545	6414045	6348544

Продолжение следует...

Продолжение таблицы

Характеристика	Неподгот.	Стоп-слова	+Низкочаст.	TF-IDF 1%	TF-IDF 2%	TF-IDF 3%
Кол. уник. ток.	278724	148677	-	148677	148677	148677
Мин. кол. ток. в док.	6	4	-	4	4	4
Модальное кол. ток. в док.	47	31	-	31	31	30
Среднее кол. ток. в док.	695	375	-	371	367	364
Медианное кол. ток. в док.	-	313	-	312	310	309
Макс. кол. ток. в док.	6514	3151	-	2903	2825	2766
Мин. кол. уник. ток. в док.	6	4	-	4	4	4
Мод. кол. уник. ток. в док.	39	27	-	27	27	30
Сред. кол. уник. ток. в док.	346	214	-	211	208	205
Мед. кол. уник. ток. в док.	-	186	-	185	183	182

Продолжение следует...

Продолжение таблицы

Характеристика	Неподгот.	Стоп-слова	+Низкочаст.	TF-IDF 1%	TF-IDF 2%	TF-IDF 3%
Макс. кол. уник. ток. в док.	2287	1353	-	1299	1262	1214

Характеристика	TF-IDF 4%	TF-IDF 5%	TF-IDF 6%.	TF-IDF 7%	TF-IDF 8%	TF-IDF 9%
Кол. док.	17340	17340	17340	17340	17340	17340
Кол. токенов	6283046	6217544	6152044	6086544	6021044	5955543
Кол. уник. ток.	148677	148677	148677	148677	148677	148677
Мин. кол. ток. в док.	4	4	4	4	4	4
Модальное кол. ток. в док.	30	30	30	30	29	29
Среднее кол. ток. в док.	360	356	352	349	345	341
Медианное кол. ток. в док.	307	306	305	303	301	299

Продолжение следует...

Продолжение таблицы

Характеристика	TF-IDF 4%	TF-IDF 5%	TF-IDF 6%	TF-IDF 7%	TF-IDF 8%	TF-IDF 9%
Макс. кол. ТОК. В ДОК.	2713	2662	2595	2545	2501	2424
Мин. кол. уник. ТОК. В ДОК.	4	4	4	4	4	4
Мод. кол. уник. ТОК. В ДОК.	27	29	29	28	28	28
Сред. кол. уник. ТОК. В ДОК.	201	198	195	192	189	186
Мед. кол. уник. ТОК. В ДОК.	181	179	177	176	174	172
Макс. кол. уник. ТОК. В ДОК.	1164	1122	1085	1047	1018	986

Характеристика	TF-IDF 10%	TF-IDF 10% + Низк.
Кол. док.	17340	17340
Кол. токенов	5890042	-
Кол. уник. ток.	148677	-
Мин. кол. ток. в док.	4	-
Модальное кол. ток. в док.	30	-
Среднее кол. ток. в док.	337	-
Медианное кол. ток. в док.	297	-
Макс. кол. ток. в док.	2391	-
Мин. кол. уник. ток. в док.	4	-
Мод. кол. уник. ток. в док.	28	-
Сред. кол. уник. ток. в док.	182	-

Продолжение следует...

Продолжение таблицы

Характеристика	TF-IDF 10%	TF-IDF 10% + Низк.
Мед. кол. уник. ток. в док.	170	-
Макс. кол. уник. ток. в док.	946	-