

Министерство науки и высшего образования Российской Федерации
федеральное государственное автономное образовательное учреждение высшего
образования
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО»

Отчет

по заданию №3
«Расстояние Левенштейна»

по дисциплине «Компьютерная лингвистика»

Автор: Лакиза Александр Николаевич

Факультет: ИКТ

Группа: К3242

Преподаватель: Чернышева Анастасия Владимировна



УНИВЕРСИТЕТ ИТМО

Санкт-Петербург 2021

Цель работы: ознакомиться с таким алгоритмом, применяемом в обработке языка, как Расстояние Левенштейна

Ссылка на исходный код и json файл с корпусом: <https://github.com/alexanderlakiza/cs224>

Вся работа была сделана в файлах `task3.py` и `task3 ui.py`. JSON-файлы, которые были созданы в процессе работы `corpus_as_dict.json`, `corpus_as_dict_of_norms.json`, `each_token_appearance.json`, `unique_norms_in_corpus.json`.

Ход Работы: Для начала посмотрим файл `task3.py`. Мне необходимо создать словарь, в котором для каждого токена из корпуса будет подставлено значение с id документов, в котором этот токен встречается. Для этого сначала надо предобработать текст.

```
1. for i in range(len(texts)):
2.     texts[i] = re.sub(r'\n', ' ', texts[i]) # Удаление символа
        переноса строки
3.
4.     while len(re.findall(r'\([^()]*\)', texts[i])) != 0: #
        Удаление скобок и пробелов
5.         texts[i] = re.sub(r'\([^()]*\)', '', texts[i])
6.         texts[i] = re.sub(r'\s+(?=[,.?!;:...])', '', texts[i])
7.         texts[i] = re.sub(r'\s+', ' ', texts[i])
8.     # Заменяем диакритические знаки
9.     texts[i] = del_diacritics(texts[i])
```

Для этого я использую методы, описанные в [Задании 2](#). Затем я записал весь корпус (где документы отредактированы) в словарь `corpus_as_dict.json` (В этом словаре ключ = заголовок, значение = отредактированный текст документа).

Далее я подгружаю свой корпус, записанный в одну строку, сделанный в Задании 2. Также чищу его на всякий случай. И получаю список всех уникальных токенов, привожу их все к начальной форме с помощью `rumorphy2`. Записываю этот список в `unique_norms_in_corpus.json`.

```
1. def del_punct(tokens_list):
2.     punct = string.punctuation + "-" + "<" + ">"
3.     return [token.lower() for token in tokens_list if token not
        in punct]
4.
5.
6. tokens = word_tokenize(corpus)
7. tokens = del_punct(tokens)
8. normed_tokens = [morph.parse(i)[0].normal_form for i in tokens] #
        Список начальных форм всех токенов
9.         if str.isalpha(morph.parse(i)[0].normal_form)]
10.     print(normed_tokens)
11.     unique_norm_tokens = list(set(normed_tokens))
12.     with open("unique_norms_in_corpus.json", "w") as f:
```

```

13.         # Запишем корпус в виде словаря в json, чтобы потом
            использовать готовый словарь
14.         # вместо использования руморphy2
15.         json.dump(unique_norm_tokens, f, ensure_ascii=False)

```

Далее я записываю корпус в еще один словарь. В этом словаре заголовок документа = ключ, а значение = список токенов в начальной форме этого документа. Записываю в *corpus_as_dict_of_norms.json*

```

1. main_normed_dict = {}
2. md_keys = list(main_dict.keys())
3. md_values = list(main_dict.values())
4. for i in range(len(main_dict)):
5.     tokened_value = md_values[i]
6.     tokened_value = word_tokenize(tokened_value)
7.     tokened_value = del_punct(tokened_value)
8.     normed_tokens = [morph.parse(i)[0].normal_form for i in
        tokened_value
9.
        if
        str.isalpha(morph.parse(i)[0].normal_form)]
10.    main_normed_dict[md_keys[i]] = normed_tokens
11.    with open("corpus_as_dict_of_norms.json", "w") as f:
12.        # Запишем корпус в виде словаря в json, чтобы потом
            использовать готовый словарь
13.        # вместо использования руморphy2
14.        json.dump(main_normed_dict, f, ensure_ascii=False)

```

И последнее, что я делаю в файле [task3.py](#): Я создаю необходимый нам словарь, в котором для каждого токена, использованного в тексте в разных формах слова, ставится в значение словарь, в котором ключ = количество использований токена, а значение = заголовки документов, в которых появляется токен.

```

1. with open("unique_norms_in_corpus.json") as f:
2.     unique_norm_tokens = json.load(f)
3.
4. with open("corpus_as_dict_of_norms.json") as f:
5.     main_normed_dict = json.load(f)
6.
7. word_usage = {}
8. for word in unique_norm_tokens:
9.     counter = 0
10.    docs_include_word = {}
11.    list_docs_include_word = []
12.    for i in range(len(main_normed_dict)):

```

```

13.         if word in list(main_normed_dict.values())[i]:
14.             # print('Номер документа:', i)
15.             c =
                list(main_normed_dict.values())[i].count(word)
16.             counter += c
17.             list_docs_include_word.append(list(main_normed_d
                ict.keys())[i])
18.             docs_include_word[counter] = list_docs_include_word
19.             word_usage[word] = docs_include_word
20.         with open("each_token_appearance.json", "w") as f:
21.             # Запишем корпус в виде словаря в json, чтобы потом
                использовать готовый словарь
22.             json.dump(word_usage, f, ensure_ascii=False)

```

Записываем этот словарь в *each_token_appearance.json*.

Стоит отметить, что в исходном коде на гитхабе многие этапы будут закомментированы. Это было сделано для того, чтобы ускорить процесс работы – вместо того, чтобы каждый раз пробегаться по всем тысячам токенов и ждать ~2 минуты, я записывал промежуточные этапы в json файлы. Поэтому их и получилось 4 новых.

Переходим к task3_ui.py.

Я чутка видоизменил функцию `check_spelling(text)`, которая была нам продемонстрирована на паре, добавив в нее пару условий.

```

1. def check_spelling(text):
2.     domain =
        "https://speller.yandex.net/services/spellservice.json"
3.     words = text.split()
4.     if len(words) == 1:
5.         request = requests.get(domain + "/checkText?text=" +
            words[0])
6.         if requests:
7.             if not request.json():
8.                 for i in range(len(words[0])):
9.                     if ord(words[0][i]) < 1040:
10.                        return 'В слове "{}" присутствуют
                            буквы не русского языка'.format(words[0])
11.                        return 'В слове "{}" нет
                            ошибки'.format(words[0])
12.                     else:
13.                        return request.json()[0]["word"],
                            request.json()[0]["s"]

```

```

14.         else:
15.             return None
16.         elif len(words) > 1:
17.             words = "+".join(words)
18.             request = requests.get(domain + "/checkText?text=" +
words)
19.             if requests:
20.                 if not request.json():
21.                     words = words.split("+")
22.                     for word in words:
23.                         for i in range(len(word)):
24.                             if ord(word[i]) < 1040:
25.                                 return 'В слове "{}"
присутствуют буквы не русского языка'.format(word)
26.                                 return 'В словах {} нет
ошибок'.format(words)
27.         else:
28.             response = [(i["word"], i["s"]) for i in
request.json()]
29.             return response
30.         else:
31.             return None
32.         return None

```

Однако в дальнейшем понял, что мне не понадобится эта функция, так как заключается немного в другом. Я решил написать функцию, в которой мы используем Расстояние Левенштейна, следующим образом: пользователь вводит слово в начальной, если оно есть среди токенов, то программа без задней мысли выводит сколько раз встречается этот токен, и в каких документах. Если же токен не встречается (в введённом слове допущена ошибка или этого токена просто нет в документах), то программа пробегается по всем токенам и выбирает тот, у которого Расстояние Левенштейна с введённым словом наименьшее.

```

1. if __name__ == "__main__":
2.     print('Данная программа покажет вам названия документов, в
которых встречается введенное вами слово\n\n')
3.     'Введите ваше слово в начальной форме: ')
4.     word = input()
5.     if word in list(word_usage.keys()):
6.         show(word)
7.     else:
8.         print('Слово "{}" в документах не
наблюдается\n'.format(word))

```

```

9.         print('Возможно вы имели в виду слово
               "{}"\n'.format(fix_flaw(word)))
10.         fixed_word = fix_flaw(word)
11.         show(fixed_word)

```

Функция show(text):

```

1. with open("each_token_appearance.json") as f:
2.     word_usage = json.load(f)
3.
4.
5. def show(word):
6.     dict_of_docs = word_usage[word]
7.     number_of_mentions = list(dict_of_docs.keys())[0]
8.     print('Слово "{}" встречается {} раз в документах:
           {}'.format(word, number_of_mentions,
9.                       dict_of_docs[number_of_me
ntions]))

```

Функция fix_flaw(text):

```

1. def fix_flaw(word):
2.     # options = check_spelling(word)[1]
3.     lev_distances = {}
4.     for i in list(word_usage.keys()):
5.         lev_distances[i] = lev.distance(word, i)
6.     """
7.     Реализация исправления ошибки через Яндекс спелл чекер
8.     for i in options:
9.         lev_distances[i] = lev.distance(word, i)
10.    """
11.    minimum =
        list(lev_distances.values()).index(min(list(lev_distances.values(
        ))))
12.    best_option = list(lev_distances.keys())[minimum]
13.    return best_option

```

Как видно из кода, в кавычках находится реализация выбора правильного варианта введенного пользователем слова через Яндекс спелл чекер. Я использовал готовую реализацию Расстояния Левенштейна из библиотеки python-Levenshtein.

Примеры работы программы:

```
alexanderlakiza@alela-KPL-W0X:~/pyfiles/cs224$ workon cs224
(cs224) alexanderlakiza@alela-KPL-W0X:~/pyfiles/cs224$ python task3_ui.py
Данная программа покажет вам названия документов, в которых встречается введенное вами слово

Введите ваше слово в начальной форме:
корабль
Слово "корабль" встречается 58 раз в документах: ['Балласт на надводных плавучих средствах', 'Бимс', 'Бортовая качка', 'Ватерлиния', 'Верфь', 'Водоизмещение', 'Вьюшка', 'Капитан судна', 'Каята', 'Килевая качка', 'Киль', 'Лавировка', 'Лоцман', 'Марсовая площадка', 'Морские термины', 'Нактоуз', 'Непотопляемость', 'Осадка', 'Палуба', 'Полубак', 'Рулевой', 'Седловатость', 'Сначала женщины и дети', 'Спасательная шлюпка', 'Судовой руль', 'Сухой док', 'Те-Солент', 'Титаник Белфаст', 'Титаник-Белфаст', 'Титаник. Легенда продолжается', 'Тоннаж', 'Трюм', 'Форштевень', 'Швартов', 'Швартовка', 'Шлюпбалка', 'Шпангоут', 'Штормтрап', 'Штраус-парк', 'Якорь']
(cs224) alexanderlakiza@alela-KPL-W0X:~/pyfiles/cs224$ python task3_ui.py
Данная программа покажет вам названия документов, в которых встречается введенное вами слово

Введите ваше слово в начальной форме:
кайтх
Слово "кайтх" в документах не наблюдается

Возможно вы имели в виду слово "уайт"

Слово "уайт" встречается 12 раз в документах: ['Джозеф Брюс Исмей', 'Номадик', 'Сначала женщины и дети', 'Те-Солент', 'Уайт Стар Лайн']
(cs224) alexanderlakiza@alela-KPL-W0X:~/pyfiles/cs224$ python task3_ui.py
Данная программа покажет вам названия документов, в которых встречается введенное вами слово

Введите ваше слово в начальной форме:
аккомпанемент
Слово "аккомпанемент" в документах не наблюдается

Возможно вы имели в виду слово "компонент"

Слово "компонент" встречается 12 раз в документах: ['Бинокль', 'Дымогарные трубы', 'Жаровые трубы', 'Латунь', 'Минеральная вода', 'Мороженое', 'Сахар', 'Тонзиллит острый', 'Электрическое поле']
(cs224) alexanderlakiza@alela-KPL-W0X:~/pyfiles/cs224$ |
```