

Московский государственный технический университет им. Н.Э. Баумана
Кафедра «Системы обработки информации и управления»

Лабораторная работа №5

по дисциплине «Методы машинного обучения»

на тему «Предобработка текста»

Выполнил: Ван Пэй

Группа: ИУ5-22М

Москва — 2021 г.

Цель лабораторной работы:

Изучение методов preprocessing текстов.

Задание:

Для произвольного предложения или текста решите следующие задачи:

1. Токенизация.
2. Частеречная разметка.
3. Лемматизация.
4. Выделение (распознавание) именованных сущностей.

```
with open('../input/hp-rus/HP1RUS.TXT', 'r', encoding='cp1251') as f:  
    data = f.read()
```

```
data[:500]
```

```
' \n Джоан Кэтлин Роулинг \n Гарри Поттер  
и философский камень \n \n Глава первая \n  
Мальчик, который выжил \n \n Мистер и мисс  
ис Дарсли, из дома номер четыре по Бирючинн  
ому проезду, могли бы \n с гордостью сказать,  
что они, слава богу, совершенно нормальные  
люди. Они были \n бы последними, от кого вы мо  
гли бы ожидать участия в чем-нибудь странн  
ом и \n таинственном, потому что они совершен  
но не одобряли подобной чепухи. \n Мистер  
Дарсли работал директором фирмы Граннинг  
з, которая выпускала'
```

```
data = data.replace('\n', '')
```

```
data[:100]
```

```
'      Джоан Кэтлин Роулинг      Гарри Поттер и  
философский камень      Глава первая      Маль  
чик, кот'
```

- ****Токенизация**

```
from nltk.tokenize import sent_tokenize
sentences = sent_tokenize(data)
```

```
[['t'], 'ttt']
```

```
[['t'], 'ttt']
```

```
type(sentences)
```

```
list
```

```
sentences[:10]
```

```
[
    Джоан Кэтлин Роулинг Гарри Поттер и
    философский камень Глава первая Маль
    чик, который выжил Мистер и миссис Дарс
    ли, из дома номер четыре по Бирючинному про
    езду, могли бы с гордостью сказать, что они,
    слава богу, совершенно нормальные люди.',
    'Они были бы последними, от кого вы могли б
    ы ожидать участия в чем-нибудь странном и т
    аинственным, потому что они совершенно не
    одобряли подобной чепухи.',
    'Мистер Дарсли работал директором фирмы Г
    раннингз, которая выпускала сверла.',
    'Он был крупный, крепкий мужчина с очень ко
    роткой шеей и очень большими усами.',
    'Миссис Дарсли была худой и светловолосой,
    зато ее шеи с лихвой хватило бы на двоих.',
    'Этим подарком природы она усердно пользо
    валась, большую часть времени шпионя за со
    седями через садовые изгороди.',
    'У Дарсли пос сы Дарсли и они были введены
```

```
from string import punctuation
punctuation
```

```
'!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~'
```

```
# sentences[2].lower()
```

```
table = str.maketrans("", "", punctuation)
simple_preprocess = [s.lower().translate(table) for s in sentences]
```

```
sentences[10]
```

```
'Дарсли содрогались от мысли, что скажут
седи, если увидят Поттеров.'
```

```
simple_preprocess[10]
```

дарсли содроғались от мысли что скажут
седи если увидят потерров'

- **Токенизация

```
: from nltk.tokenize import word_tokenize
tokenized_words = [word_tokenize(s) for s in simple_preprocess]
```

```
: tokenized_words[2:5]
```

, П В О Н X,]]
 , Н S,
 , Q P I,
 , X B S L H H O,
 , П H X B O H,
 , C,
 , П E H,
 , E E,
 , Z S L O,
 , C B E L H O B O H O C O H,
 , H,
 , X λ Π O H,
 , Q P H S,
 , H S B C H H,
 [, H H C C H C,
 , λ C S H H,]
 , Q O H P H H H H,
 , O A E H P,
 , H,
 , П E E H,
 , K O B O L K O H,
 , O A E H P,
 , C,
 , H λ ж H H H S,
 , K B E H K H H,
 , K B λ H H P H,
 , Q P H,
 [, O H,
 , C B E b H S,]
 , B P H λ C K S H S,
 , K O L O B S H,
 , L b S H H H H L Z,
 , Φ H b H P I,
 , H H b E K L O b O H,
 , b S Q O L S H,
 , H S B C H H,
 [[, H H C L E B,]

['джоан', 'кэтлин', 'роулинг', 'гарри', 'поттер']

```
','.join([word for sent in tokenized_words for word in sent][:5])
```

‘джоан, кэтлин, роулинг, гарри, поттер’

```
from wordcloud import WordCloud
import matplotlib.pyplot as plt
```

```
text = ','.join([word for sent in tokenized_words for word in sent])
```

```
wordcloud = WordCloud(height=1000, width=2000, background_color='white', colormap=
plt.figure(figsize=(20, 20))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis("off")
plt.show()
```



```
from nltk.corpus import stopwords

import nltk
nltk.download("stopwords")

russian_stopwords = stopwords.words("russian")
print(russian_stopwords)
```


Стемминг

Лемматизация

```
!pip install pymorphy2
```

Collecting pymorphy2

Downloading pymorphy2-0.9.1-py3-none-any.whl (55 kB)

 55 kB

286 kB/s eta 0:00:011

Collecting dawg-python>=0.7.1

Downloading DAWG_Python-0.7.2-py2.py3-none-any.whl (11 kB)

Requirement already satisfied: docopt>=0.6 in /opt/conda/lib/python3.7/site-packages (from pymorphy2) (0.6.2)

Collecting pymorphy2-dicts-ru<3.0, >=2.4

Downloading pymorphy2_dicts_ru-2.4.417127.4579844-py2.py3-none-any.whl (8.2 MB)

3.6 MB/s eta 0:00:01

Installing collected packages: dawg-python, pymorphy2-dicts-ru, pymorphy2

Successfully installed dawg-python-0.7.2 pymorphy2-0.9.1 pymorphy2-dicts-ru-2.4.417127.4579844

```
WARNING: You are using pip version 20.1.1; however, version 20.2.4 is available.
```

You should consider upgrading via the `'/opt/conda/bin/python3.7 -m pip install --upgrade pip'` command.

```
import pymorphy2
morph = pymorphy2.MorphAnalyzer()
```

```
removed stopwords[2766]
```

‘ д о р о г о й ’

```
morph.parse(removed_stopwords[2766])[0].normal_form
```

‘ д о р о г а ’

```
lemmatized = [morph.parse(w)[0].normal_form for w in removed_stopwords]
```

```
text = ', '.join(lemmatized)
```

```
wordcloud = WordCloud(height=1000, width=2000, background_color='white', colormap=
plt.figure(figsize=(20, 20))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis("off")
plt.show()
```



```
!wget https://storage.yandexcloud.net/natasha-navec/packs/navec_news_v1_1B_250K_
!wget https://storage.yandexcloud.net/natasha-slovnet/packs/slovnet_ner_news_v1.
```

```
--2020-11-11 09:51:41-- https://storage.yandexcloud.net/natasha-navec/packs/n
avec_news_v1_1B_250K_300d_100q.tar
Resolving storage.yandexcloud.net (storage.yandexcloud.net)... 213.180.193.24
3, 2a02:6b8::1d9
Connecting to storage.yandexcloud.net (storage.yandexcloud.net)|213.180.193.24
3|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 26634240 (25M) [application/x-tar]
Saving to: 'navec_news_v1_1B_250K_300d_100q.tar'
```

```
navec_news_v1_1B_25 100%[=====>] 25.40M 32.6MB/s in 0.8s
```

```
2020-11-11 09:51:42 (32.6 MB/s) - 'navec_news_v1_1B_250K_300d_100q.tar' save
d [26634240/26634240]
```

```
--2020-11-11 09:51:43-- https://storage.yandexcloud.net/natasha-slovnet/pack
s/slovnet_ner_news_v1.tar
Resolving storage.yandexcloud.net (storage.yandexcloud.net)... 213.180.193.24
3, 2a02:6b8::1d9
Connecting to storage.yandexcloud.net (storage.yandexcloud.net)|213.180.193.24
3|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 2385920 (2.3M) [application/octet-stream]
Saving to: 'slovnet_ner_news_v1.tar'
```

```
slovnet_ner_news_v1 100%[=====>] 2.28M 6.85MB/s in 0.3s
```

```
2020-11-11 09:51:44 (6.85 MB/s) - 'slovnet_ner_news_v1.tar' saved [2385920/2
385920]
```

```
from navec import Navec
from slovnet import NER
```

```
navec = Navec.load('./navec_news_v1_1B_250K_300d_100q.tar')
ner = NER.load('./slovnet_ner_news_v1.tar')
ner.navec(navec)
```

```
markup = ner(data)
```

```
persons = [data[s.start:s.stop] for s in markup.spans if s.type=='PER']
locations = [data[s.start:s.stop] for s in markup.spans if s.type=='LOC']
organizations = [data[s.start:s.stop] for s in markup.spans if s.type=='ORG']
```

```
# persons
```

```
f = {}  
f['key1'] = 'ppppp'  
f['key2'] = 5  
f['key3'] = [4, 5, 7]
```

```
locations[:10]
```

```
['Бирючинному проезду',  
'Бирючинного проезда',  
'Кенте',  
'Йоркшире',  
'Данди',  
'Великобританией',  
'Бирючинный проезд',  
'Бирючинного проезда',  
'Бирючинном проезде',  
'Кенте']
```

```
def cnt_objects(objects_list: list) -> dict:  
    cnt = {}  
  
    for o in objects_list:  
        if o not in cnt.keys():  
            cnt[o] = 1  
        else:  
            cnt[o] += 1  
    return cnt
```

```
sorted(cnt_objects(locations).items(), key=lambda x: x[1], reverse=True)[:5]
```

```
[('Хогвартсе', 23),  
( 'Хогвартс', 11),  
( 'Гринготтс', 11),  
( 'Лондон', 7),  
( 'Слайзерин', 7)]
```

```
sorted(cnt_objects(persons).items(), key=lambda x: x[1], reverse=True)[:5]
```

```
[('Гарри', 1207),  
( 'Рон', 322),  
( 'Хагрид', 238),  
( 'Гермиона', 162),  
( 'Дадли', 138)]
```

Список литературы

[1] Гапанюк Ю. Е. Лабораторная работа «Обработка признаков (часть2)»

[Электронный ресурс]

https://github.com/ugapanyuk/ml_course_2021/wiki/LAB_MMO__FEATURES