

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



**Лабораторная работа №5**  
**По курсу «Методы машинного обучения»**

**«Предобработка текста»**

**ИСПОЛНИТЕЛЬ:**

Чичикин Тимофей Дмитриевич  
Группа ИУ5-25М

---

**ПРОВЕРИЛ:**

Гапанюк Ю.Е.

---

Цель работы:

Изучение методов предобработки текста.

Задание:

Для произвольного предложения или текста решите следующие задачи:

- Токенизация.
- Частеречная разметка.
- Лемматизация.
- Выделение (распознавание) именованных сущностей.
- Разбор предложения.

Описание задания:

Для выполнения лабораторной работы возьмём фразу: «Предобработка данных в XML файле».

Выполнение работы:

1. Токенизация. NLTK
2. Частеречная разметка. Natasha
3. Лемматизация. Natasha
4. Выделение именованных сущностей. Natasha
5. Разбор предложения. Natasha

Вывод:

Была проделана работа по изучению методов предобработки текста, все задачи были выполнены.



```
In [ ]: !pip install numpy pandas scikit-surprise sklearn seaborn matplotlib spacy nlt
```

```
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (1.19.5)
Requirement already satisfied: pandas in /usr/local/lib/python3.7/dist-packages (1.1.5)
Requirement already satisfied: scikit-surprise in /usr/local/lib/python3.7/dist-packages (1.1.1)
Requirement already satisfied: sklearn in /usr/local/lib/python3.7/dist-packages (0.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.7/dist-packages (0.11.1)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.7/dist-packages (3.2.2)
Requirement already satisfied: spacy in /usr/local/lib/python3.7/dist-packages (2.2.4)
Requirement already satisfied: nltk in /usr/local/lib/python3.7/dist-packages (3.2.5)
Requirement already satisfied: navec in /usr/local/lib/python3.7/dist-packages (0.10.0)
Requirement already satisfied: slovnet in /usr/local/lib/python3.7/dist-packages (0.5.0)
Collecting natasha
  Downloading https://files.pythonhosted.org/packages/51/8e/ab0745100be276750fb6b8858c6180a1756696572295a74eb5aea77f3bbd/natasha-1.4.0-py3-none-any.whl (34.4MB)
    |████████████████████████████████████████| 34.4MB 109kB/s
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-packages (from pandas) (2.8.1)
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/dist-packages (from pandas) (2018.9)
Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-surprise) (1.4.1)
Requirement already satisfied: six>=1.10.0 in /usr/local/lib/python3.7/dist-packages (from scikit-surprise) (1.15.0)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-surprise) (1.0.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.7/dist-packages (from sklearn) (0.22.2.post1)
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib) (2.4.7)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.7/dist-packages (from matplotlib) (0.10.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib) (1.3.1)
Requirement already satisfied: srsly<1.1.0,>=1.0.2 in /usr/local/lib/python3.7/dist-packages (from spacy) (1.0.5)
Requirement already satisfied: thinc==7.4.0 in /usr/local/lib/python3.7/dist-packages (from spacy) (7.4.0)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.7/dist-packages (from spacy) (2.23.0)
Requirement already satisfied: plac<1.2.0,>=0.9.6 in /usr/local/lib/python3.7/dist-packages (from spacy) (1.1.3)
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.7/dist-packages (from spacy) (2.0.5)
Requirement already satisfied: catalogue<1.1.0,>=0.0.7 in /usr/local/lib/python3.7/dist-packages (from spacy) (1.0.0)
```

```

3.7/dist-packages (from spacy) (1.0.0)
Requirement already satisfied: setuptools in /usr/local/lib/python3.7/dist-pack
ages (from spacy) (56.1.0)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python
3.7/dist-packages (from spacy) (3.0.5)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/pyth
on3.7/dist-packages (from spacy) (1.0.5)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.7/
dist-packages (from spacy) (4.41.1)
Requirement already satisfied: blis<0.5.0,>=0.4.0 in /usr/local/lib/python3.7/d
ist-packages (from spacy) (0.4.1)
Requirement already satisfied: wasabi<1.1.0,>=0.4.0 in /usr/local/lib/python
3.7/dist-packages (from spacy) (0.8.2)
Requirement already satisfied: razdel in /usr/local/lib/python3.7/dist-packages
 (from slovnet) (0.5.0)
Collecting ipymarkup>=0.8.0
  Downloading https://files.pythonhosted.org/packages/bf/9b/bf54c98d50735a4a7c8
4c71e92c5361730c878ebfe903d2c2d196ef66055/ipymarkup-0.9.0-py3-none-any.whl
Collecting yargy>=0.14.0
  Downloading https://files.pythonhosted.org/packages/d3/46/bc1a17200a55f4b0608
f39ac64f1840fd4a52f9eeea462d9afecbf71246b/yargy-0.15.0-py3-none-any.whl (41kB)
|████████████████████████████████████████| 51kB 5.6MB/s
Collecting pymorphy2
  Downloading https://files.pythonhosted.org/packages/07/57/b2ff2fae3376d4f3c69
7b9886b64a54b476e1a332c67eee9f88e7f1ae8c9/pymorphy2-0.9.1-py3-none-any.whl (55k
B)
|████████████████████████████████████████| 61kB 6.4MB/s
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-pa
ckages (from requests<3.0.0,>=2.13.0->spacy) (2.10)
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr/
local/lib/python3.7/dist-packages (from requests<3.0.0,>=2.13.0->spacy) (1.2
4.3)
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/di
st-packages (from requests<3.0.0,>=2.13.0->spacy) (3.0.4)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/d
ist-packages (from requests<3.0.0,>=2.13.0->spacy) (2020.12.5)
Requirement already satisfied: importlib-metadata>=0.20; python_version < "3.8"
in /usr/local/lib/python3.7/dist-packages (from catalogue<1.1.0,>=0.0.7->spacy)
(4.0.1)
Collecting intervaltree>=3
  Downloading https://files.pythonhosted.org/packages/50/fb/396d568039d21344639
db96d940d40eb62bef704ef849b27949ded5c3bb/intervaltree-3.1.0.tar.gz
Collecting pymorphy2-dicts-ru<3.0,>=2.4
  Downloading https://files.pythonhosted.org/packages/3a/79/bea0021eeb7eeefde22
ef9e96badf174068a2dd20264b9a378f2be1cdd9e/pymorphy2_dicts_ru-2.4.417127.457984
4-py2.py3-none-any.whl (8.2MB)
|████████████████████████████████████████| 8.2MB 19.5MB/s
Collecting dawg-python>=0.7.1
  Downloading https://files.pythonhosted.org/packages/6a/84/ff1ce2071d4c650ec85
745766c0047ccc3b5036f1d03559fd46bb38b5eeb/DAWG_Python-0.7.2-py2.py3-none-any.wh
l
Requirement already satisfied: docopt>=0.6 in /usr/local/lib/python3.7/dist-pac
kages (from pymorphy2->natasha) (0.6.2)
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-packa

```

```

ges (from importlib-metadata>=0.20; python_version < "3.8"->catalogu
e<1.1.0,>=0.0.7->spacy) (3.4.1)
Requirement already satisfied: typing-extensions>=3.6.4; python_version < "3.8"
in /usr/local/lib/python3.7/dist-packages (from importlib-metadata>=0.20; pytho
n_version < "3.8"->catalogue<1.1.0,>=0.0.7->spacy) (3.7.4.3)
Requirement already satisfied: sortedcontainers<3.0,>=2.0 in /usr/local/lib/pyt
hon3.7/dist-packages (from intervaltree>=3->ipymarkup>=0.8.0->natasha) (2.4.0)
Building wheels for collected packages: intervaltree
  Building wheel for intervaltree (setup.py) ... done
  Created wheel for intervaltree: filename=intervaltree-3.1.0-py2.py3-none-an
y.whl size=26102 sha256=966842acb52b8dbc4ae3d66fd3f7e7d9f5db37be0860433922f8fc9
2f2b7d6a2
  Stored in directory: /root/.cache/pip/wheels/f3/f2/66/e9c30d3e9499e65ea2fa0d0
7c002e64de63bd0adaa49c445bf
Successfully built intervaltree
Installing collected packages: intervaltree, ipymarkup, pymorphy2-dicts-ru, daw
g-python, pymorphy2, yargy, natasha
  Found existing installation: intervaltree 2.1.0
    Uninstalling intervaltree-2.1.0:
      Successfully uninstalled intervaltree-2.1.0
Successfully installed dawg-python-0.7.2 intervaltree-3.1.0 ipymarkup-0.9.0 nat
asha-1.4.0 pymorphy2-0.9.1 pymorphy2-dicts-ru-2.4.417127.4579844 yargy-0.15.0

```

## Токенизация. NLTK

```

In [ ]: import nltk
        nltk.download('punkt')
        text1 = 'Предобработка данных в XML файле.'
        text2 = 'Меня зовут Бонд. Джеймс Бонд'

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!

In [ ]: from nltk import tokenize
        dir(tokenize)[:18]

```

```
Out[ ]: ['BlanklineTokenizer',
        'LineTokenizer',
        'MWETokenizer',
        'PunktSentenceTokenizer',
        'RegexpTokenizer',
        'ReppTokenizer',
        'SEExprTokenizer',
        'SpaceTokenizer',
        'StanfordSegmenter',
        'TabTokenizer',
        'TextTilingTokenizer',
        'ToktokTokenizer',
        'TreebankWordTokenizer',
        'TweetTokenizer',
        'WhitespaceTokenizer',
        'WordPunctTokenizer',
        '__builtins__',
        '__cached__']
```

```
In [ ]: nltk_tk_1 = nltk.WordPunctTokenizer()
        nltk_word = nltk_tk_1.tokenize(text1)
        print(nltk_word)
```

```
['Предобработка', 'данных', 'в', 'XML', 'файле', '.']
```

```
In [ ]: # Токенизация по предложениям
        nltk_tk_sents = nltk.tokenize.sent_tokenize(text1)
        print(len(nltk_tk_sents))
        nltk_tk_sents
```

```
1
```

```
Out[ ]: ['Предобработка данных в XML файле.']
```

## Частеречная разметка. Natasha

```
In [ ]: from navec import Navec
        from slovnet import Morph
```

```
In [ ]: from google.colab import drive
        drive.mount('/content/gdrive')
```

```
Mounted at /content/gdrive
```

```
In [ ]: navec = Navec.load('/content/gdrive/My Drive/MM0/navec_news_v1_1B_250K_300d_10
        n_morph = Morph.load('/content/gdrive/My Drive/MM0/slovnet_morph_news_v1.tar',
```

```
In [ ]: morph_res = n_morph.navec(navec)
```

```
In [ ]: def print_pos(markup):
        for token in markup.tokens:
            print('{ } - {}'.format(token.text, token.tag))
```

```
In [ ]: n_text1_markup = list(_ for _ in n_morph.map(nltk_tokenize))
[print_pos(x) for x in n_text1_markup]
```

```
П - PROPN|Animacy=Anim|Case=Nom|Gender=Masc|Number=Sing
р - NOUN
е - X|Foreign=Yes
д - NOUN
о - X|Foreign=Yes
б - NOUN|Animacy=Inan|Case=Loc|Gender=Masc|Number=Sing
р - X|Foreign=Yes
а - CCONJ
б - PROPN
о - NOUN|Animacy=Inan|Case=Gen|Gender=Fem|Number=Sing
т - PRON|Animacy=Inan|Case=Loc|Gender=Neut|Number=Sing
к - ADP
а - X|Foreign=Yes
- PUNCT
д - NOUN|Animacy=Inan|Case=Gen|Gender=Masc|Number=Sing
а - CCONJ
н - X|Foreign=Yes
н - X|Foreign=Yes
ы - X|Foreign=Yes
х - X|Foreign=Yes
- PUNCT
в - X|Foreign=Yes
- PUNCT
Х - X|Foreign=Yes
М - PROPN|Foreign=Yes
Л - X|Foreign=Yes
- PUNCT
ф - X|Foreign=Yes
а - CCONJ
й - ADJ|Case=Nom|Degree=Pos|Gender=Masc|Number=Sing
л - X|Foreign=Yes
е - NOUN|Animacy=Inan|Case=Gen|Gender=Masc|Number=Sing
. - PUNCT
```

```
Out[ ]: [None]
```

## Лемматизация. Natasha

```
In [ ]: from natasha import Doc, Segmenter, NewsEmbedding, NewsMorphTagger, MorphVocab
```

```
In [ ]: def n_lemmatize(text):
    emb = NewsEmbedding()
    morph_tagger = NewsMorphTagger(emb)
    segmenter = Segmenter()
    morph_vocab = MorphVocab()
    doc = Doc(text)
    doc.segment(segmenter)
    doc.tag_morph(morph_tagger)
    for token in doc.tokens:
```



```
token.lemmatize(morph_vocab)
return doc
```

```
In [ ]: n_doc1 = n_lemmatize(text1)
        {_.text: _.lemma for _ in n_doc1.tokens}
```

```
Out[ ]: {'.': '.',
        'XML': 'xml',
        'Предобработка': 'предобработка',
        'в': 'в',
        'данных': 'данные',
        'файле': 'файл'}
```

```
In [ ]: n_doc2 = n_lemmatize(text2)
        {_.text: _.lemma for _ in n_doc2.tokens}
```

```
Out[ ]: {'.': '.', 'Бонд': 'бонд', 'Джеймс': 'джеймс', 'Меня': 'я', 'зовут': 'звать'}
```

## Выделение (распознавание) именованных сущностей. Natasha

```
In [ ]: from slovnet import NER
        from ipymarkup import show_span_ascii_markup as show_markup
```

```
In [ ]: ner = NER.load('/content/gdrive/My Drive/MM0/slovnet_ner_news_v1.tar')
```

```
In [ ]: ner_res = ner.navec(navec)
```

```
In [ ]: markup_ner2 = ner(text2)
```

```
In [ ]: markup_ner2
```

```
Out[ ]: SpanMarkup(
    text='Меня зовут Бонд. Джеймс Бонд',
    spans=[Span(
        start=11,
        stop=15,
        type='PER'
    ), Span(
        start=17,
        stop=28,
        type='PER'
    )]
)
```

```
In [ ]: show_markup(markup_ner2.text, markup_ner2.spans)
```

```
Меня зовут Бонд. Джеймс Бонд
      PER— PER—————
```

# Разбор предложения. Natasha

```
In [ ]: from natasha import NewsSyntaxParser
```

```
In [ ]: emb = NewsEmbedding()  
syntax_parser = NewsSyntaxParser(emb)
```

```
In [ ]: n_doc1.parse_syntax(syntax_parser)  
n_doc1.sents[0].syntax.print()
```

```
└─ Предобработка amod  
   └─ данных  
      └─ в case  
         └─ XML  
            └─ файле  
               .
```

```
In [ ]: n_doc2.parse_syntax(syntax_parser)  
n_doc2.sents[0].syntax.print()
```

```
└─ Меня obj  
   └─ зовут  
      └─ Бонд xcomp  
         └─ . punct
```