



**Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)**

**Факультет «Информатика и системы управления»
Кафедра ИУ5 «Системы обработки информации и управления»**

**Отчёт по лабораторной работе
по дисциплине «Методы Машинного Обучения»**

Выполнил:
студент группы № ИУ5-21М
Кучеренко Михаил Александрович
_____, _____

Проверил:
к.т.н., доц., Ю.Е. Гапанюк
_____, _____

2021 г.

Лабораторная работа №5

Предобработка текста.

Набор данных - Стихотворение "Простая суть" Геннадий Эсса.

Install

```
1 pip3 install spacy
2 python3 -m spacy download ru_core_news_sm
```

Задание:

Для произвольного предложения или текста решите следующие задачи:

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

```
1 text = '''Лишь тот способен добрым быть,
2 Кому пришлось от злости выть.
3 Лишь тот узнает дружбы цену,
4 Кто испытал друзей измену.
5 И лишь тому любовь дана,
6 Кто ненависть познал сполна.
7 Кто в жизни повстречал врагов,
8 Кретинов или дураков...
9 Где разошлись дороги наши
10 Из-за ненужных пустяков.
11 Лишь тот способен дорожить
12 Тем, для чего так стоит жить.
13 И верить, и мечтать, и ждать,
14 От радости своей кричать!
15 Искать, поверить и любить.
16 Для этого нам стоит быть!
17 Поэтому ищи, терзай,
18 И ни о чем не забывай.
19 Не жди к себе пустых признаний:
20 Не будет в этом оправданий.
21 Лишь тот способен все понять,
22 И строки эти прочитать.
23 Кому не все равно по жизни
24 Простую истину познать...
25 (С) Геннадий Эсса
26 '''
27
28 limit = 18
```

```
1 from spacy.lang.ru import Russian
2 import spacy
```

Токенизация

```
1 nlp = spacy.load('ru_core_news_sm')
2 spacy_text = nlp(text)
```

```
1 for token in spacy_text[:limit]:
2     print(f'{token.text:20} - {token.pos_:7} / {token.dep_}')
```

```
1 Лишь          - PART    / advmod
2 тот           - DET     / nsubj
3 способен      - ADJ     / ROOT
4 добрым        - ADJ     / xcomp
5 быть          - VERB    / xcomp
6 ,             - PUNCT   / punct
7
8              - SPACE    / mark
```

9	Кому	- PRON	/ iobj
10	пришлось	- VERB	/ advcl
11	от	- ADP	/ case
12	злости	- NOUN	/ obl
13	выть	- VERB	/ csubj
14	.	- PUNCT	/ punct
15			
16		- SPACE	/ advmod
17	Лишь	- PART	/ advmod
18	тот	- DET	/ nsubj
19	узнает	- VERB	/ ROOT
20	дружбы	- NOUN	/ iobj

```

1 | spacy_text = nlp(text.replace('\n', ' '))
2 |
3 | for token in spacy_text[:limit]:
4 |     print(f'{token.text:20} - {token.pos_:7} / {token.dep_}')

```

1	Лишь	- PART	/ advmod
2	тот	- DET	/ nsubj
3	способен	- ADJ	/ ROOT
4	добрым	- ADJ	/ xcomp
5	быть	- VERB	/ xcomp
6	,	- PUNCT	/ punct
7		- SPACE	/ nsubj
8	Кому	- PRON	/ iobj
9	пришлось	- VERB	/ conj
10	от	- ADP	/ case
11	злости	- NOUN	/ obl
12	выть	- VERB	/ nmod
13	.	- PUNCT	/ punct
14		- SPACE	/ nsubj
15	Лишь	- PART	/ advmod
16	тот	- DET	/ nsubj
17	узнает	- VERB	/ ROOT
18	дружбы	- NOUN	/ iobj

Лемматизация

```

1 | for token in spacy_text[:limit]:
2 |     print(token, token.lemma, token.lemma_)

```

1	Лишь	15085157332802532245	лишь
2	тот	9060620132269683821	тот
3	способен	13108567891086621765	способный
4	добрым	3413372487067799558	добрый
5	быть	13421914901540782342	быть
6	,	2593208677638477497	,
7		3443622981070210256	
8	Кому	4147452868732122456	кто
9	пришлось	12476235364299464058	прийтись
10	от	7547231311137123581	от
11	злости	16727616970683336905	злость
12	выть	10742055882863096977	выть
13	.	12646065887601541794	.
14		3443622981070210256	
15	Лишь	15085157332802532245	лишь
16	тот	9060620132269683821	тот
17	узнает	9829465616891537735	узнавать
18	дружбы	9180780978176652763	дружба

Выделение (распознавание) именованных сущностей

```

1 | for entity in spacy_text.ents:
2 |     print(entity.text, entity.label_)

```

1	Лишь	PER
2	Кретинов	PER
3	Геннадий Эсса	PER

Разбор предложения

```
1 | from spacy import displacy
```

```
1 | sentence = len(spacy_text)
2 | for token in spacy_text:
3 |     if token.text == '.':
4 |         sentence = token.i
5 |         break
6 | sentence += 1
```

```
1 | displacy.render(spacy_text, style='ent', jupyter=True)
```

Лишь тот способен добрым быть, Кому пришлось от злости выть. **Лишь PER** тот узнает дружбы цену, Кто испытал друзей измену. И лишь тому любовь дана, Кто ненависть познал сполна. Кто в жизни повстречал врагов, **Кретинов PER** или дураков.... Где разошлись дороги наши Из-за ненужных пустяков. Лишь тот способен дорожить Тем, для чего так стоит жить. И верить, и мечтать, и ждать, От радости своей кричать! Искать, поверить и любить. Для этого нам стоит быть! Поэтому ищи, терзай, И ни о чем не забывай. Не жди к себе пустых признаний: Не будет в этом оправданий. Лишь тот способен все понять, И строки эти прочитать. Кому не все равно по жизни Простую истину познать...(С) **Геннадий Эсса PER**

```
1 | displacy.render(spacy_text[:sentence], style='dep', jupyter=True)
```

Лишь
PART

тот
DET

способен
ADJ

добрым
ADJ

быть,
VERB

SPACE

Кому
PRON

пришлось
VERB

от
ADP

злости
NOUN

быть.
VERB

advmod

nsubj

xcomp

xcomp

nsubj

iobj

conj

case

obl

nmod

