

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



Лабораторная работа №3
по дисциплине
«Методы машинного обучения»

Выполнил:
студент группы ИУ5И-23М

Лю Цзычжан

Москва — 2025 г.

Цель лабораторной работы: изучение методов предобработки текстов.

Требования к отчету:

Отчет по лабораторной работе должен содержать:

1. титульный лист;
2. описание задания;
3. текст программы;
4. экранные формы с примерами выполнения программы.

Задание:

Для произвольного предложения или текста решите следующие задачи:

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

```
6. import spacy
7.
8. # Загружаем модель spacy для английского языка
9. nlp = spacy.load("en_core_web_sm")
10.
11. # Пример произвольного предложения
12. text = "Apple is looking at buying U.K. startup for $1 billion."
13.
14. # Пропускаем текст через пайплайн spacy
15. doc = nlp(text)
16.
```

```
17.# Токенизация
18.print("Tokens:")
19.for token in doc:
20.    print(f"{token.text}", end=" | ")
21.print("\n")
22.
23.# Частеречная разметка
24.print("POS tagging:")
25.for token in doc:
26.    print(f"{token.text}: {token.pos_} ({token.tag_})")
27.print("\n")
28.
29.# Лемматизация
30.print("Lemmas:")
31.for token in doc:
32.    print(f"{token.text} -> {token.lemma_}")
33.print("\n")
34.
35.# Именованные сущности
36.print("Named Entities:")
37.for ent in doc.ents:
38.    print(f"{ent.text} ({ent.label_})")
39.print("\n")
40.
41.# Синтаксический разбор
42.print("Dependency Parsing:")
43.for token in doc:
44.    print(f"{token.text} --> {token.dep_} --> {token.head.text}")
45.
```

```
(venv) PS E:\BMSTU\Т\jqxx\MMO\lp5> python -u "e:\BMSTU\Т\jqxx\MMO\lp5\lab5.py"
```

```
Tokens:
```

```
Apple | is | looking | at | buying | U.K. | startup | for | $ | 1 | billion | . |
```

```
POS tagging:
```

```
Apple: PROPN (NNP)
```

```
is: AUX (VBZ)
```

```
looking: VERB (VBG)
```

```
at: ADP (IN)
```

```
buying: VERB (VBG)
```

```
U.K.: PROPN (NNP)
```

```
startup: VERB (VBD)
```

```
for: ADP (IN)
```

```
$: SYM ($)
```

```
1: NUM (CD)
```

```
billion: NUM (CD)
```

```
.: PUNCT (.)
```

```
Lemmas:
```

```
Apple -> Apple
```

```
is -> be
```

```
looking -> look
```

```
at -> at
```

```
buying -> buy
```

```
U.K. -> U.K.
```

```
startup -> startup
```

```
for -> for
```

```
$ -> $
```

```
1 -> 1
```

```
billion -> billion
```

```
. -> .
```

```
. -> .
```

Named Entities:

Apple (ORG)

U.K. (GPE)

\$1 billion (MONEY)

Lemmas:

Apple -> Apple

is -> be

looking -> look

at -> at

buying -> buy

U.K. -> U.K.

startup -> startup

for -> for

\$ -> \$

1 -> 1

billion -> billion

```
. -> .
```

Named Entities:

Apple (ORG)

U.K. (GPE)

\$1 billion (MONEY)

Apple -> Apple

is -> be

looking -> look

at -> at

buying -> buy

U.K. -> U.K.

startup -> startup

for -> for

\$ -> \$

1 -> 1

billion -> billion

```
. -> .
```

Named Entities:
Apple (ORG)
U.K. (GPE)
\$1 billion (MONEY)
looking -> look
at -> at
buying -> buy
U.K. -> U.K.
startup -> startup
for -> for
\$ -> \$
1 -> 1
billion -> billion
. -> .

Named Entities:
Apple (ORG)
U.K. (GPE)
\$1 billion (MONEY)
for -> for
\$ -> \$
1 -> 1
billion -> billion
. -> .

Named Entities:
Apple (ORG)
U.K. (GPE)
\$1 billion (MONEY)
billion -> billion
. -> .

Named Entities:
Apple (ORG)
U.K. (GPE)
\$1 billion (MONEY)

Named Entities:

Apple (ORG)

U.K. (GPE)

\$1 billion (MONEY)

Named Entities:

Apple (ORG)

U.K. (GPE)

\$1 billion (MONEY)

Dependency Parsing:

Apple (ORG)

U.K. (GPE)

\$1 billion (MONEY)

Dependency Parsing:

\$1 billion (MONEY)

Dependency Parsing:

Dependency Parsing:

Apple --> nsubj --> looking

Apple --> nsubj --> looking

is --> aux --> looking

is --> aux --> looking

looking --> ROOT --> looking

looking --> ROOT --> looking

at --> prep --> looking

at --> prep --> looking

buying --> pcomp --> at

U.K. --> nsubj --> startup

startup --> ccomp --> buying

startup --> ccomp --> buying

for --> prep --> startup

\$ --> quantmod --> billion

1 --> compound --> billion

billion --> pobj --> for

. --> punct --> looking