

# Лабораторная работа №5

Ишков Денис Олегович, ИУ5-24М, 2021г.

Тема: Предобработка текста

**Цель: Изучение методов предобработки текстов..**

Требования к отчету: Отчет по лабораторной работе должен содержать:

- титульный лист;
- описание задания;
- текст программы;
- экранные формы с примерами выполнения программы. Задание:

Для произвольного предложения или текста решите следующие задачи:

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

## Датасет

Бинарная классификация текста

<https://www.kaggle.com/blackmoon/russian-language-toxic-comments>  
(<https://www.kaggle.com/blackmoon/russian-language-toxic-comments>)

In [1]:

```
# загрузка датасета
!pip install wldhx.yadisk-direct
!curl -L $(yadisk-direct https://disk.yandex.ru/d/wedARfrtMn-Y-Q) -o labeled.csv
```

Collecting wldhx.yadisk-direct

Downloading wldhx.yadisk\_direct-0.0.6-py3-none-any.whl (4.5 kB)

Requirement already satisfied: requests in /opt/conda/lib/python3.7/site-packages (from wldhx.yadisk-direct) (2.25.1)

Requirement already satisfied: certifi&gt;=2017.4.17 in /opt/conda/lib/python3.7/site-packages (from requests-&gt;wldhx.yadisk-direct) (2020.12.5)

Requirement already satisfied: chardet&lt;5,&gt;=3.0.2 in /opt/conda/lib/python3.7/site-packages (from requests-&gt;wldhx.yadisk-direct) (4.0.0)

Requirement already satisfied: urllib3&lt;1.27,&gt;=1.21.1 in /opt/conda/lib/python3.7/site-packages (from requests-&gt;wldhx.yadisk-direct) (1.26.4)

Requirement already satisfied: idna&lt;3,&gt;=2.5 in /opt/conda/lib/python3.7/site-packages (from requests-&gt;wldhx.yadisk-direct) (2.10)

Installing collected packages: wldhx.yadisk-direct

Successfully installed wldhx.yadisk-direct-0.0.6

% Total	% Received	% Xferd	Average Speed	Time	Time	Time	Cu
rrrent							

			Dload	Upload	Total	Spent	Left Sp

eed

0	0	0	0	0	0	0	0	--:--:--	0:00:01	--:--:--
---	---	---	---	---	---	---	---	----------	---------	----------

0

100	4560k	100	4560k	0	0	1161k	0	0:00:03	0:00:03	--:--:-- 1
-----	-------	-----	-------	---	---	-------	---	---------	---------	------------

720k

## Импорт нужных библиотек

In [143]:

```

import nltk
import spacy
import numpy as np
from tqdm.notebook import tqdm
nltk.download('punkt')
from nltk import tokenize
import re
import pandas as pd
from sklearn.model_selection import train_test_split
import nltk
import string
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import SnowballStemmer
from sklearn.pipeline import Pipeline
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import ComplementNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn.metrics import precision_score, recall_score, precision_recall_curve, clas
sification_report
from matplotlib import pyplot as plt
from sklearn.metrics import plot_precision_recall_curve
import numpy as np
from sklearn.model_selection import GridSearchCV

```

[nltk\_data] Downloading package punkt to /usr/share/nltk\_data...

[nltk\_data] Package punkt is already up-to-date!

## Анализ и обработка выбросов в данных

In [4]:

```

df = pd.read_csv("labeled.csv", sep=",")
df.describe()

```

Out[4]:

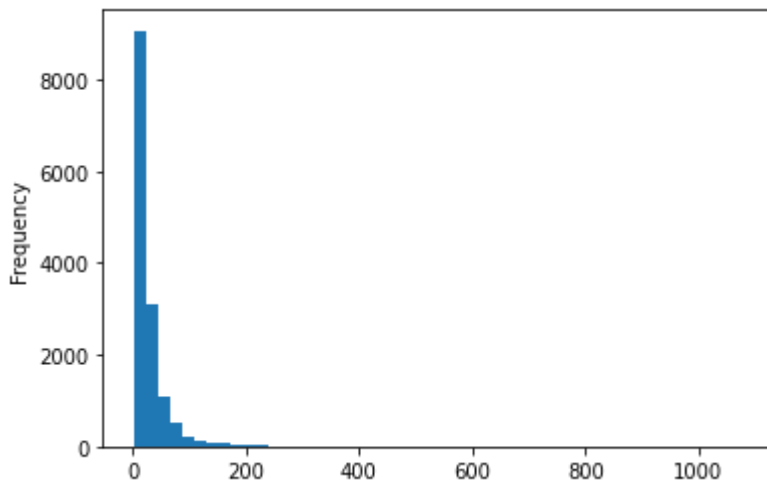
	toxic
count	14412.000000
mean	0.334860
std	0.471958
min	0.000000
25%	0.000000
50%	0.000000
75%	1.000000
max	1.000000

In [5]:

```
df.comment.str.split(' ').apply(len).plot(kind='hist', bins=50)
```

Out[5]:

&lt;AxesSubplot:ylabel='Frequency'&gt;



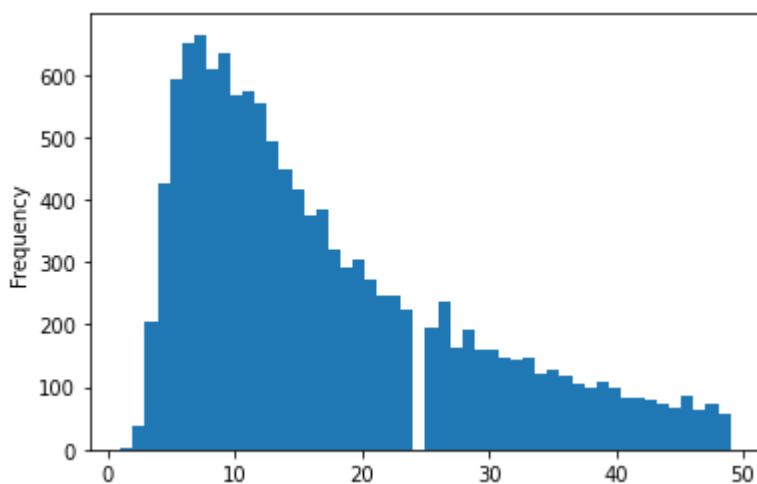
Как видно из гистограммы, количество слов сообщений в данных распределено по экспоненциальному закону. Уберём "хвост"

In [6]:

```
df = df[df.comment.str.split(' ').apply(len) < 50].copy()
df.comment.str.split(' ').apply(len).plot(kind='hist', bins=50)
```

Out[6]:

&lt;AxesSubplot:ylabel='Frequency'&gt;



In [7]:

```
df["toxic"] = df["toxic"].apply(int)
df["toxic"].value_counts()
```

Out[7]:

```
0    8102
1    4419
Name: toxic, dtype: int64
```

# 1. Токенизация

In [9]:

```
print('Токенизаторы NLTK')
for i in dir(tokenize)[:16]:
    print(i)
```

Токенизаторы NLTK  
BlanklineTokenizer  
LineTokenizer  
MWETokenizer  
PunktSentenceTokenizer  
RegexTokenizer  
ReppTokenizer  
SEExprTokenizer  
SpaceTokenizer  
StanfordSegmenter  
StanfordTokenizer  
TabTokenizer  
TextTilingTokenizer  
ToktokTokenizer  
TreebankWordTokenizer  
TweetTokenizer  
WhitespaceTokenizer

In [67]:

```
comment = df.comment.values[222]
print(comment)
cleaned_comment = re.sub('[^а-яА-Яа-зА-З0-9 \n\.]', '', comment)
print(cleaned_comment)
```

Бросила. Во время каждого занятия получала небольшие травмы, в последний раз при выполнении Сурьи Намаскар растянула связки на руке, боль была такая, что в первые секунды на перелом подумала. Теперь дома занимаюсь без напряжения со стороны тренера.

Бросила. Во время каждого занятия получала небольшие травмы в последний раз при выполнении Сурьи Намаскар растянула связки на руке боль была такая что в первые секунды на перелом подумала. Теперь дома занимаюсь без напряжения со стороны тренера.

In [22]:

```
tokenizer_wp = nltk.WhitespaceTokenizer()  
tokens = tokenizer_wp.tokenize(cleaned_comment)  
tokens
```

Out[22]:

```
['Бросила.',  
'Во',  
'время',  
'каждого',  
'занятия',  
'получала',  
'небольшие',  
'травмы',  
'в',  
'последний',  
'раз',  
'при',  
'выполнении',  
'Сурьи',  
'Намаскар',  
'растянула',  
'связки',  
'на',  
'руке',  
'боль',  
'была',  
'такая',  
'что',  
'в',  
'первые',  
'секунды',  
'на',  
'перелом',  
'подумала.',  
'Теперь',  
'дома',  
'занимаюсь',  
'без',  
'напряга',  
'со',  
'стороны',  
'тренера.']
```

In [23]:

```
tokenizer_wp = nltk.WordPunctTokenizer()
tokens = tokenizer_wp.tokenize(cleaned_comment)
tokens
```

Out[23]:

```
['Бросила',
 '.',
 'Во',
 'время',
 'каждого',
 'занятия',
 'получала',
 'небольшие',
 'травмы',
 'в',
 'последний',
 'раз',
 'при',
 'выполнении',
 'Сурьи',
 'Намаскар',
 'растянула',
 'связки',
 'на',
 'руке',
 'боль',
 'была',
 'такая',
 'что',
 'в',
 'первые',
 'секунды',
 'на',
 'перелом',
 'подумала',
 '.',
 'Теперь',
 'дома',
 'занимаюсь',
 'без',
 'напряга',
 'со',
 'стороны',
 'тренера',
 '.']
```

Делаем вывод, что WordPunctTokenizer подходит лучше

## 2. Частеречная разметка (POStagging)

In [27]:

```
!pip install natasha
```

Collecting natasha

Downloading natasha-1.4.0-py3-none-any.whl (34.4 MB)

|██| 34.4 MB 25.3 MB/s eta 0:00:01

Requirement already satisfied: ipymarkup>=0.8.0 in /opt/conda/lib/python3.7/site-packages (from natasha) (0.9.0)

Collecting slovnet>=0.3.0

Downloading slovnet-0.5.0-py3-none-any.whl (49 kB)

|██████████████████████████████████████| 49 kB 3.5 MB/s eta 0:00:01

Collecting pymorphy2

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

|██████████████████████████████████████| 55 kB 2.1 MB/s eta 0:00:01

Collecting razdel>=0.5.0

Downloading razdel-0.5.0-py3-none-any.whl (21 kB)

Collecting navec>=0.9.0

Downloading navec-0.10.0-py3-none-any.whl (23 kB)

Collecting yargy>=0.14.0

Downloading yargy-0.15.0-py3-none-any.whl (41 kB)

|██████████████████████████████████████| 41 kB 72 kB/s s eta 0:00:01

Requirement already satisfied: intervaltree>=3 in /opt/conda/lib/python3.7/site-packages (from ipymarkup>=0.8.0->natasha) (3.1.0)

Requirement already satisfied: sortedcontainers<3.0,>=2.0 in /opt/conda/lib/python3.7/site-packages (from intervaltree>=3->ipymarkup>=0.8.0->natasha) (2.3.0)

Requirement already satisfied: numpy in /opt/conda/lib/python3.7/site-packages (from navec>=0.9.0->natasha) (1.19.5)

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

Downloading pymorphy2\_dicts\_ru-2.4.417127.4579844-py2.py3-none-any.whl (8.2 MB)

|██████████████████████████████████████| 8.2 MB 42.0 MB/s eta 0:00:01

Collecting dawg-python>=0.7.1

Downloading DAWG\_Python-0.7.2-py2.py3-none-any.whl (11 kB)

Collecting docopt>=0.6

Downloading docopt-0.6.2.tar.gz (25 kB)

Building wheels for collected packages: docopt

Building wheel for docopt (setup.py) ... done

Created wheel for docopt: filename=docopt-0.6.2-py2.py3-none-any.whl size=13705 sha256=d85aee54dead03a93b816c9de6743098922292053c9314a9f03cc10fa4608bae

Stored in directory: /root/.cache/pip/wheels/72/b0/3f/1d95f96ff986c7dfff e46ce2be4062f38ebd04b506c77c81b9

Successfully built docopt

Installing collected packages: pymorphy2-dicts-ru, docopt, dawg-python, razdel, pymorphy2, navec, yargy, slovnet, natasha

Successfully installed dawg-python-0.7.2 docopt-0.6.2 natasha-1.4.0 navec-0.10.0 pymorphy2-0.9.1 pymorphy2-dicts-ru-2.4.417127.4579844 razdel-0.5.0 slovnet-0.5.0 yargy-0.15.0



In [33]:

```
from natasha import (
    Segmenter,
    MorphVocab,

    NewsEmbedding,
    NewsMorphTagger,
    NewsSyntaxParser,
    NewsNERTagger,

    PER,
    NamesExtractor,

    Doc
)
segmenter = Segmenter()
emb = NewsEmbedding()
morph_tagger = NewsMorphTagger(emb)
syntax_parser = NewsSyntaxParser(emb)

doc = Doc(cleaned_text)
doc.segment(segmenter)
print(doc.tokens)
doc.tag_morph(morph_tagger)
doc.sents[1].morph.print()
```

[DocToken(start=7, text='Бросила'), DocToken(start=7, stop=8, text='.'), DocToken(start=9, stop=11, text='Во'), DocToken(start=12, stop=17, text='время'), DocToken(start=18, stop=25, text='каждого'), DocToken(start=26, stop=33, text='занятия'), DocToken(start=34, stop=42, text='получала'), DocToken(start=43, stop=52, text='небольшие'), DocToken(start=53, stop=59, text='травмы'), DocToken(start=60, stop=61, text='в'), DocToken(start=62, stop=71, text='последний'), DocToken(start=72, stop=75, text='раз'), DocToken(start=76, stop=79, text='при'), DocToken(start=80, stop=90, text='выполнении'), DocToken(start=91, stop=96, text='Сурьи'), DocToken(start=97, stop=105, text='Намаскар'), DocToken(start=106, stop=115, text='растянула'), DocToken(start=116, stop=122, text='связки'), DocToken(start=123, stop=125, text='на'), DocToken(start=126, stop=130, text='руке'), DocToken(start=131, stop=135, text='боль'), DocToken(start=136, stop=140, text='была'), DocToken(start=141, stop=146, text='такая'), DocToken(start=147, stop=150, text='что'), DocToken(start=151, stop=152, text='в'), DocToken(start=153, stop=159, text='первые'), DocToken(start=160, stop=167, text='секунды'), DocToken(start=168, stop=170, text='на'), DocToken(start=171, stop=178, text='перелом'), DocToken(start=179, stop=187, text='подумала'), DocToken(start=187, stop=188, text='.'), DocToken(start=189, stop=195, text='Теперь'), DocToken(start=196, stop=200, text='дома'), DocToken(start=201, stop=210, text='занимаюсь'), DocToken(start=211, stop=214, text='без'), DocToken(start=215, stop=222, text='напряга'), DocToken(start=223, stop=225, text='с'), DocToken(start=226, stop=233, text='стороны'), DocToken(start=234, stop=241, text='тренера'), DocToken(start=241, stop=242, text='.')]

Во ADP

время NOUN|Animacy=Inan|Case=Acc|Gender=Neut|Number=Sing

каждого DET|Case=Gen|Gender=Neut|Number=Sing

занятия NOUN|Animacy=Inan|Case=Gen|Gender=Neut|Number=Sing

получала VERB|Aspect=Imp|Gender=Fem|Mood=Ind|Number=Sing|Tense=Past|VerbForm=Fin|Voice=Act

небольшие ADJ|Animacy=Inan|Case=Acc|Degree=Pos|Number=Plur

травмы NOUN|Animacy=Inan|Case=Acc|Gender=Fem|Number=Plur

в ADP

последний ADJ|Animacy=Inan|Case=Acc|Degree=Pos|Gender=Masc|Number=Sing

раз NOUN|Animacy=Inan|Case=Acc|Gender=Masc|Number=Sing

при ADP

выполнении NOUN|Animacy=Inan|Case=Loc|Gender=Neut|Number=Sing

Сурьи PROPN|Animacy=Inan|Case=Gen|Gender=Masc|Number=Sing

Намаскар PROPN|Animacy=Anim|Case=Gen|Gender=Masc|Number=Sing

растянула ADJ|Case=Gen|Degree=Pos|Gender=Fem|Number=Sing

связки NOUN|Animacy=Inan|Case=Gen|Gender=Fem|Number=Sing

на ADP

руке NOUN|Animacy=Inan|Case=Loc|Gender=Fem|Number=Sing

боль NOUN|Animacy=Inan|Case=Acc|Gender=Fem|Number=Sing

была AUX|Aspect=Imp|Gender=Fem|Mood=Ind|Number=Sing|Tense=Past|VerbForm=Fin|Voice=Act

такая DET|Case=Nom|Gender=Fem|Number=Sing

что CONJ

в ADP

первые ADJ|Animacy=Inan|Case=Acc|Degree=Pos|Number=Plur

секунды NOUN|Animacy=Inan|Case=Acc|Gender=Fem|Number=Plur

на ADP

перелом NOUN|Animacy=Inan|Case=Acc|Gender=Masc|Number=Sing

подумала VERB|Aspect=Perf|Gender=Fem|Mood=Ind|Number=Sing|Tense=Past|VerbForm=Fin|Voice=Act

. PUNCT

### 3. Лемматизация (Lemmatization)

In [35]:

```
morph_vocab = MorphVocab()
for token in doc.tokens:
    token.lemmatize(morph_vocab)
{_.text: _.lemma for _ in doc.tokens}
```

Out[35]:

```
{'Бросила': 'бросить',
 '.': '.',
 'Во': 'в',
 'время': 'время',
 'каждого': 'каждый',
 'занятия': 'занятие',
 'получала': 'получать',
 'небольшие': 'небольшой',
 'травмы': 'травма',
 'в': 'в',
 'последний': 'последний',
 'раз': 'раз',
 'при': 'при',
 'выполнении': 'выполнение',
 'Сурьи': 'сурья',
 'Намаскар': 'намаскар',
 'растянула': 'растянуть',
 'связки': 'связка',
 'на': 'на',
 'руке': 'рука',
 'боль': 'боль',
 'была': 'быть',
 'такая': 'такой',
 'что': 'что',
 'первые': 'первый',
 'секунды': 'секунда',
 'перелом': 'перелом',
 'подумала': 'подумать',
 'Теперь': 'теперь',
 'дома': 'дома',
 'занимаюсь': 'заниматься',
 'без': 'без',
 'напряга': 'напряга',
 'со': 'с',
 'стороны': 'сторона',
 'тренера': 'тренер'}
```

### Выделение (распознавание) именованных сущностей (NER)

In [36]:

```
ner_tagger = NewsNERTagger(emb)
doc.tag_ner(ner_tagger)
doc.tag_ner(ner_tagger)
doc.ner.print()
```

Бросила. Во время каждого занятия получала небольшие травмы в  
последний раз при выполнении Сурьи Намаскар растянула связки на руке  
PER\_\_\_\_\_

боль была такая что в первые секунды на перелом подумала. Теперь дома  
занимаюсь без напряжения со стороны тренера.

In [74]:

```
text = '''Курская битва (5 июля – 23 августа 1943 года; также известна как Битва на Курской дуге) – совокупность стратегических оборонительной (5–23 июля) и наступательных (12 июля – 23 августа) операций Красной армии в Великой Отечественной войне с целью сорвать крупное наступление сил вермахта и разгромить его стратегическую группировку[9]. По своим масштабам, задействованным силам и средствам, напряжённости, результатам и военно-политическим последствиям является одним из ключевых сражений Второй мировой войны и Великой Отечественной войны. В историографии считается самым крупным (величайшим) танковым сражением в истории[10][11][12][13][14][Прим. 1]. В нём участвовали около 2 млн человек, 6 тысяч танков, 4 тысячи самолётов; сражение проложило «путь к великим советским наступательным действиям 1944–45 годов».
```

Сражение является важнейшей частью стратегического плана летне-осенней кампании 1943 года, согласно советской и российской историографии, включает в себя: Курскую стратегическую оборонительную операцию (5–23 июля), Орловскую (12 июля – 18 августа) и Белгородско-Харьковскую (3–23 августа) стратегические наступательные операции. Битва продолжалась 50 дней. Немецкая сторона наступательную часть сражения называла операция «Цитадель». В результате наступления по плану «Кутузов» потерпела поражение орловская группировка немецких войск, а занимаемый ею орловский стратегический плацдарм был ликвидирован. По итогам операции «Румянцев» потерпела поражение белгородско-харьковская группировка немцев, и соответствующий плацдарм также был ликвидирован[15]. Коренной перелом в ходе Великой Отечественной войны, начатый под Сталинградом, был завершён в Курской битве и сражении за Днепр, а в последовавшей Тегеранской конференции по инициативе Ф. Рузвельта уже обсуждался составленный им лично «2 месяца тому назад план расчленения Германии на пять государств».

После завершения битвы стратегическая инициатива окончательно перешла на сторону Красной армии, которая продолжала освобождать страну от немецких захватчиков и до окончания войны проводила в основном наступательные операции. Вермахт в ходе отступления с территории СССР проводил тактику «выжженной земли»[16].

23 августа – день разгрома советскими войсками немецких войск в Курской битве – является одним из дней воинской славы России. Белгород, Курск и Орёл стали первыми городами России, которым присвоено почётное звание «Город воинской славы»[17].'''

```
doc2 = Doc(text)
doc2.segment(segmenter)
doc2.tag_morph(morph_tagger)
#doc2.sents[1].morph.print()
for token in doc.tokens:
    token.lemmatize(morph_vocab)
doc2.tag_ner(ner_tagger)
doc2.ner.print()
```

Курская битва (5 июля – 23 августа 1943 года; также известна как Битва на Курской дуге) – совокупность стратегических оборонительной (5–23 июля) и наступательных (12 июля – 23 августа) операций Красной армии в  
ORG\_\_\_\_\_

Великой Отечественной войне с целью сорвать крупное наступление сил вермахта и разгромить его стратегическую группировку[9]. По своим масштабам, задействованным силам и средствам, напряжённости, результатам и военно-политическим последствиям является одним из ключевых сражений Второй мировой войны и Великой Отечественной войны. В историографии считается самым крупным (величайшим) танковым сражением в истории[10][11][12][13][14][Прим. 1]. В нём участвовали около 2 млн человек, 6 тысяч танков, 4 тысячи самолётов; сражение проложило «путь к великим советским наступательным действиям 1944–45 годов».

Сражение является важнейшей частью стратегического плана летне-осенней кампании 1943 года, согласно советской и российской историографии, включает в себя: Курскую стратегическую оборонительную операцию (5–23 июля), Орловскую (12 июля – 18 августа) и Белгородско-Харьковскую  
LOC\_\_\_\_\_ LOC\_\_\_\_\_

(3–23 августа) стратегические наступательные операции. Битва продолжалась 50 дней. Немецкая сторона наступательную часть сражения называла операция «Цитадель».

В результате наступления по плану «Кутузов» потерпела поражение  
ORG\_\_\_\_\_

орловская группировка немецких войск, а занимаемый ею орловский стратегический плацдарм был ликвидирован. По итогам операции «Румянцев» потерпела поражение белгородско-харьковская группировка немцев, и соответствующий плацдарм также был ликвидирован[15].

Коренной перелом в ходе Великой Отечественной войны, начатый под Сталинградом, был завершён в Курской битве и сражении за Днепр, а в  
LOC\_\_\_\_\_ LOC—

последовавшей Тегеранской конференции по инициативе Ф. Рузвельта уже  
PER\_\_\_\_\_

обсуждался составленный им лично «2 месяца тому назад план расчленения Германии на пять государств».  
LOC\_\_\_\_\_

После завершения битвы стратегическая инициатива окончательно перешла на сторону Красной армии, которая продолжала освобождать страну от  
ORG\_\_\_\_\_

немецких захватчиков и до окончания войны проводила в основном наступательные операции. Вермахт в ходе отступления с территории СССР  
LOC—

проводил тактику «выжженной земли»[16].

23 августа – день разгрома советскими войсками немецких войск в Курской битве – является одним из дней воинской славы России.  
LOC—

Белгород, Курск и Орёл стали первыми городами России, которым  
LOC\_\_\_\_\_ LOC— LOC— LOC—  
присвоено почётное звание «Город воинской славы»[17].

In [84]:

```
names_extractor = NamesExtractor(morph_vocab)
for span in doc2.spans:
    span.normalize(morph_vocab)
{_.text: _.normal for _ in doc2.spans if _.text != _.normal}
```

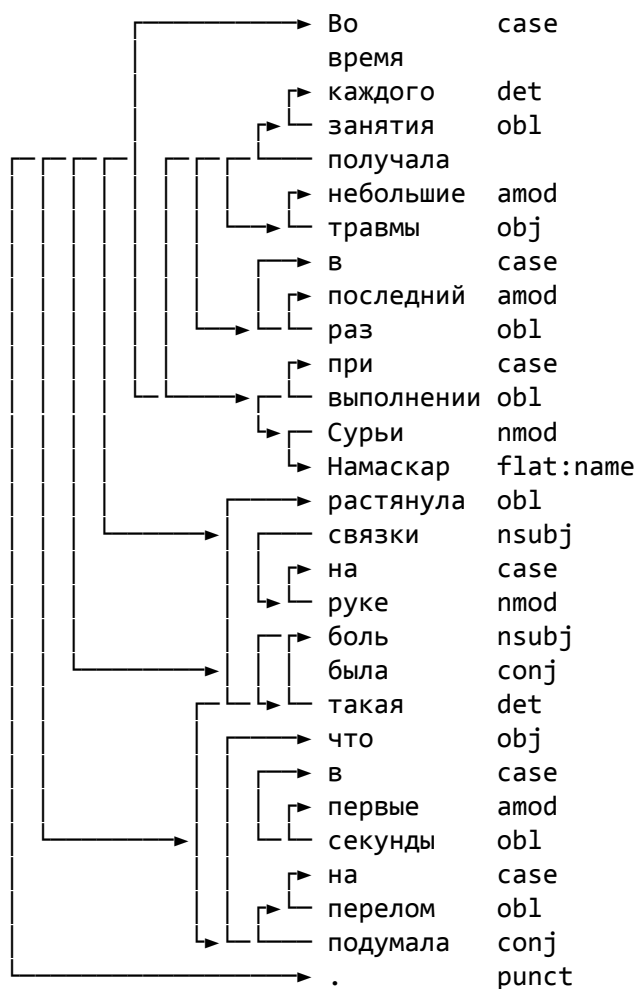
Out[84]:

```
{'Красной армии': 'Красная армия',
 'Орловскую': 'Орловская',
 'Сталинградом': 'Сталинград',
 'Ф. Рузвельта': 'Ф. Рузвельт',
 'Германии': 'Германия',
 'России': 'Россия',
 'Орёл': 'Орел'}
```

## 5. Разбор предложения

In [86]:

```
doc.parse_syntax(syntax_parser)
doc.sents[1].syntax.print()
```



# Лабораторная работа №6

Ишков Денис Олегович, ИУ5-24М, 2021г.

## Тема: Классификация текста

### Цель: Изучение методов классификации текстов.

Требования к отчету: Отчет по лабораторной работе должен содержать:

- титульный лист;
- описание задания;
- текст программы;
- экранные формы с примерами выполнения программы. Задание:

Для произвольного набора данных, предназначенного для классификации текстов, решите задачу классификации текста двумя способами:

- Способ 1. На основе CountVectorizer или TfidfVectorizer.
- Способ 2. На основе моделей word2vec или Glove или fastText.
- Сравните качество полученных моделей.

Для поиска наборов данных в поисковой системе можно использовать ключевые слова "datasets for text classification".

### Разделим данные на обучающую и тестовую выборки

In [88]:

```
train_df, test_df = train_test_split(df, test_size=500, random_state=69, stratify=df.toxic)
```

In [90]:

```
test_df["toxic"].value_counts(), train_df["toxic"].value_counts()
```

Out[90]:

```
(0    324
 1    176
Name: toxic, dtype: int64,
0    7778
 1    4243
Name: toxic, dtype: int64)
```

### Предобработка текста



In [91]:

```

sentence_example = df.iloc[-1]["comment"]
tokens = word_tokenize(sentence_example, language="russian")
tokens_without_punctuation = [i for i in tokens if i not in string.punctuation]
russian_stop_words = stopwords.words("russian")
tokens_without_stop_words_and_punctuation = [i for i in tokens_without_punctuation if i
not in russian_stop_words]
snowball = SnowballStemmer(language="russian")
stemmed_tokens = [snowball.stem(i) for i in tokens_without_stop_words_and_punctuation]

```

In [92]:

```

print(f"Исходный текст: {sentence_example}")
print("-----")
print(f"Токены: {tokens}")
print("-----")
print(f"Токены без пунктуации: {tokens_without_punctuation}")
print("-----")
print(f"Токены без пунктуации и стоп слов: {tokens_without_stop_words_and_punctuation}")
print("-----")
print(f"Токены после стемминга: {stemmed_tokens}")
print("-----")

```

Исходный текст: До сих пор пересматриваю его видео. Орамбо кстати на своем канале пилит похожий контент, но качеством похуже, там же и Шуран не редко светится, храню хрупкую надежду что когда-то он вернется, такая годнота ведь.

-----

Токены: ['До', 'сих', 'пор', 'пересматриваю', 'его', 'видео', '.', 'Орамб', 'о', 'кстати', 'на', 'своем', 'канале', 'пилит', 'похожий', 'контент', ',', 'но', 'качеством', 'похуже', ',', 'там', 'же', 'и', 'Шуран', 'не', 'редк', 'о', 'светится', ',', 'храню', 'хрупкую', 'надежду', 'что', 'когда-то', 'о', 'н', 'вернется', ',', 'такая', 'годнота', 'ведь', '.']

-----

Токены без пунктуации: ['До', 'сих', 'пор', 'пересматриваю', 'его', 'видео', 'Орамбо', 'кстати', 'на', 'своем', 'канале', 'пилит', 'похожий', 'контент', 'но', 'качеством', 'похуже', 'там', 'же', 'и', 'Шуран', 'не', 'редко', 'светится', 'храню', 'хрупкую', 'надежду', 'что', 'когда-то', 'он', 'вернется', 'такая', 'годнота', 'ведь']

-----

Токены без пунктуации и стоп слов: ['До', 'сих', 'пор', 'пересматриваю', 'видео', 'Орамбо', 'кстати', 'своем', 'канале', 'пилит', 'похожий', 'контент', 'качеством', 'похуже', 'Шуран', 'редко', 'светится', 'храню', 'хрупкую', 'надежду', 'когда-то', 'вернется', 'такая', 'годнота']

-----

Токены после стемминга: ['до', 'сих', 'пор', 'пересматрива', 'виде', 'орамб', 'кстат', 'сво', 'канал', 'пил', 'похож', 'контент', 'качеств', 'похуж', 'шура', 'редк', 'свет', 'хран', 'хрупк', 'надежд', 'когда-т', 'вернет', 'так', 'годнот']

-----

In [96]:

```
segmenter = Segmenter()
morph_vocab = MorphVocab()

emb = NewsEmbedding()
morph_tagger = NewsMorphTagger(emb)
syntax_parser = NewsSyntaxParser(emb)
ner_tagger = NewsNERTagger(emb)
```

In [99]:

```
snowball = SnowballStemmer(language="russian")
russian_stop_words = stopwords.words("russian")

def tokenize_sentence(sentence: str, remove_stop_words: bool = True):
    tokens = word_tokenize(sentence, language="russian")
    tokens = [i for i in tokens if i not in string.punctuation]
    if remove_stop_words:
        tokens = [i for i in tokens if i not in russian_stop_words]
    doc = Doc(' '.join(tokens))
    doc.segment(segmenter)
    doc.tag_morph(morph_tagger)
    doc.parse_syntax(syntax_parser)
    doc.tag_ner(ner_tagger)
    for token in doc.tokens:
        token.lemmatize(morph_vocab)
    tokens = [t.lemma for t in doc.tokens] #[snowball.stem(i) for i in tokens]
    return tokens

tokenize_sentence(sentence_example)
```

Out[99]:

```
['до',
 'сей',
 'пора',
 'пересматривать',
 'видео',
 'орамбо',
 'кстати',
 'свой',
 'канал',
 'пилить',
 'похожий',
 'контент',
 'качество',
 'плохой',
 'шуран',
 'редко',
 'светиться',
 'хранить',
 'хрупкий',
 'надежда',
 'когда-то',
 'вернуться',
 'такой',
 'годнота']
```

# Классификатор 1: Complement Naive Bayes

## CountVectorizer

In [101]:

```

grid_pipeline = Pipeline([
    ("vectorizer", CountVectorizer(tokenizer=lambda x: tokenize_sentence(x, remove_stop
_words=True))),
    ("model",
    GridSearchCV(
        ComplementNB(alpha=1.0, norm=False),
        param_grid={'alpha': [0.35+5e-3*i for i in range(20)]+\
                             [0.5+5e-3*i for i in range(60)]+\
                             [0.5, 0.75, 1., 10., 100., 1e3, 1e4, 1e5],
                    'norm': [True, False]},
        cv=5,
        verbose=0,
        scoring='roc_auc',
    )
])

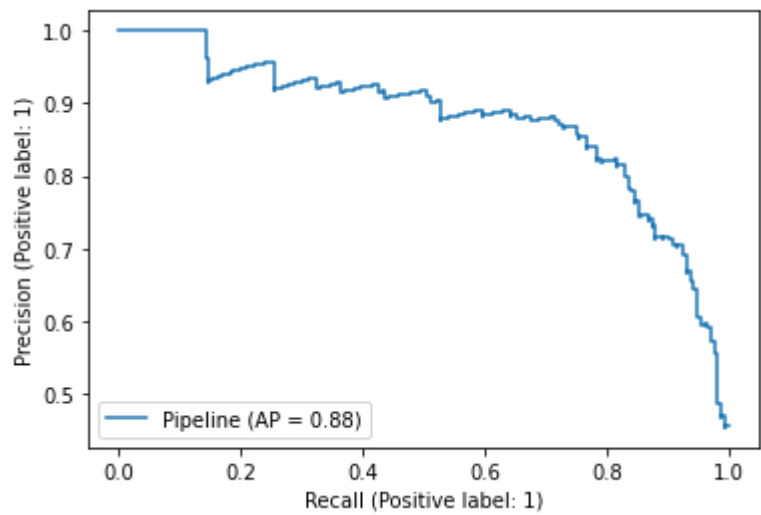
grid_pipeline.fit(train_df["comment"], train_df["toxic"])
print(grid_pipeline['model'].best_params_)
model_pipeline = Pipeline([
    ("vectorizer", CountVectorizer(tokenizer=lambda x: tokenize_sentence(x, remove_stop
_words=True))),
    ("model", ComplementNB(**grid_pipeline['model'].best_params_))
])

model_pipeline.fit(train_df["comment"], train_df["toxic"])
prec_c_10, rec_c_10, thresholds_c_10 = precision_recall_curve(y_true=test_df["toxic"],
                                                             probas_pred=model_pipeline
e.predict_proba(test_df["comment"])[ :, 1])
plot_precision_recall_curve(estimator=model_pipeline, X=test_df["comment"], y=test_df[
"toxic"])
print(classification_report(y_true=test_df["toxic"],
                           y_pred=model_pipeline.predict(test_df["comment"]),
                           digits=4))

```

{'alpha': 0.575, 'norm': False}

	precision	recall	f1-score	support
0	0.8794	0.9228	0.9006	324
1	0.8438	0.7670	0.8036	176
accuracy			0.8680	500
macro avg	0.8616	0.8449	0.8521	500
weighted avg	0.8669	0.8680	0.8664	500



TfidfVectorizer

In [100]:

```

grid_pipeline = Pipeline([
    ("vectorizer", TfidfVectorizer(tokenizer=lambda x: tokenize_sentence(x, remove_stop
_words=True))),
    ("model",
    GridSearchCV(
        ComplementNB(alpha=1.0, norm=False),
        param_grid={'alpha': [0.3+5e-3*i for i in range(60)]+\
            [0.5, 0.75, 1., 10., 100., 1e3, 1e4, 1e5],
                    'norm': [True, False]},
        cv=5,
        verbose=0,
        scoring='roc_auc',
    )
])

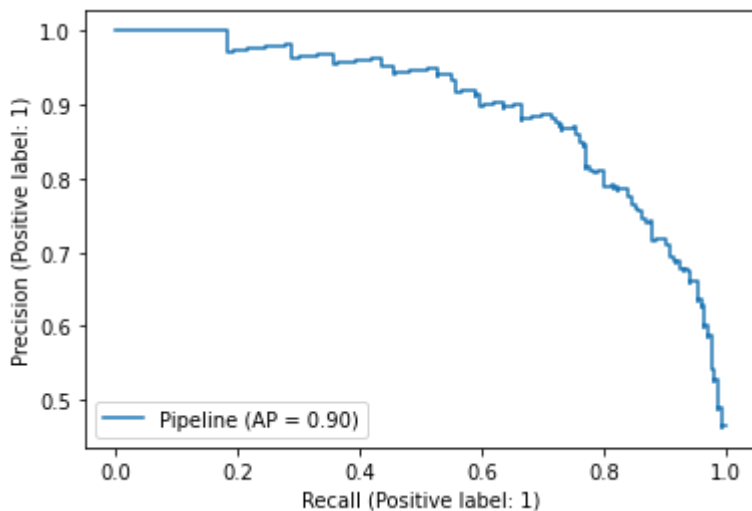
grid_pipeline.fit(train_df["comment"], train_df["toxic"])
print(grid_pipeline['model'].best_params_)
model_pipeline = Pipeline([
    ("vectorizer", TfidfVectorizer(tokenizer=lambda x: tokenize_sentence(x, remove_stop
_words=True))),
    ("model", ComplementNB(**grid_pipeline['model'].best_params_))
])

model_pipeline.fit(train_df["comment"], train_df["toxic"])
prec_c_10, rec_c_10, thresholds_c_10 = precision_recall_curve(y_true=test_df["toxic"],
                                                                probas_pred=model_pipeline
                                                                .predict_proba(test_df["comment"])[ :, 1])
plot_precision_recall_curve(estimator=model_pipeline, X=test_df["comment"], y=test_df[
"toxic"])
print(classification_report(y_true=test_df["toxic"],
                            y_pred=model_pipeline.predict(test_df["comment"]),
                            digits=4))

```

```
{'alpha': 0.35, 'norm': False}
```

	precision	recall	f1-score	support
0	0.8665	0.9414	0.9024	324
1	0.8716	0.7330	0.7963	176
accuracy			0.8680	500
macro avg	0.8690	0.8372	0.8493	500
weighted avg	0.8683	0.8680	0.8650	500



## Классификация с помощью эмбедингов Navec

In [102]:

```
# загрузка
from navec import Navec
!wget https://storage.yandexcloud.net/natasha-navec/packs/navec_hudlit_v1_12B_500K_300d_100q.tar
```

```
--2021-05-23 12:27:17-- https://storage.yandexcloud.net/natasha-navec/packs/navec_hudlit_v1_12B_500K_300d_100q.tar
Resolving storage.yandexcloud.net (storage.yandexcloud.net)... 213.180.193.243, 2a02:6b8::1d9
Connecting to storage.yandexcloud.net (storage.yandexcloud.net)|213.180.193.243|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 53012480 (51M) [application/x-tar]
Saving to: 'navec_hudlit_v1_12B_500K_300d_100q.tar'
```

```
navec_hudlit_v1_12B 100%[=====>] 50.56M 15.9MB/s in 4.2s
```

```
2021-05-23 12:27:22 (12.0 MB/s) - 'navec_hudlit_v1_12B_500K_300d_100q.tar'
saved [53012480/53012480]
```

In [107]:

```
path = 'navec_hudlit_v1_12B_500K_300d_100q.tar'
navec = Navec.load(path)
navec['навек'].shape
```

Out[107]:

(300,)

In [138]:

```
def embed_sentence(sentence):
    s = np.zeros((50, 300))
    try:
        for i, token in enumerate(tokenize_sentence(sentence)):
            emb = navec.get(token)
            if emb is not None:
                s[i] = emb
    except:
        print(sentence)
    pbar.update(1)
    return s

pbar = tqdm()
train_embeddings = np.array(train_df.comment.apply(embed_sentence).to_list())
test_embeddings = np.array(test_df.comment.apply(embed_sentence).to_list())
```

Когда пикрил детей резал, то пидарашки тоже ничего не слышали и не видели. Хотя несколько детей у него орали пару суток, пока он, перед убийством, на д ними рофлил. Просто у рузких особенность такая. Кланяться всем, кто представляет опасность и смаковать мучения более слабых попутно высирая на них доносы.

СРОЧНА!!!!!!!!!!!!!!!!!!!!!!!!!!!! ИЗБИЕНИЕ И!З!Б!И!Е!Н!И!Е БЕСПЛАТНА И БЕЗ СМС ПОДПИСЫВАЙТЕСЬ НА КОНАЛ СТАВЬТЕ ЛАЙК ДОНАТЬТЕ ПРОДАЙТЕ ПОЧКУ СДЕЛАЙТЕ БОЧКУ ПОСМОТРИТЕ ВИДЕО Я КАПСОМ НАПИСАЛ ЧТОБЫ ВЫ ПОСМОТРЕЛИ СМОТРИТЕ ЕЩЕ Е БУЧИХ ВОСКЛИЦАТЕЛЬНЫХ ЗНАКОВ НАСТАВИЛ БЕСПЛАТНО БЕЗ СМС СЛЫШИТЕ БЕСПЛАТНА И БЕЗ С М С СМОТРЕТЬ ОНЛАЙН ОФЛАЙН ХУЙЛАЙН

Например вот тут: [https://www.google.ru/url?sa=t&source=web&rct=j&url=http://www.consultant.ru/document/cons\\_doc\\_LAW\\_34683\\_aed7d03df679e3376974dadd131b899dc6966650\\_ved\\_2ahUKEwiar-2jj8\\_gAhVE-yoKHR3ECB4QFjAAegQIBBAB\\_usg\\_A0vVaw090bjWoB9D2SDDZBVcDLUF](https://www.google.ru/url?sa=t&source=web&rct=j&url=http://www.consultant.ru/document/cons_doc_LAW_34683_aed7d03df679e3376974dadd131b899dc6966650_ved_2ahUKEwiar-2jj8_gAhVE-yoKHR3ECB4QFjAAegQIBBAB_usg_A0vVaw090bjWoB9D2SDDZBVcDLUF)

## Стратегии (Reshape, Max, Mean, Min)

### Reshape



In [148]:

```
grid_pipeline = Pipeline([
    ("model",
     GridSearchCV(
         LogisticRegression(random_state=42, max_iter=200),
         param_grid={'C': [1e-5, 1e-3, 0.01, 0.1]},
         cv=4,
         verbose=3,
         scoring='roc_auc',
     )
])

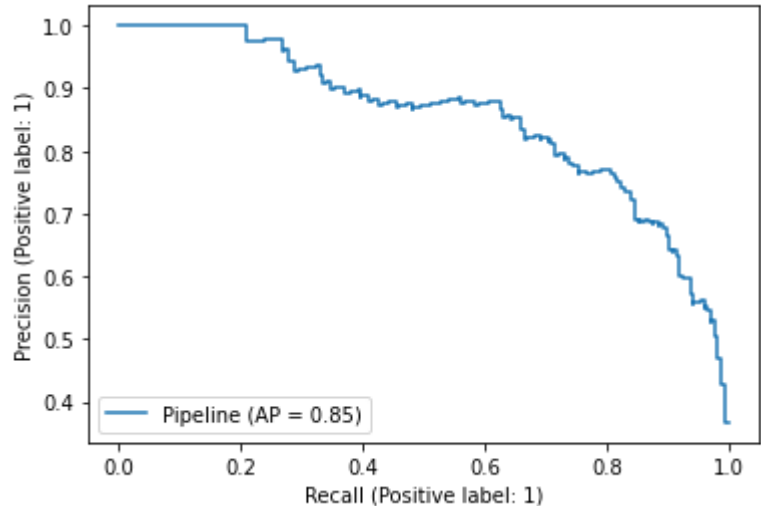
# X_tr = train_embeddings.reshape(-1, 50*300)
# X_t = test_embeddings.reshape(-1, 50*300)

grid_pipeline.fit(X_tr, train_df["toxic"])
print(grid_pipeline['model'].best_params_)
model_pipeline = Pipeline([
    ("model", LogisticRegression(**grid_pipeline['model'].best_params_),)
])

model_pipeline.fit(X_tr, train_df["toxic"])
plot_precision_recall_curve(estimator=model_pipeline, X=X_t, y=test_df["toxic"])
print(classification_report(y_true=test_df["toxic"],
                             y_pred=model_pipeline.predict(X_t),
                             digits=4))
```

```
Fitting 4 folds for each of 4 candidates, totalling 16 fits
[CV 1/4] END .....C=1e-05; total time=
3.4s
[CV 2/4] END .....C=1e-05; total time=
3.9s
[CV 3/4] END .....C=1e-05; total time=
3.4s
[CV 4/4] END .....C=1e-05; total time=
3.7s
[CV 1/4] END .....C=0.001; total time=
5.8s
[CV 2/4] END .....C=0.001; total time=
5.5s
[CV 3/4] END .....C=0.001; total time=
5.8s
[CV 4/4] END .....C=0.001; total time=
5.8s
[CV 1/4] END .....C=0.01; total time=
8.9s
[CV 2/4] END .....C=0.01; total time=
9.7s
[CV 3/4] END .....C=0.01; total time=
9.1s
[CV 4/4] END .....C=0.01; total time=
9.3s
[CV 1/4] END .....C=0.1; total time=
18.5s
[CV 2/4] END .....C=0.1; total time=
17.7s
[CV 3/4] END .....C=0.1; total time=
19.2s
[CV 4/4] END .....C=0.1; total time=
19.6s
{'C': 0.01}
```

	precision	recall	f1-score	support
0	0.8462	0.9167	0.8800	324
1	0.8188	0.6932	0.7508	176
accuracy			0.8380	500
macro avg	0.8325	0.8049	0.8154	500
weighted avg	0.8365	0.8380	0.8345	500



In [156]:

```
# MEAN
grid_pipeline = Pipeline([
    ("model",
     GridSearchCV(
         LogisticRegression(random_state=42, max_iter=300),
         param_grid={'C': [4.5, 5, 5.5, 6.]},
         cv=4,
         verbose=3,
         scoring='roc_auc',
     )
])

X_tr = train_embeddings.mean(axis=1)
X_t = test_embeddings.mean(axis=1)

grid_pipeline.fit(X_tr, train_df["toxic"])
print(grid_pipeline['model'].best_params_)
model_pipeline = Pipeline([
    ("model", LogisticRegression(**grid_pipeline['model'].best_params_),)
])

model_pipeline.fit(X_tr, train_df["toxic"])
plot_precision_recall_curve(estimator=model_pipeline, X=X_t, y=test_df["toxic"])
print(classification_report(y_true=test_df["toxic"],
                           y_pred=model_pipeline.predict(X_t),
                           digits=4))
```

Fitting 4 folds for each of 4 candidates, totalling 16 fits

```
[CV 1/4] END .....C=4.5; total time=
0.5s
[CV 2/4] END .....C=4.5; total time=
0.6s
[CV 3/4] END .....C=4.5; total time=
0.4s
[CV 4/4] END .....C=4.5; total time=
0.4s
[CV 1/4] END .....C=5; total time=
0.4s
[CV 2/4] END .....C=5; total time=
0.6s
[CV 3/4] END .....C=5; total time=
0.6s
[CV 4/4] END .....C=5; total time=
0.5s
[CV 1/4] END .....C=5.5; total time=
0.5s
[CV 2/4] END .....C=5.5; total time=
0.4s
[CV 3/4] END .....C=5.5; total time=
0.5s
[CV 4/4] END .....C=5.5; total time=
0.5s
[CV 1/4] END .....C=6.0; total time=
0.3s
[CV 2/4] END .....C=6.0; total time=
0.6s
[CV 3/4] END .....C=6.0; total time=
0.6s
[CV 4/4] END .....C=6.0; total time=
0.5s
{'C': 5}
```

```
/opt/conda/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:7
65: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max\_iter) or scale the data as shown i  
n:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG)

	precision	recall	f1-score	support
0	0.8626	0.9105	0.8859	324
1	0.8165	0.7330	0.7725	176
accuracy			0.8480	500
macro avg	0.8395	0.8217	0.8292	500
weighted avg	0.8463	0.8480	0.8460	500

