

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

Лабораторная работа №6

по дисциплине «Методы машинного обучения»

на тему «Классификация текста»

Выполнил: Ван Пэй

Группа: ИУ5-22М

Москва — 2021 г.

Цель лабораторной работы:

Изучение методов классификации текстов.

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".

Описание решения:

Использован способ 2: На основе моделей. Начало модели соответствует Лаб 5.

3. Базовые векторные представления слов и текстов

Мешок слов aka Bag of words aka BOW)

```
] : with open('../input/preprocessed1/lemmatized_tokens (1).txt', 'r') as
    data = f.read().split('\n')

tokenized_lemmatized_words = [i.split(',') for i in data]
tokenized_lemmatized_words[:5]
```

```
[['джоан',
  'кэтли́н',
  'роули́нг',
  'га́рри',
  'потте́р',
  'фило́софский',
  'каме́нь',
  'гла́ва',
  'перва́я',
  'ма́льчик',
  'кото́рый',
  'выжи́л',
  'мисте́р',
  'мисси́с',
  'да́рсли',
  'дома́',
  'номе́р',
  'четы́ре',
  'бирю́чинному',
  'пое́зду',
  'мо́гли',
  'гордо́стью',
  'сказа́ть',
  'сла́ва',
  'бо́гу',
  'соверше́нно',
  'норма́льные',
  'люди́'],
 ['последними́',
  'кого́',
  'мо́гли',
  'ожи́датель',
  ]]
```

```
import gensim
from gensim import corpora

dictionary = corpora.Dictionary()

corpus = [dictionary.doc2bow(sent, allow_update=True) for sent in tokens]
```

```
# тут можно посмотреть что лежит в словаре
# dictionary.token2id
```

```
for doc in corpus[:10]:
    print([[dictionary[id], freq] for id, freq in doc])
```

```
[[['бирючинному', 1], ['богу', 1], ['выжил', 1],
['гарри', 1], ['глава', 1], ['гордостью', 1],
['дарсли', 1], ['джоан', 1], ['дома', 1], ['камень', 1],
['который', 1], ['кэтлин', 1], ['люди', 1],
['мальчик', 1], ['миссис', 1], ['мистер', 1],
['могли', 1], ['номер', 1], ['нормальные', 1],
['первая', 1], ['поттер', 1], ['проезду', 1],
['роулинг', 1], ['сказать', 1], ['слава', 1],
['совершенно', 1], ['философский', 1], ['четыре', 1]]
[['могли', 1], ['совершенно', 1], ['кого', 1],
['одобряли', 1], ['ожидать', 1], ['подобной', 1],
['последними', 1], ['странном', 1], ['таинственным', 1],
['участия', 1], ['чемнибудь', 1], ['чепухи', 1]]
[['дарсли', 1], ['мистер', 1], ['выпускала', 1],
['граннингз', 1], ['директором', 1], ['которая', 1],
['работал', 1], ['сверла', 1], ['фирмы', 1]]
[['большими', 1], ['короткой', 1], ['крепкий', 1],
['крупный', 1], ['мужчина', 1], ['очень', 2],
['усами', 1], ['шеей', 1]]
[['дарсли', 1], ['миссис', 1], ['двоих', 1], ['зато', 1],
['лихвой', 1], ['светловолосой', 1], ['хватило', 1],
['худой', 1], ['шеи', 1]]
[['большую', 1], ['времени', 1], ['изгороди', 1],
['подарком', 1], ['пользовалась', 1], ['природы', 1],
['садовые', 1], ['соседями', 1], ['усердно', 1],
['часть', 1], ['шпионя', 1], ['этим', 1]]
[['дарсли', 1], ['всем', 1], ['дадли', 1], ['лучш
```

```
len(dictionary.keys())
```

12848

```
from gensim.models import TfidfModel
```

```
tfidf = TfidfModel(corpus)
```

```
for doc in tfidf[corpus[:10]]:
```

```
    print([[dictionary[id], round(freq, 2)] for id, freq in doc])
```

```
[['бирючинному', 0.25], ['богу', 0.23], ['выжил', 0.21], ['гарри', 0.05], ['глава', 0.17], ['гордостью', 0.23], ['дарсли', 0.12], ['джоан', 0.25], ['дома', 0.16], ['камень', 0.14], ['который', 0.12], ['кэтлин', 0.25], ['люди', 0.15], ['мальчик', 0.15], ['миссис', 0.14], ['мистер', 0.13], ['могли', 0.15], ['номер', 0.18], ['нормальные', 0.25], ['первая', 0.22], ['поттер', 0.13], ['проезду', 0.25], ['роулинг', 0.25], ['сказать', 0.14], ['слава', 0.22], ['совершенно', 0.18], ['философский', 0.2], ['четыре', 0.16]]
[['могли', 0.19], ['совершенно', 0.23], ['кого', 0.23], ['одобряли', 0.32], ['ожидать', 0.26], ['подобной', 0.32], ['последними', 0.32], ['странном', 0.32], ['таинственным', 0.32], ['участия', 0.32], ['чемнибудь', 0.29], ['чепухи', 0.32]]
[['дарсли', 0.19], ['мистер', 0.21], ['выпускала', 0.39], ['граннингз', 0.39], ['директором', 0.39], ['которая', 0.24], ['работал', 0.39], ['сверла', 0.34], ['фирмы', 0.39]]
[['большими', 0.32], ['короткой', 0.36], ['крепкий', 0.39], ['крупный', 0.36], ['мужчина', 0.32], ['очень', 0.34], ['усами', 0.36], ['шеей', 0.39]]
[['дарсли', 0.19], ['миссис', 0.22], ['двоих', 0.39], ['зато', 0.33], ['лихвой', 0.39], ['светловол
```

Word2Vec(word to vec) ¶

```
with open('../input/preprocessed/with_pos_tag (3).txt', 'r') as f:
    data = f.read().split('\n')
```

```
with_pos_tags = [sent.split(',') for sent in data]
```

```
with_pos_tags[:5]
```

```
[('джоан_NOUN',
  'кэтлинка_NOUN',
  'роулинга_NOUN',
  'гарри_NOUN',
  'поттер_NOUN',
  'философский_NOUN',
  'камень_NOUN',
  'глава_NOUN',
  'первый_ADJF',
  'мальчик_NOUN',
  'который_ADJF',
  'выжить_INFN',
  'мистер_NOUN',
  'миссис_NOUN',
  'дарслить_INFN',
  'из_PREP',
  'дом_NOUN',
  'номер_NOUN',
  'четыре_NUMR',
  'по_PREP',
  'брючинный_ADJF',
  'проезд_NOUN',
  'мочь_NOUN',
  'бы_PRCL',
  'гордость_NOUN',
  'сказать_INFN',
  'что_CONJ',
  '..._NOUN']
```

```
import gensim.downloader as api
model = api.load("word2vec-ruscorpora-300") # download the model
```

```
model.most_similar(positive=['женщина_NOUN', 'король_NOUN'])
```

```
[('королева_NOUN', 0.7313905358314514),
 ('герцог_NOUN', 0.6502389311790466),
 ('принцесса_NOUN', 0.6266286373138428),
 ('герцогиня_NOUN', 0.6240381598472595),
 ('королевство_NOUN', 0.6094207763671875),
 ('зюдерманландский_ADJ', 0.6084389686584473),
 ('дурлахский_ADJ', 0.608166515827179),
 ('ульрик::элеонора_NOUN', 0.6073106527328491),
 ('максимилианов_NOUN', 0.6057005524635315),
 ('принц_NOUN', 0.5984029173851013)]
```

Список литературы

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

[Электронный ресурс]

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