



**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 №6

## Классификация текста.

Набор данных - [20 newsgroups text dataset](#)

Классы: 20

Выборка: 18846

## Install

```
1 pip3 install gensim
2 pip3 install spacy
3 python3 -m spacy download en_core_web_sm
```

## 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".

```
1 import nltk
2 import spacy
3 import numpy as np
4 import sklearn
5 from sklearn.datasets import fetch_20newsgroups
6 from nltk import tokenize
7 import re
8 import pandas as pd
```

```
1 nltk.download('punkt')
```

```
1 [nltk_data] Downloading package punkt to /Users/snipghost/nltk_data...
2 [nltk_data] Unzipping tokenizers/punkt.zip.
```

```
1 True
```

```
1 categories = ["sci.crypt", "sci.electronics", "talk.religion.misc"]
2 newsgroups_train = fetch_20newsgroups(subset='train',
   categories=categories)
3 newsgroups_test = fetch_20newsgroups(subset='test', categories=categories)
```

```
1 unique, frequency = np.unique(newsgroups_train.target, return_counts=True)
```

```
1 for l, f in zip(unique, frequency):
2     print(f'value: {l}, count: {f}')
```

```
1 value: 0, count: 595
2 value: 1, count: 591
3 value: 2, count: 377
```

```
1 from spacy.lang.en import English
2 import spacy
3 from nltk.corpus import stopwords
4
5 nlp = spacy.load("en_core_web_sm", disable=["parser", "ner"])
6 nltk.download('stopwords')
7 stopwords_eng = set(stopwords.words('english'))
```

```
1 [nltk_data] Downloading package stopwords to
2 [nltk_data]      /Users/snipghost/nltk_data...
3 [nltk_data] Unzipping corpora/stopwords.zip.
```

```

1 def prepare(t):
2     # t = ' '.join([i.strip().lower() for i in t.split(' ')])
3     t = re.sub(r'^a-zA-Z0-9 \n]', ' ', t)
4     t = re.sub('\s+', ' ', t)
5     t = ' '.join([token.lemma_.lower() for token in nlp(t) if token not in
stopwords_eng])
6     return t
7
8 texts = newsgroups_train.data
9 texts_array = []
10
11 for text in texts:
12     prepared_text = prepare(text)
13     texts_array.append(prepared_text)

```

```

1 print(len(texts_array), texts_array[0])

```

```

1 1563 from mhalldlynxdacnortheasternedu mark hald subject re dayton hamfest
organization northeastern university boston ma 02115 usa distribution usa
lines 13 i book a hotel red roof inn last week in cincinnati blue ash which
be at the northern tip of the metro cincy area i choose it for a few reason
1 all hotel in and near dayton be book solid 2 this hotel be only cost
28night 3 it be one of about 4 room leave on the night i reserve 4
cincinnati probably have more to to at night than dayton i intend to hit
the riverboat entertainment at dusk if anyone have other suggestion for
nightlife please let i know of other hot spot thanks mark

```

```

1 test_texts = newsgroups_test.data
2 test_texts_arr = []
3
4 for text in test_texts:
5     prepared_text = prepare(text)
6     test_texts_arr.append(prepared_text)

```

```

1 print(len(test_texts_arr), test_texts_arr[0])

```

```
1 1040 from markkcyresswestsuncom mark kampe subject re cybele and
transgendersexualism organization sunsoft south lines 29 distribution world
replyto markkcyresswestsuncom nntppostinghost sagredo in article
260493115730ravenaimsuncedu fhuntmeduncedu freb hunt write be there some
relation between the name cybele and the phenemenon of the sibyl your
paragraph above seem to indicate there might be the oed give the etymology
of sibyl as come from the ancient greek sigma iota beta upsilon lambda
lambda alpha s i b ih l l a which be claim to come from the doric sigma
iota omicron beta upsilon lambda lambda alpha s i o b ih l l a which if i
read it properly in turn come from the attican athenian theta epsilon
omicron beta omicron upsilon lambda eta th eh o b o ih l ae i do nt know
much about attis but it would nt surprise i to learn that this god be tie
to the athenian capital alpha tau tau iota kappa upsilon sigma a t t i k u
s the oed do not list any etymology for cybele since that be a proper noun
but i suggest that the greek spelling of that word would be much close to
the anticedant of sibyl than the two word be now perhaps cybele be a french
or latin spelling
```

## Способ 1.

TF-IDF + CountVectorizer

```
1 from sklearn.feature_extraction.text import TfidfVectorizer,
CountVectorizer
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.metrics import classification_report
```

```
1 tfidf_vectorizer = TfidfVectorizer()
2
3 train_feature_matrix_tfidf = tfidf_vectorizer.fit_transform(texts_array)
4 test_feature_matrix__tfidf = tfidf_vectorizer.transform(test_texts_arr)
```

```
1 count_vectorizer = CountVectorizer()
2
3 train_feature_matrix_count = count_vectorizer.fit_transform(texts_array)
4 test_feature_matrix_count = count_vectorizer.transform(test_texts_arr)
```

```
1 target_values_train = newsgroups_train.target
2 target_values_test = newsgroups_test.target
```

```
1 # sklearn.metrics.SCORERS.keys()
```

```

1 knn_tfidf = KNeighborsClassifier()
2 parameters = {'n_neighbors': [2, 3, 5, 7, 9, 11]}
3
4 knn_tfidf_grid = GridSearchCV(knn_tfidf, parameters,
    scoring='balanced_accuracy', verbose=4, cv=5)

```

```

1 knn_tfidf_grid.fit(train_feature_matrix_count, target_values_train)
2
3 print('best param of n_neighbors',
    knn_count_grid.best_params_['n_neighbors'])
4 best_knn_tfidf =
    KNeighborsClassifier(n_neighbors=knn_count_grid.best_params_['n_neighbors']
    )
5
6 best_knn_tfidf.fit(train_feature_matrix_tfidf, target_values_train)
7 best_pred_knn = best_knn_tfidf.predict(test_feature_matrix__tfidf)
8
9 print(classification_report(target_values_test, best_pred_knn))

```

```

1 Fitting 5 folds for each of 6 candidates, totalling 30 fits
2 [CV] n_neighbors=2 .....
3 [CV] ..... n_neighbors=2, score=0.714, total= 0.1s
4 [CV] n_neighbors=2 .....
5 [CV] ..... n_neighbors=2, score=0.715, total= 0.1s
6 [CV] n_neighbors=2 .....

```

```

1 [Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
    workers.
2 [Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.1s remaining:
    0.0s
3 [Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.2s remaining:
    0.0s

```

```

1 [CV] ..... n_neighbors=2, score=0.722, total= 0.1s
2 [CV] n_neighbors=2 .....
3 [CV] ..... n_neighbors=2, score=0.703, total= 0.1s
4 [CV] n_neighbors=2 .....
5 [CV] ..... n_neighbors=2, score=0.709, total= 0.1s
6 [CV] n_neighbors=3 .....

```

```

1 [Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 0.2s remaining:
    0.0s

```

```

1 [CV] ..... n_neighbors=3, score=0.764, total= 0.1s
2 [CV] n_neighbors=3 .....

```

```
3 [CV] ..... n_neighbors=3, score=0.759, total= 0.1s
4 [CV] n_neighbors=3 .....
5 [CV] ..... n_neighbors=3, score=0.727, total= 0.1s
6 [CV] n_neighbors=3 .....
7 [CV] ..... n_neighbors=3, score=0.710, total= 0.1s
8 [CV] n_neighbors=3 .....
9 [CV] ..... n_neighbors=3, score=0.752, total= 0.1s
10 [CV] n_neighbors=5 .....
11 [CV] ..... n_neighbors=5, score=0.758, total= 0.1s
12 [CV] n_neighbors=5 .....
13 [CV] ..... n_neighbors=5, score=0.759, total= 0.1s
14 [CV] n_neighbors=5 .....
15 [CV] ..... n_neighbors=5, score=0.721, total= 0.2s
16 [CV] n_neighbors=5 .....
17 [CV] ..... n_neighbors=5, score=0.736, total= 0.1s
18 [CV] n_neighbors=5 .....
19 [CV] ..... n_neighbors=5, score=0.723, total= 0.2s
20 [CV] n_neighbors=7 .....
21 [CV] ..... n_neighbors=7, score=0.711, total= 0.1s
22 [CV] n_neighbors=7 .....
23 [CV] ..... n_neighbors=7, score=0.715, total= 0.1s
24 [CV] n_neighbors=7 .....
25 [CV] ..... n_neighbors=7, score=0.730, total= 0.1s
26 [CV] n_neighbors=7 .....
27 [CV] ..... n_neighbors=7, score=0.706, total= 0.1s
28 [CV] n_neighbors=7 .....
29 [CV] ..... n_neighbors=7, score=0.724, total= 0.1s
30 [CV] n_neighbors=9 .....
31 [CV] ..... n_neighbors=9, score=0.702, total= 0.1s
32 [CV] n_neighbors=9 .....
33 [CV] ..... n_neighbors=9, score=0.716, total= 0.1s
34 [CV] n_neighbors=9 .....
35 [CV] ..... n_neighbors=9, score=0.720, total= 0.1s
36 [CV] n_neighbors=9 .....
37 [CV] ..... n_neighbors=9, score=0.680, total= 0.1s
38 [CV] n_neighbors=9 .....
39 [CV] ..... n_neighbors=9, score=0.716, total= 0.1s
40 [CV] n_neighbors=11 .....
41 [CV] ..... n_neighbors=11, score=0.651, total= 0.1s
42 [CV] n_neighbors=11 .....
43 [CV] ..... n_neighbors=11, score=0.701, total= 0.1s
44 [CV] n_neighbors=11 .....
45 [CV] ..... n_neighbors=11, score=0.706, total= 0.1s
46 [CV] n_neighbors=11 .....
47 [CV] ..... n_neighbors=11, score=0.673, total= 0.1s
48 [CV] n_neighbors=11 .....
49 [CV] ..... n_neighbors=11, score=0.713, total= 0.1s
50 best param of n_neighbors 3
```

```
1 [Parallel(n_jobs=1)]: Done 30 out of 30 | elapsed: 3.1s finished
```

	precision	recall	f1-score	support
0	0.73	0.96	0.83	396
1	0.97	0.68	0.80	393
2	0.89	0.87	0.88	251
accuracy			0.83	1040
macro avg	0.86	0.84	0.84	1040
weighted avg	0.86	0.83	0.83	1040

```
1 knn_count = KNeighborsClassifier()
2 parameters = {'n_neighbors': [2, 3, 5, 7, 9, 11]}
3
4 knn_count_grid = GridSearchCV(knn_count, parameters,
    scoring='balanced_accuracy', verbose=4, cv=5)
```

```
1 knn_count_grid.fit(train_feature_matrix_count, target_values_train)
2
3 print('best param of n_neighbors',
    knn_count_grid.best_params_['n_neighbors'])
4 best_knn_count =
    KNeighborsClassifier(n_neighbors=knn_count_grid.best_params_['n_neighbors']
    )
5 print(best_knn_count)
6 best_knn_count.fit(train_feature_matrix_count, target_values_train)
7 best_knn_pred_count = best_knn_count.predict(test_feature_matrix_count)
8
9 print(classification_report(target_values_test, best_knn_pred_count))
```

```
1 Fitting 5 folds for each of 6 candidates, totalling 30 fits
2 [CV] n_neighbors=2 .....
3 [CV] ..... n_neighbors=2, score=0.714, total= 0.1s
4 [CV] n_neighbors=2 .....
```

```
1 [Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent
    workers.
2 [Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.1s remaining:
    0.0s
```



```
1 [CV] ..... n_neighbors=2, score=0.715, total= 0.1s
2 [CV] n_neighbors=2 .....
3 [CV] ..... n_neighbors=2, score=0.722, total= 0.1s
4 [CV] n_neighbors=2 .....
5 [CV] ..... n_neighbors=2, score=0.703, total= 0.1s
6 [CV] n_neighbors=2 .....
```

```
1 [Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.2s remaining:
  0.0s
2 [Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 0.4s remaining:
  0.0s
```

```
1 [CV] ..... n_neighbors=2, score=0.709, total= 0.1s
2 [CV] n_neighbors=3 .....
3 [CV] ..... n_neighbors=3, score=0.764, total= 0.1s
4 [CV] n_neighbors=3 .....
5 [CV] ..... n_neighbors=3, score=0.759, total= 0.1s
6 [CV] n_neighbors=3 .....
7 [CV] ..... n_neighbors=3, score=0.727, total= 0.2s
8 [CV] n_neighbors=3 .....
9 [CV] ..... n_neighbors=3, score=0.710, total= 0.1s
10 [CV] n_neighbors=3 .....
11 [CV] ..... n_neighbors=3, score=0.752, total= 0.1s
12 [CV] n_neighbors=5 .....
13 [CV] ..... n_neighbors=5, score=0.758, total= 0.1s
14 [CV] n_neighbors=5 .....
15 [CV] ..... n_neighbors=5, score=0.759, total= 0.1s
16 [CV] n_neighbors=5 .....
17 [CV] ..... n_neighbors=5, score=0.721, total= 0.1s
18 [CV] n_neighbors=5 .....
19 [CV] ..... n_neighbors=5, score=0.736, total= 0.1s
20 [CV] n_neighbors=5 .....
21 [CV] ..... n_neighbors=5, score=0.723, total= 0.1s
22 [CV] n_neighbors=7 .....
23 [CV] ..... n_neighbors=7, score=0.711, total= 0.1s
24 [CV] n_neighbors=7 .....
25 [CV] ..... n_neighbors=7, score=0.715, total= 0.1s
26 [CV] n_neighbors=7 .....
27 [CV] ..... n_neighbors=7, score=0.730, total= 0.1s
28 [CV] n_neighbors=7 .....
29 [CV] ..... n_neighbors=7, score=0.706, total= 0.2s
30 [CV] n_neighbors=7 .....
31 [CV] ..... n_neighbors=7, score=0.724, total= 0.1s
32 [CV] n_neighbors=9 .....
33 [CV] ..... n_neighbors=9, score=0.702, total= 0.1s
34 [CV] n_neighbors=9 .....
35 [CV] ..... n_neighbors=9, score=0.716, total= 0.1s
36 [CV] n_neighbors=9 .....
```

```

37 [CV] ..... n_neighbors=9, score=0.720, total= 0.1s
38 [CV] n_neighbors=9 .....
39 [CV] ..... n_neighbors=9, score=0.680, total= 0.2s
40 [CV] n_neighbors=9 .....
41 [CV] ..... n_neighbors=9, score=0.716, total= 0.1s
42 [CV] n_neighbors=11 .....
43 [CV] ..... n_neighbors=11, score=0.651, total= 0.2s
44 [CV] n_neighbors=11 .....
45 [CV] ..... n_neighbors=11, score=0.701, total= 0.1s
46 [CV] n_neighbors=11 .....
47 [CV] ..... n_neighbors=11, score=0.706, total= 0.1s
48 [CV] n_neighbors=11 .....
49 [CV] ..... n_neighbors=11, score=0.673, total= 0.1s
50 [CV] n_neighbors=11 .....
51 [CV] ..... n_neighbors=11, score=0.713, total= 0.1s
52 best param of n_neighbors 3
53 KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
54                      metric_params=None, n_jobs=None, n_neighbors=3, p=2,
55                      weights='uniform')

```

```

1 [Parallel(n_jobs=1)]: Done 30 out of 30 | elapsed: 3.4s finished

```

		precision	recall	f1-score	support
1					
2					
3	0	0.60	0.74	0.66	396
4	1	0.67	0.61	0.64	393
5	2	0.73	0.57	0.64	251
6					
7	accuracy			0.65	1040
8	macro avg	0.67	0.64	0.65	1040
9	weighted avg	0.66	0.65	0.65	1040

## Способ 2.

```

1 import tqdm
2 from gensim.models import Word2Vec
3 import gensim.downloader

```

```

1 /Users/snipghost/opt/anaconda3/lib/python3.7/site-
  packages/gensim/similarities/__init__.py:15: UserWarning: The
  gensim.similarities.levenshtein submodule is disabled, because the optional
  Levenshtein package <https://pypi.org/project/python-Levenshtein/> is
  unavailable. Install Levenshtein (e.g. `pip install python-Levenshtein`) to
  suppress this warning.
2 warnings.warn(msg)

```

```

1 gensim.downloader.info()
2 # glove_vectors = gensim.downloader.load('glove-twitter-25')
3 glove_vectors = gensim.downloader.load('glove-wiki-gigaword-50')

```

```

1 [=====] 100.0% 66.0/66.0MB
  downloaded

```

```

1 class GloveTokenizer:
2     def __init__(self, glove_tokenizer):
3         self.glove = glove_tokenizer
4         self.token_length = 800
5         self.embedding_size = 50
6
7     def __getitem__(self, word):
8         try:
9             vector = glove_vectors.get_vector(word).reshape(1,
10 self.embedding_size)
11         except KeyError as e:
12             vector = np.zeros((1, self.embedding_size))
13         return vector
14
15     def __padd(self, sentence):
16         padded_sentence = np.zeros((self.token_length,
17 self.embedding_size))
18         for i, token in enumerate(sentence):
19             padded_sentence[i] = token
20         return padded_sentence
21
22     def tokenize(self, sentence):
23         encoded_sentence = []
24         sentence = sentence.strip(' ').split(' ')
25         for i in sentence:
26             token = self.__getitem__(i)
27             encoded_sentence.append(token)

```

```
27 |         return np.array(self.__padd(encoded_sentence), dtype=np.float16)
```

```
1 | tokenizer = GloveTokenizer(glove_vectors)
```

```
1 | def prepare(t):
2 |     # t = ' '.join([i.strip().lower() for i in t.split(' ')])
3 |     t = re.sub(r'^a-zA-Z0-9 ', '', t)
4 |     t = re.sub('\s+', ' ', t)
5 |     lemmas = [token.lemma_.lower() for token in nlp(t) if token not in
stopwords_eng]
6 |     t = ' '.join(lemmas)
7 |     vectors = tokenizer.tokenize(t)
8 |     return vectors, len(lemmas)
9 |
10 | vectors_array_train = []
11 | labels_train = []
12 |
13 | for enum, text, label in zip(range(len(newsgroups_train.data)),
newsgroups_train.data, newsgroups_train.target):
14 |     try:
15 |         vector, length = prepare(text)
16 |         # print(vector, vector.shape)
17 |         vectors_array_train.append(vector)
18 |         labels_train.append(label)
19 |     except IndexError as e:
20 |         print(enum, e)
21 |         continue
22 |
23 |
24 | vectors_array_train = np.array(vectors_array_train)
25 | print(vectors_array_train.shape)
26 | train_data = vectors_array_train.reshape((-1,
vectors_array_train.shape[1]*vectors_array_train.shape[2]))
27 | train_data.shape
```

```
1 | 18 index 800 is out of bounds for axis 0 with size 800
2 | 47 index 800 is out of bounds for axis 0 with size 800
3 | 186 index 800 is out of bounds for axis 0 with size 800
4 | 204 index 800 is out of bounds for axis 0 with size 800
5 | 211 index 800 is out of bounds for axis 0 with size 800
6 | 214 index 800 is out of bounds for axis 0 with size 800
7 | 261 index 800 is out of bounds for axis 0 with size 800
8 | 279 index 800 is out of bounds for axis 0 with size 800
9 | 313 index 800 is out of bounds for axis 0 with size 800
10 | 318 index 800 is out of bounds for axis 0 with size 800
11 | 330 index 800 is out of bounds for axis 0 with size 800
12 | 334 index 800 is out of bounds for axis 0 with size 800
```

13	355 index 800 is out of bounds for axis 0 with size 800
14	372 index 800 is out of bounds for axis 0 with size 800
15	389 index 800 is out of bounds for axis 0 with size 800
16	405 index 800 is out of bounds for axis 0 with size 800
17	418 index 800 is out of bounds for axis 0 with size 800
18	420 index 800 is out of bounds for axis 0 with size 800
19	425 index 800 is out of bounds for axis 0 with size 800
20	442 index 800 is out of bounds for axis 0 with size 800
21	446 index 800 is out of bounds for axis 0 with size 800
22	457 index 800 is out of bounds for axis 0 with size 800
23	463 index 800 is out of bounds for axis 0 with size 800
24	485 index 800 is out of bounds for axis 0 with size 800
25	549 index 800 is out of bounds for axis 0 with size 800
26	585 index 800 is out of bounds for axis 0 with size 800
27	637 index 800 is out of bounds for axis 0 with size 800
28	640 index 800 is out of bounds for axis 0 with size 800
29	654 index 800 is out of bounds for axis 0 with size 800
30	655 index 800 is out of bounds for axis 0 with size 800
31	670 index 800 is out of bounds for axis 0 with size 800
32	697 index 800 is out of bounds for axis 0 with size 800
33	703 index 800 is out of bounds for axis 0 with size 800
34	747 index 800 is out of bounds for axis 0 with size 800
35	755 index 800 is out of bounds for axis 0 with size 800
36	756 index 800 is out of bounds for axis 0 with size 800
37	758 index 800 is out of bounds for axis 0 with size 800
38	820 index 800 is out of bounds for axis 0 with size 800
39	834 index 800 is out of bounds for axis 0 with size 800
40	877 index 800 is out of bounds for axis 0 with size 800
41	892 index 800 is out of bounds for axis 0 with size 800
42	906 index 800 is out of bounds for axis 0 with size 800
43	949 index 800 is out of bounds for axis 0 with size 800
44	955 index 800 is out of bounds for axis 0 with size 800
45	959 index 800 is out of bounds for axis 0 with size 800
46	979 index 800 is out of bounds for axis 0 with size 800
47	985 index 800 is out of bounds for axis 0 with size 800
48	1041 index 800 is out of bounds for axis 0 with size 800
49	1115 index 800 is out of bounds for axis 0 with size 800
50	1120 index 800 is out of bounds for axis 0 with size 800
51	1122 index 800 is out of bounds for axis 0 with size 800
52	1124 index 800 is out of bounds for axis 0 with size 800
53	1125 index 800 is out of bounds for axis 0 with size 800
54	1138 index 800 is out of bounds for axis 0 with size 800
55	1140 index 800 is out of bounds for axis 0 with size 800
56	1150 index 800 is out of bounds for axis 0 with size 800
57	1153 index 800 is out of bounds for axis 0 with size 800
58	1165 index 800 is out of bounds for axis 0 with size 800
59	1185 index 800 is out of bounds for axis 0 with size 800
60	1186 index 800 is out of bounds for axis 0 with size 800
61	1191 index 800 is out of bounds for axis 0 with size 800

```
62 1201 index 800 is out of bounds for axis 0 with size 800
63 1209 index 800 is out of bounds for axis 0 with size 800
64 1260 index 800 is out of bounds for axis 0 with size 800
65 1262 index 800 is out of bounds for axis 0 with size 800
66 1295 index 800 is out of bounds for axis 0 with size 800
67 1321 index 800 is out of bounds for axis 0 with size 800
68 1339 index 800 is out of bounds for axis 0 with size 800
69 1426 index 800 is out of bounds for axis 0 with size 800
70 1489 index 800 is out of bounds for axis 0 with size 800
71 1502 index 800 is out of bounds for axis 0 with size 800
72 1545 index 800 is out of bounds for axis 0 with size 800
73 1551 index 800 is out of bounds for axis 0 with size 800
74 1561 index 800 is out of bounds for axis 0 with size 800
75 (1489, 800, 50)
```

```
1 (1489, 40000)
```

```
1 vectors_array_test = []
2 labels_test= []
3
4 for enum, text, label in zip(range(len(newsgroups_test.data)),
newsgroups_test.data, newsgroups_test.target):
5     try:
6         vector, length = prepare(text)
7         vectors_array_test.append(vector)
8         labels_test.append(label)
9     except IndexError as e:
10        print(enum, e)
11        continue
```

```
1 46 index 800 is out of bounds for axis 0 with size 800
2 56 index 800 is out of bounds for axis 0 with size 800
3 69 index 800 is out of bounds for axis 0 with size 800
4 106 index 800 is out of bounds for axis 0 with size 800
5 122 index 800 is out of bounds for axis 0 with size 800
6 184 index 800 is out of bounds for axis 0 with size 800
7 282 index 800 is out of bounds for axis 0 with size 800
8 286 index 800 is out of bounds for axis 0 with size 800
9 341 index 800 is out of bounds for axis 0 with size 800
10 376 index 800 is out of bounds for axis 0 with size 800
11 402 index 800 is out of bounds for axis 0 with size 800
12 439 index 800 is out of bounds for axis 0 with size 800
```

```
13 455 index 800 is out of bounds for axis 0 with size 800
14 459 index 800 is out of bounds for axis 0 with size 800
15 505 index 800 is out of bounds for axis 0 with size 800
16 613 index 800 is out of bounds for axis 0 with size 800
17 651 index 800 is out of bounds for axis 0 with size 800
18 668 index 800 is out of bounds for axis 0 with size 800
19 794 index 800 is out of bounds for axis 0 with size 800
20 808 index 800 is out of bounds for axis 0 with size 800
21 884 index 800 is out of bounds for axis 0 with size 800
22 997 index 800 is out of bounds for axis 0 with size 800
23 999 index 800 is out of bounds for axis 0 with size 800
24 1027 index 800 is out of bounds for axis 0 with size 800
25 1034 index 800 is out of bounds for axis 0 with size 800
```

```
1 vectors_array_test = np.array(vectors_array_test)
2 test_data = vectors_array_test.reshape((-1,
    vectors_array_test.shape[1]*vectors_array_test.shape[2]))
3 test_data.shape
```

```
1 (1015, 40000)
```

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.model_selection import GridSearchCV
3 from sklearn.metrics import classification_report
4
5
6 glove_knn_clf = KNeighborsClassifier()
7 parameters = {'n_neighbors': [2, 3, 5, 7, 9]}
8
9 glove_clf_grid = GridSearchCV(glove_knn_clf, parameters, verbose=4, cv=3,
10                               scoring='balanced_accuracy', n_jobs=-1)
11
12 glove_clf_grid.fit(train_data, labels_train)
13
14 print('best param of n_neighbors',
    glove_clf_grid.best_params_['n_neighbors'])
```

```
1 Fitting 3 folds for each of 5 candidates, totalling 15 fits
```

```
1 [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
2 [Parallel(n_jobs=-1)]: Done 12 out of 15 | elapsed: 4.7min remaining:
  1.2min
3 [Parallel(n_jobs=-1)]: Done 15 out of 15 | elapsed: 5.7min finished
```

```
1 best param of n_neighbors 7
```

```
1 best_glove_knn =
  KNeighborsClassifier(n_neighbors=glove_clf_grid.best_params_['n_neighbors']
  )
2 best_glove_knn.fit(train_data, labels_train)
3 best_pred_glove_knn = best_glove_knn.predict(test_data[:800])
4
5 print(classification_report(labels_test[:800], best_pred_glove_knn))
```

	precision	recall	f1-score	support
0	0.49	0.54	0.51	316
1	0.53	0.50	0.51	305
2	0.48	0.42	0.45	179
accuracy			0.50	800
macro avg	0.50	0.49	0.49	800
weighted avg	0.50	0.50	0.50	800

```
1
```