

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ Н. Э. БАУМАНА  
Факультет информатики и систем управления  
Кафедра теоретической информатики и компьютерных технологий

Лабораторная работа №13  
по курсу «Информационный поиск»  
«Анализ тональности»

Выполнил:  
студент группы ИУ9-21М  
Беляев А. В.

Проверила:  
Лукашевич Н. В.

Москва 2019

# 1 Цель работы

Даны 2 коллекции цитат - обучающая и тестовая. Необходимо провести анализ тональности цитат и оценить работу классификатора из пакета `scikit-learn` на разных типах векторизации:

- TF-IDF
- частоты
- булевские вектора

В качестве метода классификации был выбран алгоритм AdaBoost (Adaptive Boosting). В основе алгоритма лежат т.н. «weak learners» – слабые классификаторы, которые случайным образом «угадывают» результат. Такие классификаторы объединены в «комитеты классификации». На каждой итерации алгоритма проверяются неверно классифицированные объекты, чтобы на следующей итерации алгоритма комитет классификации «сфокусировал свое внимание» на них.

## 1.1 Ход работы

```
1 from datetime import datetime
2 from xml.dom import minidom
3 import nltk
4 import numpy as np
5 from nltk.corpus import stopwords
6 from sklearn.ensemble import AdaBoostClassifier
7 from sklearn.feature_extraction.text import TfidfVectorizer,
  ↳ CountVectorizer
8 from sklearn import metrics
9 from lab16 import clean_up_sentence, normalize_sentence
10
11 TRAIN_DATA = 'news_eval_train.xml'
12 TEST_DATA = 'news_eval_test.xml'
13
14 TONALITY = ['+', '-', '0']
15 BOOL_VECT_TRUE = 1
16 BOOL_VECT_FALSE = 0
17
18 nltk.download('stopwords')
19 STOP_WORDS = stopwords.words('russian')
20 STOP_WORDS_EXCEPTIONS = [
21     'все', 'ничего', 'никогда', 'наконец', 'больше', 'хорошо', 'лучше',
22     'всегда', 'конечно', 'всю', 'такой', 'впрочем', 'так', 'можно',
23     ↳ 'даже']
24 STOP_WORDS = list(set(STOP_WORDS) - set(STOP_WORDS_EXCEPTIONS))
25
26 class Cite:
```

```

26     def __init__(self, speech_: str, eval_: str):
27         self.speech = speech_.strip()
28         self.evaluation = eval_.strip()
29         self.tokenized = tokenize(self.speech)
30         self.tokenized_str = ' '.join(self.tokenized)
31
32
33     # [[1, 2, 3], [4], [5, 6]] -> [1, 2, 3, 4, 5, 6]
34     def flatten(lst: list) -> list:
35         flat_list = []
36         for sublist in lst:
37             for item in sublist:
38                 flat_list.append(item)
39         return flat_list
40
41
42     def tokenize(text: str) -> list:
43         def remove_stopwords(sentence: str) -> list:
44             words = sentence.split(' ')
45             return list(filter(lambda w: w not in STOP_WORDS, words))
46
47         sentences = nltk.sent_tokenize(text)
48         clean_sents = list(map(lambda s: clean_up_sentence(s), sentences))
49         norm_sents = list(map(lambda s: normalize_sentence(s), clean_sents))
50         no_stopwords_sents = list(map(lambda s: remove_stopwords(s),
51                                       ↪ norm_sents))
52         return flatten(no_stopwords_sents)
53
54     def classify(vectorizer, train_cites: list, test_cites: list,
55               ↪ boolean_vectorizer=False):
56         tokenized_train = [c.tokenized_str for c in train_cites]
57         x_train = vectorizer.fit_transform(tokenized_train).toarray()
58         if boolean_vectorizer:
59             x_train = np.where(x_train > 0, BOOL_VECT_TRUE, BOOL_VECT_FALSE)
60         y_train = [c.evaluation for c in train_cites]
61
62         tokenized_test = [c.tokenized_str for c in test_cites]
63         x_test = vectorizer.transform(tokenized_test).toarray()
64         y_test = [c.evaluation for c in test_cites]
65
66         classifier = AdaBoostClassifier(n_estimators=10)
67         classifier.fit(x_train, y_train)
68         predicted = classifier.predict(x_test)

```

```

69     print(metrics.classification_report(y_test, predicted))
70     print(f'classified: {datetime.now()}')
71
72
73 def parse(filename: str) -> list:
74     xmldoc = minidom.parse(filename)
75     document = xmldoc.getElementsByTagName('document')[0]
76     sentences = document.getElementsByTagName('sentence')
77
78     citations = []
79     for s in sentences:
80         speech =
81             ↪ s.getElementsByTagName('speech')[0].childNodes[0].nodeValue
82         evaluation =
83             ↪ s.getElementsByTagName('evaluation')[0].childNodes[0].nodeValue.strip()
84         if evaluation in TONALITY:
85             citations.append(Cite(speech, evaluation))
86     return citations
87
88 def main():
89     train_cites = parse(TRAIN_DATA)
90     test_cites = parse(TEST_DATA)
91
92     print(f'classify: {datetime.now()}')
93     vectorizer = TfidfVectorizer()
94     # vectorizer = CountVectorizer()
95
96     classify(vectorizer, train_cites, test_cites,
97             ↪ boolean_vectorizer=False)
98
99 if __name__ == '__main__':
100     main()

```

## 2 Результаты

Входные XML файлы были разобраны, цитаты с неоднозначной тональностью были исключены, а остальные были нормализованы.

Из-за того, что классификаторы, состоящие в комитете классификации не намного лучше случайного гадания (как следует из документации), алгоритм часто используется в сочетании с другими методами. В данной работе AdaBoost был применен отдельно от других алгоритмов. Можно выделить его плюс по сравнению с другими методами - главным настраиваемым параметром метода является `n_estimators` - количество классификаторов в комитете.

Были проверены результаты работы с комитетом из 10, 100 и 300 классификато-

ров.

На листингах ниже представлены результаты классификации с использованием разных способов векторизации.

TF-IDF:

1	===== N = 10 =====
2	precision recall f1-score support
3	+ 0.56 0.30 0.39 1448
4	- 0.45 0.89 0.59 1890
5	0 0.00 0.00 0.00 1235
6	
7	micro avg 0.47 0.47 0.47 4573
8	macro avg 0.34 0.40 0.33 4573
9	weighted avg 0.36 0.47 0.37 4573
10	===== N = 100 =====
11	+ 0.54 0.48 0.51 1448
12	- 0.50 0.76 0.61 1890
13	0 0.40 0.14 0.21 1235
14	
15	micro avg 0.50 0.50 0.50 4573
16	macro avg 0.48 0.46 0.44 4573
17	weighted avg 0.49 0.50 0.47 4573
18	===== N = 300 =====
19	+ 0.52 0.50 0.51 1448
20	- 0.54 0.66 0.59 1890
21	0 0.40 0.28 0.33 1235
22	
23	micro avg 0.51 0.51 0.51 4573
24	macro avg 0.48 0.48 0.48 4573
25	weighted avg 0.49 0.51 0.49 4573

Частотная векторизация:

1	===== N = 10 =====
2	precision recall f1-score support
3	+ 0.57 0.31 0.40 1448
4	- 0.45 0.89 0.59 1890
5	0 0.00 0.00 0.00 1235
6	
7	micro avg 0.47 0.47 0.47 4573
8	macro avg 0.34 0.40 0.33 4573
9	weighted avg 0.36 0.47 0.37 4573
10	===== N = 100 =====
11	+ 0.53 0.55 0.54 1448
12	- 0.52 0.74 0.61 1890
13	0 0.44 0.14 0.21 1235

14					
15	micro avg	0.52	0.52	0.52	4573
16	macro avg	0.50	0.48	0.46	4573
17	weighted avg	0.50	0.52	0.48	4573
18	===== N = 300 =====				
19	+	0.55	0.54	0.54	1448
20	-	0.57	0.69	0.63	1890
21	0	0.41	0.28	0.33	1235
22					
23	micro avg	0.53	0.53	0.53	4573
24	macro avg	0.51	0.50	0.50	4573
25	weighted avg	0.52	0.53	0.52	4573

Булевская векторизация:

1	===== N = 10 =====				
2		precision	recall	f1-score	support
3	+	0.57	0.31	0.40	1448
4	-	0.45	0.89	0.59	1890
5	0	0.00	0.00	0.00	1235
6					
7	micro avg	0.47	0.47	0.47	4573
8	macro avg	0.34	0.40	0.33	4573
9	weighted avg	0.36	0.47	0.37	4573
10	===== N = 100 =====				
11	+	0.56	0.51	0.53	1448
12	-	0.51	0.78	0.62	1890
13	0	0.44	0.14	0.21	1235
14					
15	micro avg	0.52	0.52	0.52	4573
16	macro avg	0.50	0.47	0.45	4573
17	weighted avg	0.51	0.52	0.48	4573
18	===== N = 300 =====				
19	+	0.54	0.54	0.54	1448
20	-	0.58	0.68	0.63	1890
21	0	0.45	0.32	0.37	1235
22					
23	micro avg	0.54	0.54	0.54	4573
24	macro avg	0.52	0.52	0.52	4573
25	weighted avg	0.53	0.54	0.53	4573

Поскольку комитет из 10 слабых классификаторов дает посредственные результаты, а комитет из 300 слишком долго работает (порядка 30 минут), считаем комитет из 100 классификаторов оптимальным.

Исключим из `STOP_WORDS` слова, содержащиеся в `STOP_WORDS_EXCEPTIONS`, чтобы по-другому учитывать эмоциональную окраску цитат.

Проверим результаты для Count Vectorizer'a, т.к. он дал наилучший результат по отношению к временным затратам. В комитете 100 классификаторов:

1		<b>precision</b>	<b>recall</b>	<b>f1-score</b>	<b>support</b>
2	+	0.53	0.53	0.53	1448
3	-	0.52	0.75	0.61	1890
4	0	0.43	0.12	0.19	1235
5					
6	micro avg	0.51	0.51	0.51	4573
7	macro avg	0.49	0.47	0.45	4573
8	weighted avg	0.50	0.51	0.47	4573

### 3 Выводы

Был произведен анализ тональности с помощью метода AdaBoost пакета scikit-learn. Выбранный метод показал сравнительно слабые результаты, как и было заявлено в документации к нему. При попытке подогнать параметры метода под результат, ожидаемо выигрываем в одних результатах и проигрываем в других.

Тем не менее, этот метод может быть хорошим дополнением к другим методам, т.к. слабо подвержен проблеме переобучения.