

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



## **Рубежный контроль № 2**

**По курсу «методы машинного обучения в АСОИУ»**

**Выполнил:**

студент ИУ5-23М  
Семенов И.А.

**Проверил:**

Гапанюк Ю.Е.

Подпись:

11.05.2024

---

Москва, 2024

## Задание

Необходимо решить задачу классификации текстов на основе любого выбранного Вами датасета (кроме примера, который рассматривался в лекции). Классификация может быть бинарной или многоклассовой. Целевой признак из выбранного Вами датасета может иметь любой физический смысл, примером является задача анализа тональности текста.

Необходимо сформировать два варианта векторизации признаков - на основе CountVectorizer и на основе TfidfVectorizer.

В качестве классификаторов необходимо использовать два классификатора по варианту для Вашей группы:

Группа	Классификатор №1	Классификатор №2
ИУ5-23М, ИУ5И-23М	<u>LinearSVC</u>	<u>LogisticRegression</u>

Для каждого метода необходимо оценить качество классификации. Сделайте вывод о том, какой вариант векторизации признаков в паре с каким классификатором показал лучшее качество.

## Решение

Загружаем датасет

```
| from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
| from sklearn.model_selection import train_test_split
| from sklearn.svm import LinearSVC
| from sklearn.linear_model import LogisticRegression
| from sklearn.metrics import accuracy_score
| import pandas as pd
| import time

| df = pd.read_csv('data.csv')
```

	yy-mm-dd	sense	from	to	description	name
0	2023-12-19 06:40:41	{'class': 'negative', 'polarity': -0.1, 'subje...	CryptoNews	altcoin	Grayscale CEO Michael Sonnenshein believes the...	Grayscale CEO Calls for Simultaneous Approval ...
1	2023-12-19 06:03:24	{'class': 'neutral', 'polarity': 0.0, 'subject...	CryptoNews	blockchain	In an exclusive interview with CryptoNews, Man...	Indian Government is Actively Collaborating Wi...

Разделяем данные на тестовую и обучающую выборку и векторизуем их с помощью CountVectorizer и TfidfVectorizer

```
X, Y = df['text'], df['source']
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=42)

time_arr = []
```

```
# CountVectorizer
count_vect = CountVectorizer()
X_train_counts = count_vect.fit_transform(X_train)
X_test_counts = count_vect.transform(X_test)
```

```
# TfidfVectorizer
tfidf_vect = TfidfVectorizer()
X_train_tfidf = tfidf_vect.fit_transform(X_train)
X_test_tfidf = tfidf_vect.transform(X_test)
```

## Обучаем классификаторы CountVectorizer

```
# LinearSVC
gbc = LinearSVC()
start_time = time.time()
gbc.fit(X_train_counts, y_train)
train_time = time.time() - start_time
time_arr.append(train_time)
pred_gbc_counts = gbc.predict(X_test_counts)
print("Точность (CountVectorizer + LinearSVC):", accuracy_score(y_test, pred_gbc_counts))

# Logistic Regression
lr = LogisticRegression(max_iter=1500)
start_time = time.time()
lr.fit(X_train_counts, y_train)
train_time = time.time() - start_time
time_arr.append(train_time)
pred_lr_counts = lr.predict(X_test_counts)
print("Точность (CountVectorizer + LogisticRegression):", accuracy_score(y_test, pred_lr_counts))
```

Точность (CountVectorizer + LinearSVC): 0.605543710021322

Точность (CountVectorizer + LogisticRegression): 0.6268656716417911

## Обучаем классификаторы TfidfVectorizer

```
# Произведем обучения двух классификаторов (по варианту) для TfidfVectorizer

# LinearSVC
gbc = LinearSVC()
start_time = time.time()
gbc.fit(X_train_tfidf, y_train)
train_time = time.time() - start_time
time_arr.append(train_time)
pred_gbc_tfidf = gbc.predict(X_test_tfidf)
print("Точность (TfidfVectorizer + LinearSVC):", accuracy_score(y_test, pred_gbc_tfidf))

# Logistic Regression
lr = LogisticRegression(max_iter=1500)
start_time = time.time()
lr.fit(X_train_tfidf, y_train)
train_time = time.time() - start_time
time_arr.append(train_time)
pred_lr_tfidf = lr.predict(X_test_tfidf)
print("Точность (TfidfVectorizer + LogisticRegression):", accuracy_score(y_test, pred_lr_tfidf))
```

Точность (TfidfVectorizer + LinearSVC): 0.5906183368869936  
Точность (TfidfVectorizer + LogisticRegression): 0.605543710021322

Выводим отсортированные данные

```
from tabulate import tabulate

data = [
    ["(CountVectorizer + LogisticRegression)", accuracy_score(y_test, pred_lr_counts), time_arr[0]],
    ["(CountVectorizer + LinearSVC)", accuracy_score(y_test, pred_gbc_counts), time_arr[1]],
    ["(TfidfVectorizer + LogisticRegression)", accuracy_score(y_test, pred_lr_tfidf), time_arr[2]],
    ["(TfidfVectorizer + LinearSVC)", accuracy_score(y_test, pred_gbc_tfidf), time_arr[3]]
]

sorted_data = sorted(data, key=lambda x: x[1], reverse=True)
```

Связка	Точность валидации	Время обучения
(CountVectorizer + LogisticRegression)	0.626866	0.13026
(CountVectorizer + LinearSVC)	0.605544	1.3938
(TfidfVectorizer + LogisticRegression)	0.605544	0.10843
(TfidfVectorizer + LinearSVC)	0.590618	0.770637