

Наивный Байесовский классификатор.

В этом ноутбуке мы разберём один из самых простых методов классификации - наивный байесовский классификатор. Однако в некоторых задачах он работает даже лучше других, более сложных моделей. В любом случае, наивный байесовский классификатор содержит в себе важные теоретические идеи, поэтому с ним в любом случае полезно ознакомиться.

Применение при детекции спама.

Данные для решения задачи детекции спама можно сделать следующим образом: взять набор размеченных текстовых сообщений, часть которых размечена как спам, а остальные - как не спам, зафиксировать словарь (самый простой вариант - взять все слова, встречающиеся в наборе текстовых сообщений) и преобразовать текстовые данные в целочисленные, посчитав для каждого слова из словаря, встречается ли оно в данном сообщении. А на этих данных уже можно обучить наивный байесовский классификатор.

При реализации класса для наивного байесовского классификатора надо помнить один очень важный на практике момент: произведение вероятностей большого количества чисел может очень быстро сравняться с нулем при вычислении на компьютере, так как компьютеру может не хватить вычислительной точности. Поэтому при реализации стоит использовать логарифмы вероятностей.

Применим наивный байесовский классификатор к конкретному датасету

<https://archive.ics.uci.edu/ml/datasets/SMS+Spam+Collection>
(<https://archive.ics.uci.edu/ml/datasets/SMS+Spam+Collection>).

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import accuracy_score, f1_score
5 from sklearn.model_selection import train_test_split
```

1. Чтение данных.

In [32]:

```
1 labels = []
2 messages = []
3
4 with open('SMSSpamCollection.txt', 'r') as fin:
5     for line in fin:
6         cur_label, cur_message = line.split('\t')
7         labels.append(cur_label)
8         messages.append(cur_message)
```

In [33]:

```
1 raw_df = pd.DataFrame()
2 raw_df['data'] = messages
3 raw_df['label'] = labels
4 raw_df.head()
```

Out[33]:

	data	label
0	Go until jurong point, crazy.. Available only ...	ham
1	Ok lar... Joking wif u oni...\n	ham
2	Free entry in 2 a wkly comp to win FA Cup fina...	spam
3	U dun say so early hor... U c already then say...	ham
4	Nah I don't think he goes to usf, he lives aro...	ham

В датасете метки бывают 2 видов:

- ham - означает, что сообщение - не спам,
- spam - означает, что сообщение - спам.

2. Предобработка данных.

Очевидно, что сразу в таком виде нельзя передавать данные наивному байесовскому классификатору. Их надо привести к численному виду.

Столбец `label` привести к численному виду можно очень просто.

In [34]:

```
1 raw_df['label'] = (raw_df['label'] == 'spam') * 1
2 raw_df.head()
```

Out[34]:

	data	label
0	Go until jurong point, crazy.. Available only ...	0
1	Ok lar... Joking wif u oni...\n	0
2	Free entry in 2 a wkly comp to win FA Cup fina...	1
3	U dun say so early hor... U c already then say...	0
4	Nah I don't think he goes to usf, he lives aro...	0

Для преобразования текстовых сообщений воспользуемся `CountVectorizer`, работающий по принципу мешка слов (bag of words). Он имеет следующие гиперпараметры:

- `max_df` -- максимальная доля сообщений, в которых может встретиться слово из словаря (такой параметр может быть полезен для борьбы со стоп-словами). То есть в словарь не включаются слишком частые слова.
- `min_df` -- минимальная доля сообщений, в которых может встретиться слово из словаря. То есть в словарь не включаются слишком редкие слова.

- `max_features` -- максимальное возможное число слов в словаре (берётся `max_features` наиболее частых слов).
- `stop_words` -- можно просто взять и задать стоп-слова, которые не будут добавлены в словарь ни при каких обстоятельствах.

In [35]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
```

In [36]:

```
1 vectorizer = CountVectorizer(min_df=0.03)
2 transformed_data = vectorizer.fit_transform(messages).toarray()
```

Напечатаем весь мешок слов и их количество:

In [37]:

```
1 print(len(vectorizer.get_feature_names()))
2 print(vectorizer.get_feature_names())
```

```
66
['all', 'am', 'and', 'are', 'at', 'be', 'but', 'call', 'can', 'come',
'day', 'do', 'for', 'free', 'from', 'get', 'go', 'good', 'got', 'gt',
'have', 'he', 'how', 'if', 'in', 'is', 'it', 'its', 'just', 'know', 'l
ike', 'll', 'love', 'lt', 'me', 'my', 'no', 'not', 'now', 'of', 'ok',
'on', 'only', 'or', 'out', 'send', 'so', 'text', 'that', 'the', 'the
n', 'there', 'this', 'time', 'to', 'up', 'ur', 'want', 'was', 'we', 'w
hat', 'when', 'will', 'with', 'you', 'your']
```

Посмотрим на преобразованные данные

которая восстановит априорное распределение по переданным в функцию данным.

В нашей текущей задаче для признаков, описывающих количество вхождений каждого слова из словаря в сообщение, логично использовать `MultinomialNB`. Однако после мы сравним точность предсказаний `MultinomialNB` с точностью предсказаний `BernoulliNB` для бинарных признаков (каждый признак является индикатором того, присутствует ли данное слово из словаря в сообщении).

In [39]:

```
1 from sklearn.naive_bayes import MultinomialNB
2
3 multinomial_nb = MultinomialNB()
```

In [40]:

```
1 transformed_data
```

Out[40]:

```
array([[0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       ...,
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 1, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0]])
```

Как обычно, разделим данные на обучающую выборку и на тестовую.

In [41]:

```
1 X_train, X_test, y_train, y_test \
2   = train_test_split(transformed_data, raw_df['label'],
3                       random_state=42)
```

Обучаем модель и смотрим качество на тестовой выборке

In [42]:

```
1 multinomial_nb.fit(X_train, y_train)
2 predictions = multinomial_nb.predict(X_test)
3
4 print('accuracy:', accuracy_score(predictions, y_test))
5 print('f1 score:', f1_score(predictions, y_test))
```

```
accuracy: 0.9440459110473458
f1 score: 0.7845303867403315
```

Результат получился весьма неплохой.

А теперь посмотрим, как с этой же задачей справится наивный байесовский классификатор на бинарных данных.

In [43]:

```
1 X_train = (X_train > 0) * 1
2 X_test = (X_test > 0) * 1

[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 1 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
```

In [44]:

```
1 from sklearn.naive_bayes import BernoulliNB
2
3 bernoulli_nb = BernoulliNB()
```

In [45]:

```
1 bernoulli_nb.fit(X_train, y_train)
2 predictions = bernoulli_nb.predict(X_test)
3
4 print('accuracy:', accuracy_score(predictions, y_test))
5 print('f1 score:', f1_score(predictions, y_test))
```

```
accuracy: 0.9497847919655668
f1 score: 0.8055555555555556
```

Вывод.

Результат получился достаточно неожиданный. Наивный байесовский классификатор, обученный на бинаризованных данных показал более высокую точность классификации.

4. Большой размер словаря

А теперь посмотрим, что будет, если мы возьмём другое количество слов для словаря.

It's gonna be huge!

In [48]:

```
1 huge_vectorizer = CountVectorizer()
2 huge_data = huge_vectorizer.fit_transform(messages).toarray()
3 print(huge_data.shape)
```

```
(5574, 8713)
```

In [49]:

```
1 X_train, X_test, y_train, y_test = train_test_split(
2     huge_data, raw_df['label'], random_state=42
3 )
```

In [50]:

```
1 multinomial_nb.fit(X_train, y_train)
2 predictions = multinomial_nb.predict(X_test)
3
4 print('accuracy:', accuracy_score(predictions, y_test))
5 print('f1 score:', f1_score(predictions, y_test))
```

accuracy: 0.9849354375896701
f1 score: 0.9454545454545454

In [51]:

```
1 bernoulli_nb.fit(X_train, y_train)
2 predictions = bernoulli_nb.predict(X_test)
3
4 print('accuracy:', accuracy_score(predictions, y_test))
5 print('f1 score:', f1_score(predictions, y_test))
```

accuracy: 0.9806312769010043
f1 score: 0.9247910863509748

Вывод.

От увеличения количества рассматриваемых слов в данном случае точность предсказаний возрасла как для наивного байесовского классификатора над категориальными признаками, так и для классификатора над бинарными признаками.

Бонусная часть

1. Рассмотрите различные способы построения словаря для классификации, например, установив некоторые параметры класса `CountVectorizer`.
2. Попытайтесь улучшить точность классификации (на тестовой выборке) наивного байесовского классификатора за счёт изменения гиперпараметров классификатора.
3. Решите задачу детекции спама при помощи некоторого другого известного классификатора: логистической регрессии, kNN и сравните точность предсказаний с наивным байесовским классификатором.