

Join GitHub today

Dismiss

GitHub is home to over 40 million developers working together to host and review code, manage projects, and build software together.

Sign up

Branch: master ▼

Find file

Copy path

[Machine-learning-python](#) / Анализ текстов.py



sergeybelov должен теперь работать и давать правильный ответ на ДЗ

32f6878 on 2 Mar 2017

[1 contributor](#)

Raw

BlameHistory



107 lines (87 sloc) | 7.03 KB

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Tue Feb 28 15:08:44 2017
4
5  @author: tehn-11
6  """
7
8  #=====
9  # Анализ текстов
10 #=====
11 from sklearnimportdatasets
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 from sklearn.model_selection import GridSearchCV
14 from sklearn.svm import SVC
15 import numpy as np
16 import pandas as pd
17 from sklearn.model_selection import KFold
18
19 #=====
20 # Для начала вам потребуется загрузить данные.
21 # В этом задании мы воспользуемся одним из датасетов, доступных в scikit-learn'e – 20 newsgroup
22 # Для этого нужно воспользоваться модулем datasets:
23 #=====
24 newsgroups = datasets.fetch_20newsgroups(
25     subset='all',
26     categories=['alt.atheism', 'sci.space']
```

```

27         )
28
29     #=====
30     # Вычислите TF-IDF-признаки для всех текстов.
31     # Обратите внимание, что в этом задании мы предлагаем вам вычислить TF-IDF по всем данным.
32     # При таком подходе получается, что признаки на обучающем множестве используют информацию из те
33     # поскольку мы не используем значения целевой переменной из теста.
34     # На практике нередко встречаются ситуации, когда признаки объектов тестовой выборки известны н
35     # и поэтому можно ими пользоваться при обучении алгоритма.
36     #=====
37
38     #В Scikit-Learn это реализовано в классе
39     #sklearn.feature_extraction.text.TfidfVectorizer.
40     #Преобразование обучающей выборки нужно делать с помощью функции fit_transform, тестовой – с po
41     y_train = newsgroups.target#Класс
42     X_train = newsgroups.data#Характеристики
43
44     vectorizer = TfidfVectorizer()
45     dataMatrix=vectorizer.fit_transform(X_train).toarray()#матрица объектов по словам, в ячейках ве
46
47     #idf = vectorizer.idf_
48     words=vectorizer.get_feature_names()
49     #tf_idf=dict(zip(words, idf))#токены с весами
50
51
52     #=====
53     # Подберите минимальный лучший параметр C из множества [10^-5, 10^-4, ... 10^4, 10^5] для SVM с
54     # при помощи кросс-валидации по 5 блокам. Укажите параметр random_state=241 и для SVM, и для KF
55     # В качестве меры качества используйте долю верных ответов (accuracy).
56     #=====
57
58     grid = {'C': np.power(10.0, np.arange(-5, 6))}
59     cv = KFold(n_splits=5, shuffle=True, random_state=241)
60     clf = SVC(kernel='linear', random_state=241)
61     gs = GridSearchCV(clf, grid, scoring='accuracy', cv=cv)
62     gs.fit(dataMatrix, y_train)#внимание очень долго работает на тестовых данных
63
64     #записываем параметры в массив
65     #validationTest=dict(zip(gs.cv_results_['mean_test_score'], gs.cv_results_['params']))#получаем
66     #validationTestData=pd.DataFrame(data=validationTest).transpose()#создаем датафрейм с транспонир
67     #validationTestData.sort_index(ascending=[False],inplace=True)
68
69     #У GridSearchCV есть поле best_estimator_,
70     #которое можно использовать, чтобы не обучать заново классификатор с оптимальным параметром.
71     bestC=gs.best_estimator_.C
72
73     #устарело
74     #for a in gs.grid_scores_:
75     #     validationTest[a.mean_validation_score]=a.parameters# – оценка качества по кросс-валидаци
76
77     #=====
78     # Обучите SVM по всей выборке с оптимальным параметром C, найденным на предыдущем шаге.

```

```

79 #=====
80 clf = SVC(C=bestC,kernel='linear', random_state=241)
81 result=clf.fit(dataMatrix, y_train)
82
83
84 #=====
85 # Найдите 10 слов с наибольшим абсолютным значением веса (веса хранятся в поле coef_ у svm.SVC)
86 # Они являются ответом на это задание. Укажите эти слова через запятую или пробел, в нижнем регистре
87 #в лексикографическом порядке.
88 #=====
89 weights=[]
90 for element in result.coef_.T:#с транспонированием, иначе все коэффициенты идут строкой
91     weights.append(abs(element[0]))
92
93 combine=dict(zip(words,weights))
94 bestWords=pd.DataFrame(data=combine,index=[0]).transpose()#с транспонированием, иначе все коэффициенты
95 bestWords.columns = ['weights']#переименуем колонки
96 bestWords.sort_values(['weights'], ascending=[False],inplace=True)
97
98 topTenWordsCollection=bestWords.head(10).index.values#берем значения 10 индексов (там слова)
99 newWords=pd.DataFrame(data=topTenWordsCollection)
100 newWords.columns = ['words']#переименуем колонки
101 newWords.words.astype(str)
102 newWords=newWords.apply(lambda x: x[0].lower(),axis =1)#приводим к нижнему регистру, имя колонки
103 newWords.sort_values(0, ascending=[True],inplace=True)#сортируем в лексикографическом порядке.
104
105 collection=newWords.values#получаем значения массива
106 print (' '.join(map(lambda x: x, collection)))#сливаем в строку каждый элемент

```