

資料探勘研究與實務 HW3 - Sentiment Analysis 1

資管碩一 0853412 吳宛儒

1. 資料前處理

先載入需要的套件：

```
8 import pandas as pd
9 from nltk.corpus import stopwords
10 from sklearn.feature_extraction.text import CountVectorizer
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 from sklearn.ensemble import AdaBoostClassifier
13 import xgboost as xgb
14 from sklearn.metrics import classification_report
```

- 讀取"training_label.txt"與"testing_label.txt"並利用分割符號切割字串、建立 train&test 之 DataFrame
- 去除停頓詞 stop words

使用套件：

```
8 import pandas as pd
```

定義 load data 的 function：

```
16 #%%
17 def load_data(dataset):
18     if dataset == 'train':
19         data = pd.read_csv('training_label.txt', sep='\t', header=None)
20     elif dataset == 'test':
21         data = pd.read_csv('testing_label.txt', sep='\t', header=None)
22     data.columns = ['text']
23     data['label'] = data['text'].str[0]
24     data['data'] = data['text'].str[10:]
25     del data['text']
26     return data
```

開始 load data，定義 train 以及 test set，train set 只取 train data 前 10000 筆。透過 split_data function 將 data 跟 label 切開來。

```
28 #%%
29 # Load data
30 train = load_data('train')
31 train = train[:][:10000]
32 test = load_data('test')
33
34 def split_data(dataset):
35     X, y = [], []
36     X = dataset['data'].values
37     y = dataset['label'].values
38     return X, y
39
40 trainX, trainy, testX, testy = [], [], [], []
41 trainX, trainy = split_data(train)
42 testX, testy = split_data(test)
```

接著定義 stop words，考量到自己定義 stop words 的 list 有點不切實際，使用到了 nltk 對於 english 的 stop words 定義，將 stop words 設定為 english 的 stop words，並添加其他符號。

註解：這邊 stop word 可以預先定義(如下圖)，也可以在後面的 CountVectorizer 或 TfidfVectorizer 中定義(都會寫出來，但最後跑程式的時候是寫在 Tfidf 中)

```
9 from nltk.corpus import stopwords
45 stop_words = set(stopwords.words('english'))
46
47 for w in ['!', ',', '.', '?', '-s', '-ly', '</s>', 's']:
48     stop_words.add(w)
```

CountVectorizer

原本有使用 CountVectorizer，但是後來在網路上查到其功能 TfidfVectorizer 之 transform 已有包含在內，因此這部分的程式碼在最後並沒有使用。

```
10 from sklearn.feature_extraction.text import CountVectorizer
50 ###
51 # 這個部分 tf-idf 的 transform 有包含，所以可以不用做這個 block
52 cv = CountVectorizer(stop_words='english')
53 X = cv.fit_transform(trainX)
54 cool_sw = X.toarray()
55
56 feature_sw = cv.get_feature_names()
```

c. 文字探勘前處理，將文字轉換成向量，像是常見的方法 tf-idf、word2vec... 等

使用套件：TfidfVectorizer from sklearn

這邊選擇使用 tf-idf 來將文字轉成向量，要設定 stop words 為 english 的 stop words。

```
11 from sklearn.feature_extraction.text import TfidfVectorizer
60 # tf-idf
61 vectorizer = TfidfVectorizer(stop_words='english')
62 ldf_train = vectorizer.fit_transform(trainX)
63 ldf_test = vectorizer.transform(testX)
64 tfidf_feature = vectorizer.get_feature_names()
65 print(ldf_train.shape)
66
```

2. 建模：

AdaBoost

使用套件：AdaBoostClassifier from sklearn

```
12 from sklearn.ensemble import AdaBoostClassifier
```

建立 classifier，將剛才 tf-idf 做出來的資料拿進去 classifier 做 fit，最後 predict test data 的 label 並印出。

```
67 # ada
68 ada = AdaBoostClassifier(n_estimators=100, random_state=0)
69 ada.fit(ldf_train, trainy)
70 pred_ada = ada.predict(ldf_test)
```

xgboost

先安裝 XGBoost。

```
13 import xgboost as xgb
```

建立 XGBoost Classifier，接著用 tf-idf 結果 fit，接著再 predict。

```
72 # xgbc
73 xgbc = xgb.XGBClassifier()
74 xgbc.fit(ldf_train, trainy)
75 xgbc.score(ldf_test, testy)
76 pred_xgbc = xgbc.predict(ldf_test)
```

3. 評估模型

利用 "testing_label.txt" 的資料對 2. 所建立的模型進行測試，並計算 Accuracy、Precision、Recall、F-measure

使用套件：sklearn.metrics 中的 classification_report

```
14 from sklearn.metrics import classification_report
```

```
78 print('Adaboost:', classification_report(testy, pred_ada))
79 print('Adaboost:', classification_report(testy, pred_xgbc))
```

AdaBoost 結果：

Adaboost:		precision	recall	f1-score	support
	0	0.70	0.43	0.53	37
	1	0.69	0.87	0.77	53
	accuracy			0.69	90
	macro avg	0.69	0.65	0.65	90
	weighted avg	0.69	0.69	0.67	90

XGBoost 結果：

XGboost:		precision	recall	f1-score	support
	0	0.45	0.24	0.32	37
	1	0.60	0.79	0.68	53
	accuracy			0.57	90
	macro avg	0.53	0.52	0.50	90
	weighted avg	0.54	0.57	0.53	90