

資料探勘研究與實務 Data Mining

0853426 陳紹雲

一、 資料前處理-讀入資料

利用迴圈以行依序讀入串列中，分別對 Train, Test Data 做處理，最後轉成 DataFrame

```
In [1]: 1 import numpy as np
2 import pandas as pd
3
4 res = [];count = 0
5 with open('data/training_label.txt', 'r' , encoding='utf-8') as fn:
6     for line in fn:
7         line=line.strip('\n')
8         if line != "":
9             line_list = str(line).split("+++$$$+++")
10             line_list[1] = line_list[1].strip()
11             res.append(line_list)
12             count += 1
13             if(count>=10000):
14                 break
```

二、 資料前處理- Data Cleaning

主要對資料做三種保留方式，

1. 用 bs4 套件去除 html tags。
2. 用 re 抓 emoji，保留在 text
3. 最後再把所有單詞轉成小寫把 emoji 附加在最後面

```

4 def preprocessor(text):
5     # remove HTML tags
6     text = BeautifulSoup(text, 'html.parser').get_text()
7     |
8     # regex for matching emoticons, keep emoticons, ex: :), :-P, :-D
9     r = '(:|;|=|X)(?-)?(?:\)|\(|D|P)'
10    emoticons = re.findall(r, text)
11    text = re.sub(r, '', text)
12
13    # convert to lowercase and append all emoticons behind (with space in between)
14    # replace('-', '') removes nose of emoticons
15    text = re.sub('[\W]+', ' ', text.lower()) + ' '.join(emoticons).replace('-', '')
16    return text

```

三、 資料前處理- Word Stemming

製作一個可以拆解句子成 token 的函式，且跟 StopWord 比對

```
4 def tokenizer_stem(text):
5     porter = PorterStemmer()
6     return [porter.stem(word) for word in re.split('\s+', text.strip()) \
7             if word not in stop and re.match('[a-zA-Z]+', text)] # re.match只比對字串頭
```

四、 CountVectorizer

利用字頻建立 Token of Bags.

```
5 doc = train['review']
6 count = CountVectorizer(preprocessor=preprocessor, tokenizer=tokenizer_stem)
7 doc_bag = count.fit_transform(doc).toarray() # every data features
```

五、 TfidfVectorizer

建立 idf scores vector，其中考慮 term-frequency (TF) as BoW 也考

慮到 the document-frequency (DF)

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf = TfidfVectorizer(preprocessor=preprocessor, tokenizer=tokenizer_stem)
3 tfidf.fit(doc)
```

六、 Feature Hashing

將每個詞彙散列到具有固定數量的 bucket 的 table 中來減少維度詞彙空

間。(降低內存所需，但是犧牲掉 IDF 加權)

```
In [9]: 1 from sklearn.feature_extraction.text import HashingVectorizer
2
3 hashvec = HashingVectorizer(n_features=2**10,
4                             preprocessor=preprocessor,
5                             tokenizer=tokenizer_stem)
6
7 doc_hash = hashvec.transform(doc)
8 print(doc_hash.shape)

(10000, 1024)
```

七、利用 Pipeline 建模

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.ensemble import AdaBoostClassifier
3 from sklearn.datasets import make_classification
4 from sklearn.model_selection import cross_val_score
5 from xgboost import XGBClassifier
6
7 names = ['AdaBoostClassifier+preprocess+hash', 'AdaBoostClassifier+preprocess', 'XGBClassifier+preprocess+hash', 'XGBClassifier+preprocess']
8 pipe1 = Pipeline([('vect', HashingVectorizer(n_features=2**10,
9                                           preprocessor=preprocessor,
10                                          tokenizer=tokenizer_stem)),
11                  ('clf', AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None,
12                                             learning_rate=1.0, n_estimators=50, random_state=0))])
13
14 pipe2 = Pipeline([('vect', TfidfVectorizer(preprocessor=preprocessor,
15                                           tokenizer=tokenizer_stem)),
16                  ('clf', AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None,
17                                             learning_rate=1.0, n_estimators=50, random_state=0))])
18
19 pipe3 = Pipeline([('vect', HashingVectorizer(n_features=2**10,
20                                           preprocessor=preprocessor,
21                                          tokenizer=tokenizer_stem)),
22                  ('clf', XGBClassifier(n_estimators=50, random_state=0))])
23
24 pipe4 = Pipeline([('vect', TfidfVectorizer(preprocessor=preprocessor,
25                                           tokenizer=tokenizer_stem)),
26                  ('clf', XGBClassifier(n_estimators=50, random_state=0))])
```

七、 交叉驗證(10-fold)

```
1 print('[auc (10-fold cv)]')
2 for name, clf in zip(names, [pipe1, pipe2, pipe3, pipe4]):
3     scores = cross_val_score(estimator=clf, X=train['review'], y=train['sentiment'], \
4                             cv=10, scoring='roc_auc')
5     print('%s: %.3f (+/-%.3f)' % (name, scores.mean(), scores.std()))
```

```
[auc (10-fold cv)]
AdaBoostClassifier+preprocess+hash: 0.710 (+/-0.014)
AdaBoostClassifier+preprocess: 0.724 (+/-0.014)
XGBClassifier+preprocess+hash: 0.704 (+/-0.014)
XGBClassifier+preprocess: 0.710 (+/-0.018)
```

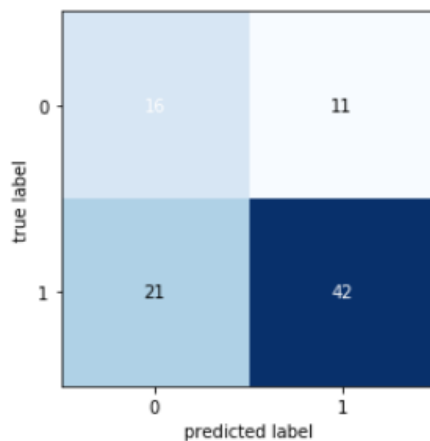
這邊發現 Adaboost 加上有做過預處理的效果最好。

八、 結果(抓兩個模型+預處理的做比較)

AdaBoostClassifier+preprocess

```
1 pipe2.fit(X=train['review'], y=train['sentiment'])
2 y_pred = pipe2.predict(test['review']).astype('int64')
3 y_test = test['sentiment'].values.astype('int64')
4 shen_bing_how_shy(y_pred, y_test)
```

Accuracy : 0.644
F1_Score : 0.724
Precision : 0.432
Recall : 0.593



XGBClassifier+preprocess

```
1 pipe4.fit(X=train['review'], y=train['sentiment'])
2 y_pred = pipe3.predict(test['review']).astype('int64')
3 y_test = test['sentiment'].values.astype('int64')
4 shen_bing_how_shy(y_pred, y_test)
```

[XGBClassifier+preprocess]
Accuracy : 0.622
F1_Score : 0.691
Precision : 0.486
Recall : 0.545

