

Resumen

The goal of this project was to predict sentiment for the given Twitter post using Python Data science packages. I implemented to data processing approaches; Kfold train-test split method and sklearn train-test-split method. I used 5 different types of data obtained from main data and different classification methods. Method score and accuracy is reported.

Contents

1 Used Python Libraries

2 Evaluating Function

3 Preprocessing

- 3.1 Loading Data
- 3.2 Data prepration

4 Input Data Generating

5 Results

1. Used Python Libraries

In my code I used the following modules:

```
1 import os
2 import nltk
3 import numpy as np
4 import re as regex
5 import pandas as pd
6 from collections import Counter
7 import matplotlib.pyplot as plt
8 from sklearn.naive_bayes import BernoulliNB
9 from sklearn.metrics import f1_score,
    precision_score, recall_score, accuracy_score
10 from sklearn.model_selection import train_test_split,
    cross_val_score, GridSearchCV,
    RandomizedSearchCV
11 from sklearn import metrics, svm, model_selection,
    feature_extraction, linear_model, tree, ensemble,
    neural_network
```

2. Evaluating Function

I developed the following code for evaluating classification methods over the data. This is based on the similarity if applying different classification methods in sklearn. In this function **dev** I used "Kfold" from sklearn model-selection module to split data into training and test and k is set to 10. In

the following I calculated accuracy score with different measures including: "metrics.accuracy-score", "precision-score", "recall-score", "accuracy-score" and "f1-score" in kfold case. At the end I returned the average of the score called kfold score. In this **dev** the train-test-split method from "model-selection" is also applied. I computed method score "cls.score".

```
1 def cl-ts-evaluation(cls, k, data, target, method=
    None, test_ratio=0.25):
2     metric_ac_score = []
3     precision = []
4     recall = []
5     accuracy = []
6     f1 = []
7     if method == 'kfold':
8         kfold = model_selection.KFold(n_splits=k)
9         for ind_train, ind_test in kfold.split(data):
10            cls.fit(data[ind_train], target[
    ind_train])
11            y_pred = cls.predict(data[ind_test]) #
12            metric_ac_score.append(metrics.
    accuracy_score(y_pred, target[ind_test]))
13            precision.append(precision_score(y_pred,
    target[ind_test], average=None, pos_label=None)
    )
14            recall.append(recall_score(y_pred,
    target[ind_test], average=None, pos_label=None))
15            accuracy.append(accuracy_score(y_pred,
    target[ind_test]))
16            f1.append(f1_score(y_pred, target[
    ind_test], average=None, pos_label=None))
17            return [np.mean(metric_ac_score), np.
    mean(precision), np.mean(recall), np.mean(
    accuracy), np.mean(f1)]
18
19     else:
20         for i in range(k):
21             x_train, x_test, y_train, y_test =
    model_selection.train_test_split(data, target,
    test_size=test_ratio)
22             cls.fit(x_train, y_train)
23             metric_ac_score.append(cls.score(x_test
    , y_test))
24
25     return np.mean(metric_ac_score)
```

3. Preprocessing

3.1. Loading Data

Loading data is simple: first we provide the address of the data file in the following lines:

```
1 fileAddress=os.path.join(os.getcwd(), '
    twitter_sentiment_corpus.csv')
2 print('Start to loading first data from this file: \
    n', fileAddress)
```

Then the .csv is red:

```
1 tweetsData = pd.read_csv(fileAddress)
```

Finally I used stemmer

```
1 stemmer=nltk.PorterStemmer()
2 df["TweetText"] = list(map(lambda str: stemmer.stem(
    str.lower()), df["TweetText"]))
```

3.2. Data prepration

In this section I first removed the improper data i.e. the data in which the emotion is not equal to "positive", "negative" or "neutral":

```
1 tweetsData = tweetsData.drop(tweetsData[tweetsData['
    Sentiment'] == 'irrelevant'].index)
```

The I tried to remove unavailable data, numbers and spaces by the following commands respectively:

```
1 df = df[df["TweetText"] != "Not Available"]
2 r_num = regex.compile(r"s?[0-9]+\.[0-9]*")
3 df.loc[:, "TweetText"].replace(r_num, "", inplace=
    True)
```

Also the punctuation is removed by the following command:

```
1 for remove in map(lambda r: regex.compile(regex.
    escape(r)), [",", ":", "\", "=", "&", ";", "%",
    "$", "@", "%", "^", "*", "(", ")", "[", "]",
    "|", "/", "\\", ">", "<", "-", "!", "?", ".",
    ",", "_", "—", "#"]):
2 df.loc[:, "TweetText"].replace(remove, "", inplace=
    True)
```

The distribution of the emotions is computed by counting number of each emotion in data:

```
1 negative = len(df[df["Sentiment"] == "negative"])
2 positive = len(df[df["Sentiment"] == "positive"])
3 neutral = len(df[df["Sentiment"] == "neutral"])
4 t_size = negative + positive + neutral
5 dist_negative = negative/t_size
6 dist_positive = positive/t_size
7 dist_neutral = neutral/t_size
```

In order to enumerate the target, in the following I changed the emotions into numbers; positive to 1, negative to -1, and neutral to 0:

```
1 row_index = tweetsData[tweetsData['Sentiment'] == '
    positive'].index
2 tweetsData.loc[row_index, 'Sentiment'] = 1
3 row_index = tweetsData[tweetsData['Sentiment'] == '
    neutral'].index
4 tweetsData.loc[row_index, 'Sentiment'] = 0
5 row_index = tweetsData[tweetsData['Sentiment'] == '
    negative'].index
6 tweetsData.loc[row_index, 'Sentiment'] = -1
```

Determining the data and target:

```
1 tweetData = np.array(tweetsData['TweetText'])
2 tweetTarget = np.array(tweetsData['Sentiment'].
    values, dtype="|S6")
```

constructing a stopwords by using corpus.stopwords.words from nltk package:

```
1 stopwords=nltk.corpus.stopwords.words("english")
```

4. Input Data Generating

In order to prepare input data for processing, which helps to increase the accuracy, I used vectorizatio method. One proper vectorization method in sklearn is TfidfVectorizer vectorized from feature-extraction.text. This process is done in the following lines and we get **Data-1**:

```
1 # TfidfVectorizer
2 tfidf=feature_extraction.text.TfidfVectorizer(use_idf
    =True, sublinear_tf=True, stop_words=stopwords) #
    'english'
3 data_1=tfidf.fit_transform(tweetData)
```

Data-2 is generated by using bigram TfidfVectorizer method; I set ngram-range=(1, 2) and also token pattern r'+' . Regular expression denoting what constitutes a "token", only used if analyzer == 'word'.

When building the vocabulary ignore terms that have a document frequency strictly lower than the threshold min-df=1.

```
1 tfidf2=feature_extraction.text.TfidfVectorizer(
    use_idf=True, sublinear_tf=True, ngram_range=(1,
    2), token_pattern=r'\b\w+\b', min_df=1,
    stop_words=stopwords) #'english'
2 data_2=tfidf2.fit_transform(tweetData)
```

Data-3 is generated by using CountVectorizer method which Convert a collection of text documents to a matrix of token counts. This implementation produces a sparse representation of the counts using scipy.sparse.csr-matrix. Prameters are left as default.

```
1 countVectorize = feature_extraction.text.
    CountVectorizer(stop_words=stopwords) #'english'
2 data_3=countVectorize.fit_transform(tweetData)
```

Data-4 is generated by using bigram CountVectorizer method in which the bigram parameters are considered: Token pattern and threshold is set as Data-2 case:

```
1 # bigram vectorizer
2 bigram_vectorizer = feature_extraction.text.
    CountVectorizer(ngram_range=(1, 2), min_df=1,
    stop_words=stopwords) #'english'
3 data_4 = bigram_vectorizer.fit_transform(tweetData)
```

Data-5 is generated by using bigram HashingVectorizer method. HashingVectorizer converts a collection of text documents to a matrix of token occurrences.

It turns a collection of text documents into a scipy.sparse matrix holding token occurrence counts.

This text vectorizer implementation uses the hashing trick to find the token string name to feature integer index mapping.

This strategy has several advantages:

- it is very low memory scalable to large datasets as there is no need to store a vocabulary dictionary in memory
- it is fast to pickle and un-pickle as it holds no state besides the constructor parameters
- it can be used in a streaming (partial fit) or parallel pipeline as there is no state computed during fit.

[Ref. from sklearn website]

```
1 # HashingVectorizer
2 hashingVectorizer = feature_extraction.text.
  HashingVectorizer(n_features=100)
3 data_5 = hashingVectorizer.fit_transform(tweetData)
```

The following part of the code computes classification and produces the results:

```
1 for meth in method:
2     i=0
3     for model in models:
4         for dat in data:
5             print("Model: {}".format(str(type(
6                 model).__name__, data_ind[i], meth))
7             print("-----")
8             print("Data: {}".format(data_ind[i]))
9             print("-----")
10            print("Method: {}".format(meth))
11            print("-----")
12            accuracy = cl_ts_evaluation(cls=model,
13                                       k=10, data=dat, target=tweetTarget, method=meth)
14            if meth == 'kfold':
15                print("----- Results")
16                print("")
17                print("metric_ac_score: " + str(
18                    accuracy[0]))
19                print("Precision: " + str(accuracy
20                    [1]))
21                print("Recall: " + str(accuracy[2]))
22                print("Accuracy: " + str(accuracy
23                    [3]))
24                print("F1: " + str(accuracy[4]))
25                print("=====")
26            else:
27                print('Accuracy of the method: \t\t'
28                    '{0:.4f}'.format(accuracy))
29                print("*****")
30            i+=1
```

5. Results

In this report I applied the following classification methods:

```
1 models = [linear_model.LogisticRegression(),
2           svm.LinearSVC(),
3           tree.DecisionTreeClassifier(max_depth=10),
4           ensemble.RandomForestClassifier(),
5           BernoulliNB()]
```

The following results are obtained in the case of using train-test-split method of sklearn: In the following I report three maximum scores that I obtained during computation:

```
1 1. Model: LinearSVC
2 -----
3 Data: data_1
4 -----
5 Accuracy of the method: 0.7669
6 *****
7 2. Model: LogisticRegression
8 -----
9 Data: data_4
10 -----
11 Accuracy of the method: 0.7614
12 *****
13 3. Model: LinearSVC
14 -----
15 Data: data_4
16 -----
17 Accuracy of the method: 0.7600
18 *****
19 4. Model: LinearSVC
20 -----
21 Data: data_2
22 -----
23 Accuracy of the method: 0.7595
24 *****
25 5. Model: LogisticRegression
26 -----
27 Data: data_3
28 -----
29 Accuracy of the method: 0.7592
30 *****
31 *****
```

In the following I report the results of train test splitting by Kfold. I used different approaches here which all are reporting the accuracy of the classification. First I want to describe the scores;

recall score: The recall is the ratio $\frac{tp}{tp+fn}$ where tp is the number of true positives and fn the number of false negatives. The recall is intuitively the ability of the classifier to find all the positive samples.

metrics.accuracy_score: Accuracy classification score. In multilabel classification, this function computes subset accuracy: the set of labels predicted for a sample must exactly match the corresponding set of labels in y-true.

metrics.accuracy_score: The precision is the ratio $\frac{tp}{(tp+fp)}$ where tp is the number of true positives and fp the number of false positives. The precision is intuitively the ability of the classifier not to label as positive a sample that is negative.

F1 score: Compute the F1 score, also known as balanced F-score or F-measure. The F1 score can be interpreted as a weighted average of the precision and recall, where an F1 score reaches its best value at 1 and worst score at 0. The relative contribution of precision and recall to the F1 score are equal. The formula for the F1 score is:

$$F1 = 2 * \frac{(\text{precision} * \text{recall})}{(\text{precision} + \text{recall})}$$

In the multi-class and multi-label case, this is the weighted average of the F1 score of each class.

```
2
3 =====
4 Model: LogisticRegression==>
5 -----
6 Data: data_4
7 -----
8 Method: kfold
9 -----
10 Precision: 0.63963963964
11 =====
12 Model: LinearSVC==>
13 -----
14 Data: data_2
15 -----
16 Method: kfold
17 -----
18 Precision: 0.63888888889
19 =====
20 Model: LinearSVC==>
21 -----
22 Data: data_4
```

```
23 -----
24 Method: kfold
25 -----
26 Precision: 0.617486338798
27 =====
28 Model: RandomForestClassifier==>
29 -----
30 Data: data_3
31 -----
32 Method: kfold
33 -----
34 Precision: 0.604166666667
35 =====
36 Model: RandomForestClassifier==>
37 -----
38 Data: data_4
39 -----
40 Method: kfold
41 -----
42 Precision: 0.607142857143
```