

Report for Project1

MENG Qian 20466688

It is a classification problem in this project. The data set is given, which is consist of 3 csv files. We need to train a model on the training set and use the test data to predict the label of them, and then through the accuracy of the result to measure the model's quality.

Obviously, we do not know which classifier is the best one which has the highest accuracy. To carry out this project, we need to choose some of classifier at first that suitable for solving such problems. After experiments, we choose the model which the accuracy is highest as our model.

Data description

The data set consists of the training data, the labeled training data and the test data. Through command lines such as ‘head’, ‘wc’, we know that train data is 3220*57; labeled train data is 3220*57; the test data is 1380*57.

The first 10 data points of training data (incomplete),

[illegible]

The first 10 data points of labeled training data,

0.0
1.0
0.0
0.0
0.0
0.0
0.0
1.0
1.0
0.0

The first 10 data points of test data (incomplete),

[illegible]

Data import and processing

```
# import data
data_train = pd.read_csv(r'/Users/mengqian/Desktop/DLproject1/traindata.csv')
label_train = pd.read_csv(r'/Users/mengqian/Desktop/DLproject1/trainlabel.csv')
label_test = pd.read_csv(r'/Users/mengqian/Desktop/DLproject1/testdata.csv')

# data processing
train_scaled = preprocessing.scale(data_train)
test_scaled = preprocessing.scale(label_test)
train_scaled = pd.DataFrame(data=train_scaled)
label_train = np.ravel(label_train)
clf.fit(train_scaled, label_train)
```

Select models

For solving this problem, I chose three classifiers to label the test data.

```
models = [  
    # Naive Bayes classifier for multivariate Bernoulli models  
    (BernoulliNB(), "Naive_Bayes"),  
  
    # Classifier implementing the k-nearest neighbors vote  
    (KNeighborsClassifier(algorithm='kd_tree'), "k_nearest"),  
  
    # C-Support Vector Classification  
    (SVC(probability=True, kernel='rbf'), "svc")  
]
```

Evaluate models

To avoid over-fitting, using k-fold cross-validation method to evaluate the accuracy of models.

```
# use Cross Validation to estimate the accuracy of the model
scores = cross_val_score(clf, train_scaled, label_train, cv=5)
print(scores)
```

```
# print the result and save as a csv file
result = clf.predict(test_scaled)
res = map(lambda x: int(x), result)
file_name = "/Users/mengqian/Desktop/DLproject1/test_label_%s.csv" % model_name
np.savetxt(file_name, res, delimiter=",")
print(res)
```

```
[ 0.90232558  0.90993789  0.89440994  0.89269051  0.90357698]
[0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1,
[ 0.88062016  0.90062112  0.90372671  0.9066874  0.88491446]
[0, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1,
[ 0.9255814  0.93322981  0.93012422  0.9377916  0.91912908]
[0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1,
```

Conclusion

From the output, it can be included that the svm of C-Support Vector Classification method has the highest accuracy, so I choose it as my model in this project.