

Report of Project 1

He Yajing 20459312

November 12, 2017

1 Introduction

For this project, I think I can use many kinds of classifiers to classify the training data. But they have different training accuracies. Thus, I consider it is significant to identify the classifier with a very high accuracy. Therefore, the process of my project includes data reading, classifier selection, model training and label prediction.

2 Data Reading

2.1 Training Data

For training process, there are two parts of data, one is training data and the other is training label. I read training data line by line and convert them into a numpy array with normalization. Also, for training label, I read it line by line and convert into a numpy array. Instead of normalization, I change the label from string type into int type.

```
#read training data
train_data = []
with open("traindata.csv") as rd:
    data = rd.readlines()
    for line in data:
        train_data.append(line.split(","))
train_data = np.array(train_data)
train_data = preprocessing.normalize(train_data)
#print(train_data)

#read label data
label_data = []
with open("trainlabel.csv") as ld:
    label = ld.readlines()
label_data = np.array(label)
for i in range(len(label_data)):
    if label_data[i] == '0.0\n':
        label_data[i] = 0
    else:
        label_data[i] = 1
label_data = label_data.astype(np.int32)
# for i in range(len(label_data)):
#     print(label_data[i])
```

Figure 1: Read training data and training label.

2.2 Test Data

For test data reading, it is similar to training data. I read the file line by line at first and convert into the numpy array with normalization.

3 Classifier Selection

After data is ready, I will firstly use training data to test different models and their own accuracies. Therefore, I choose the SVM model, logistic regression model, bagging classifier, gra-

```

37
38 #read test data
39 test_data = []
40 with open("testdata.csv") as rd:
41     data = rd.readlines()
42     for line in data:
43         test_data.append(line.split(","))
44 test_data = np.array(test_data)
45 test_data = preprocessing.normalize(test_data)
46 #print(test_data)
47

```

Figure 2: Read test data.

dient boosting classifier, decision tree classifier, random forest classifier, MLP classifier, ensemble learning and so on. For some models, I also need to adjust different parameters to make the model highly accurate. Here comes the different accuracies for different models.

```
The accuracy of SVM model is 0.604658385093
The accuracy of LogisticRegression model is 0.68850931677
The accuracy of BaggingClassifier model is 0.90559062112
The accuracy of GBDT model is 0.97546838509
The accuracy of Decision Tree model is 0.999689440994
The accuracy of RandomForestClassifier model is 0.99847204969
The accuracy of ensemble voting model is 0.99847204969
The training accuracy of AdaClassifier model is 0.80291925456
killy@kali:~/BDT_00008$ python3 1
```

Figure 3: Accuracies of different classifiers

When I use ensemble model, it will take more than one classifiers to train the data and combine all their responses to generate the final results, which will improve the accuracy significantly. Therefore, I try to use ensemble model to combine the three classifiers with high accuracies. Among lots of possible combinations, I find that if I combine Logistic Regression, Random Forest and Decision Tree models and use ensemble voting model decide the final results, the accuracy will be extremely high. Also, it is important to adjust the weight of each classifier in ensemble model, so I try to put larger weight for more accurate classifier and higher accuracy turns out. Finally, I choose ensemble voting including Logistic Regression, Random Forest and Decision Tree models to do this project.

4 Model Training and Prediction

After I decide what model to use, I use normalized training dataset to train the model. For programming, there is a package called scikit-learn which include many classifiers and I can import needed classifiers as well as adjust parameters to get better performance. Finally, I get a model with 0.9997 accuracy. Then I use this model to predict test data by using predict function. This function will compute the probabilities of classification and put test data into a class with larger probability. Eventually, I put all predictions into a csv file to record them. The following is the code for classification model including weights for all classifiers.

```

# Fit the training accuracy of ensemble voting model is
rf1 = LogisticRegression(random_state=0)
rf2 = RandomForestClassifier(random_state=0)
voting = VotingClassifier(estimators=[(rf1, 'rf1'), (rf2, 'rf2'), (gbc, 'gbc')], voting='soft', weights=[1,1,2])
model = VotingClassifier(estimators=[(rf1, 'rf1'), (rf2, 'rf2'), (gbc, 'gbc')], voting='soft', weights=[1,1,2])
# Fit the model
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
acc = model.score(X_train, y_train, sample_weight=1)
# Print the training accuracy of ensemble voting model is
print(acc)

# Predict the test data
test = model.predict(X_test)
acc_test = test_acc(y_test, test)

# Print the test accuracy
print(acc_test)

# Predict the probability of test data
print(model.predict_proba(X_test))

```

Figure 4: Classification Model