

Author:Ruchan-Pi  
UID:500055496

Group 174

Update Time: 2022/4/7

Hardware: Intel(R) Core(TM) i7-8550U CPU @ 1.80GHz (8 CPUs), ~2.0GHz

Memory: 16384MB RAM, Available OS Memory: 16256MB RAM

# 1. input data

```
In [3]: import pandas as pd
import os
import numpy as np
plt.rcParams['figure.figsize'] = (10, 10)
plt.rcParams['display.max_columns'] = 10)

data_train_df = pd.read_csv('1./input/train/train.csv')
data_test_df = pd.read_csv('1./input/test/test_input.csv', index_col=0)

In [4]: # train.csv including feature and label using for training model.
data_train_df = pd.read_csv('1./input/train/train.csv')

# test_train.csv includes 5000 samples used for label prediction. Test samples do not have labels.
data_test_df = pd.read_csv('1./input/test/test_input.csv', index_col=0)

In [5]: data_train_df.head()

Out[5]:
```

| id | v1 | v2 | v3 | v4 | ... | v781 | v782 | v783 | v784 | label |
|----|----|----|----|----|-----|------|------|------|------|-------|
| 0  | 1  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 1     |
| 1  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0     |
| 2  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0     |
| 3  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 4     |
| 4  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 8     |

5 rows x 786 columns

```
In [6]: data_test_df.head()

Out[6]:
```

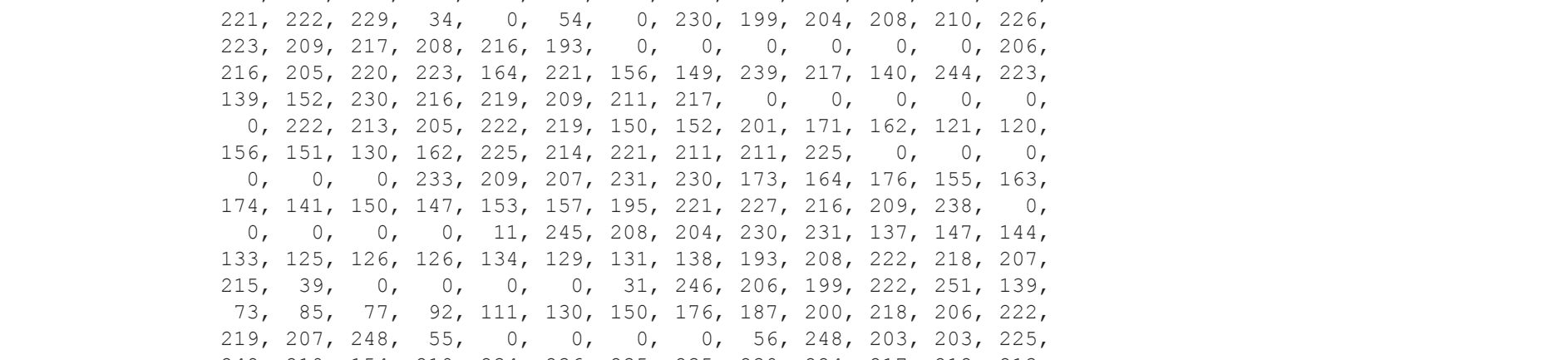
| v1 | v2 | v3 | v4 | v5 | ... | v780 | v781 | v782 | v783 | v784 |
|----|----|----|----|----|-----|------|------|------|------|------|
| 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 1  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 2  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 3  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 4  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |

5 rows x 784 columns

```
In [7]: # selecting input feature
data_train_feature = data_train_df.loc[:, "v1":"v784"].to_numpy()
# selecting test feature
data_test_feature = data_test_df.loc[:, "v1":"v784"].to_numpy()
# selecting output label
data_train_label = data_train_df['label'].to_numpy()
plt.show()
```

```
In [8]: import matplotlib.pyplot as plt
data_train_feature = data_train_feature.reshape((data_train_feature.shape[0], 28, 28))
plt.imshow(data_train_feature[0], cmap=plt.cm.gray)
plt.title('class = %s' % str(data_train_label[0])) # = Pullover
plt.show()

class 2 Pullover
```



0 5 10 15 20 25

```
In [9]: data_test_df.head()

Out[9]:
```

| v1 | v2 | v3 | v4 | v5 | ... | v780 | v781 | v782 | v783 | v784 |
|----|----|----|----|----|-----|------|------|------|------|------|
| 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 1  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 2  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 3  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 4  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |

5 rows x 784 columns

# 2. prediction

```
In [10]: # prediction on the test data
# normalization

def norm(data):
    # min = data.min(axis=0)
    # max = data.max(axis=0)
    # norm_result = (data-min)/(max-min)
    # return norm_result

train_norm = norm(data_train_feature)

In [11]: data_train_feature[0]

Out[11]:
```

| id | v1 | v2 | v3 | v4 | v5 | ... | v780 | v781 | v782 | v783 | v784 |
|----|----|----|----|----|----|-----|------|------|------|------|------|
| 0  | 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 1  | 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 2  | 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 3  | 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 4  | 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |

5 rows x 784 columns

```
In [12]: train_norm[0]

Out[12]:
```

| id | v1         | v2         | v3         | v4         | v5         | ... | v780       | v781       | v782       | v783       | v784       |
|----|------------|------------|------------|------------|------------|-----|------------|------------|------------|------------|------------|
| 0  | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | ... | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 |
| 1  | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | ... | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 |
| 2  | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | ... | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 |
| 3  | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | ... | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 |
| 4  | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | ... | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 |

5 rows x 784 columns

```
In [13]: # PCA (after normalization)
train_norm = train_norm.reshape((train_norm.shape[0], -1))

# from sklearn.decomposition import PCA
# pca = PCA(n_components = 200)
# train_norm_pca = pca.fit_transform(train_norm)
# test_norm_pca = pca.fit_transform(test_norm)
```

```
In [14]: # separate data into training and validation set (train:test=8:2)
from sklearn.model_selection import train_test_split
X_train, X_valis, y_train, y_valis = train_test_split(train_norm, data_train_label, test_size = 0.01, random_state=0)
```

```
In [15]: # KNN classifier after parameter tuning
from time import time
from sklearn.metrics import accuracy_score
from sklearn.neighbors import KNeighborsClassifier

# K-fold method
from sklearn.model_selection import StratifiedKFold

# set parameter
para_grid_knn = {'n_neighbors': (3, 4, 5, 6, 7, 8), 'p': (1, 2)}

# 10-fold method
kf = StratifiedKFold(n_splits=10, shuffle=True)

# score data
acc_knn = []
for p in para_grid_knn["p"]:
    acc_mean = []
    for n in para_grid_knn["n_neighbors"]:
        acc = []
        start = time()
        knn = KNeighborsClassifier(n_neighbors = n_neighbors, p = p)
        knn.fit(X_train, y_train)
        y_pred = knn.predict(X_valis)
        acc.append(accuracy_score(y_pred, y_valis))
        time_list.append(time()-start)
    print("n_neighbors = %d, p = %d, Accuracy = %.2f%%, running time = %.2fs" % (n_neighbors, p, acc_mean, time_knn.append(time_list)))

n_neighbors = 3, p = 1, Accuracy = 86.33%, running time = 1.96s
n_neighbors = 4, p = 1, Accuracy = 86.00%, running time = 1.92s
n_neighbors = 5, p = 1, Accuracy = 85.33%, running time = 1.92s
n_neighbors = 6, p = 1, Accuracy = 85.00%, running time = 1.93s
n_neighbors = 7, p = 1, Accuracy = 85.67%, running time = 1.93s
n_neighbors = 8, p = 1, Accuracy = 85.00%, running time = 1.92s
n_neighbors = 3, p = 2, Accuracy = 86.00%, running time = 0.26s
n_neighbors = 4, p = 2, Accuracy = 86.00%, running time = 0.25s
n_neighbors = 5, p = 2, Accuracy = 85.33%, running time = 0.22s
n_neighbors = 6, p = 2, Accuracy = 86.33%, running time = 0.24s
n_neighbors = 7, p = 2, Accuracy = 85.33%, running time = 0.26s
n_neighbors = 8, p = 2, Accuracy = 85.67%, running time = 0.24s
```

```
In [16]: plt.figure(figsize = (15,12))

plt.subplot(2,2,1)
plt.plot(para_grid_knn["n_neighbors"], acc_knn[0])
plt.plot(para_grid_knn["n_neighbors"], acc_knn[1])
plt.legend(["p = 1", "p = 2"])
plt.xticks(para_grid_knn["n_neighbors"])
plt.title("The accuracy of KNN")
plt.ylabel("Accuracy(%)")
plt.xlabel("n_neighbors")

plt.subplot(2,2,2)
plt.plot(para_grid_knn["n_neighbors"], time_knn[0])
plt.plot(para_grid_knn["n_neighbors"], time_knn[1])
plt.legend(["p = 1", "p = 2"])
plt.xticks(para_grid_knn["n_neighbors"])
plt.title("Running time of KNN")
plt.ylabel("Time(s)")
plt.xlabel("n_neighbors")

plt.show()
```



The accuracy of KNN

Running time of KNN

# 3.1 KNN classifier

```
In [17]: # KNN with 5 fold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import KFold

KNN_score=[]
KNN_time=[]
kf = KFold(n_splits=5)
knn_score_sum=0

for train_index, test_index in kf.split(X_train, y_train):
    this_train_x, this_train_y = X_train[train_index], y_train[train_index]
    this_test_x, this_test_y = X_train[test_index], y_train[test_index]
    SVMClassifier = KNeighborsClassifier(n_neighbors = 6, p = 1)
    SVMClassifier.fit(this_train_x, this_train_y)
    prediction = SVMClassifier.predict(this_test_x)
    score = accuracy_score(this_test_y, prediction)
    print(score)
    KNN_score.append(score)
    start=time.time()
    prediction = SVMClassifier.predict(test_norm)
    end=time.time()
    print('Finish classifying, it takes %s' % (end-start)+'s')
    KNN_time.append(end-start)
    knn_score_sum = knn_score_sum + score

0.8511784517845117
Finish classifying, it takes 156.861338470459s
0.8430974430974431
Finish classifying, it takes 181.5887253152341s
0.8577441077441077
Finish classifying, it takes 154.0821062980652s
0.847979797979798
Finish classifying, it takes 155.8311786842346s
0.8501681501681502
Finish classifying, it takes 175.1950855255127s
```

```
In [18]: knn_score_sum/5

Out[18]: 0.85003700336695
```

# 3.2 SVM classifier

```
In [19]: ##SVM with 5 fold
from sklearn.model_selection import KFold
from sklearn.metrics import accuracy_score
from sklearn.model_selection import KFold

SVM_score=[]
SVM_time=[]
kf = KFold(n_splits=5)
svm_score_sum=0

for train_index, test_index in kf.split(X_train, y_train):
    this_train_x, this_train_y = X_train[train_index], y_train[train_index]
    this_test_x, this_test_y = X_train[test_index], y_train[test_index]
    SVMClassifier = svm.LinearSVC(max_iter=5000, dual=True)
    SVMClassifier.fit(this_train_x, this_train_y)
    prediction = SVMClassifier.predict(this_test_x)
    score = accuracy_score(this_test_y, prediction)
    print(score)
    SVM_score.append(score)
    start=time.time()
    prediction = SVMClassifier.predict(test_norm)
    end=time.time()
    print('Finish classifying, it takes %s' % (end-start)+'s')
    SVM_time.append(end-start)
    svm_score_sum = svm_score_sum + score

0.8388988888888889
It takes 0.009974241256713867s
Finish classifying, it takes 0.0094778375024414s
0.8493750244149375
Finish classifying, it takes 0.0080382823940918s
0.8373737373737373
It takes 0.00703954696532734s
0.8424375243752438
Finish classifying, it takes 0.00900062942504803s
```

```
In [20]: svm_score_sum/5

Out[20]: 0.8354208754208754
```

# 3.3 logistic regression

```
In [21]: # Logistic Regression
from sklearn.linear_model import LogisticRegression

# set parameter
para_grid_lr = {'max_iter': (100, 200, 300), 'penalty': ('l1', 'l2')}

# score data
time_lr = []
acc_lr = []

for penalty in para_grid_lr["penalty"]:
    acc_mean = []
    time_list = []
    for max_iter in para_grid_lr["max_iter"]:
        acc = []
        start = time()
        lr = LogisticRegression(max_iter = max_iter, penalty = penalty, C = 1, solver = 'liblinear')
        lr.fit(X_train, y_train)
        y_pred = lr.predict(X_valis)
        acc.append(accuracy_score(y_pred, y_valis))
        time_list.append(time()-start)
    print("max_iter = %d, penalty = %s, Accuracy = %.2f%%, running time = %.2fs" % (max_iter, penalty, acc_mean, time_lr.append(time_list)))
    acc_lr.append(acc_mean)

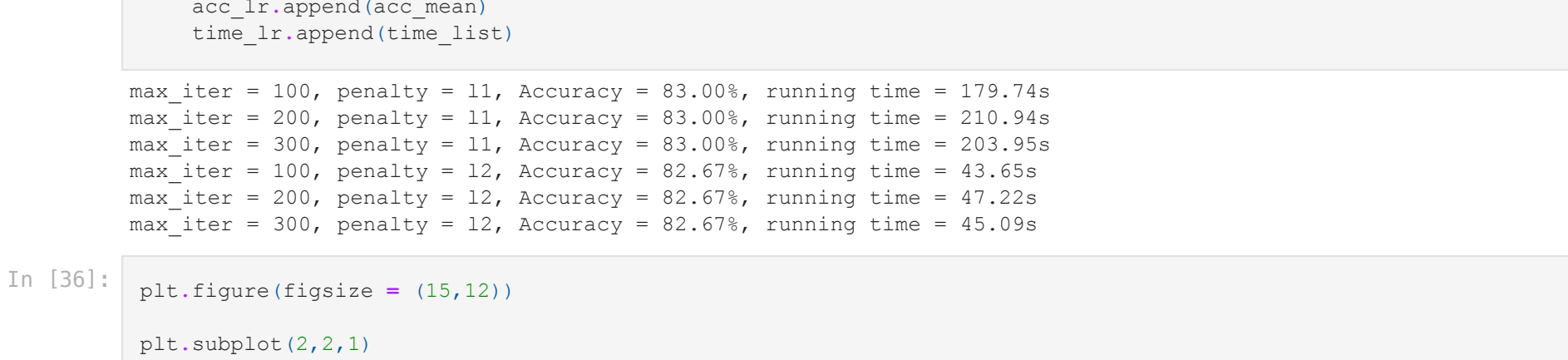
max_iter = 100, penalty = l1, Accuracy = 83.00%, running time = 179.74s
max_iter = 200, penalty = l1, Accuracy = 82.67%, running time = 210.94s
max_iter = 300, penalty = l1, Accuracy = 83.00%, running time = 203.95s
max_iter = 100, penalty = l2, Accuracy = 82.67%, running time = 43.65s
max_iter = 200, penalty = l2, Accuracy = 82.67%, running time = 47.22s
max_iter = 300, penalty = l2, Accuracy = 82.67%, running time = 45.09s
```

```
In [22]: plt.figure(figsize = (15,12))

plt.subplot(2,2,1)
plt.plot(para_grid_lr["max_iter"], acc_lr[0])
plt.plot(para_grid_lr["max_iter"], acc_lr[1])
plt.legend(["Penalty = l1", "Penalty = l2"])
plt.xticks(para_grid_lr["max_iter"])
plt.title("The accuracy of LR")
plt.ylabel("Accuracy(%)")
plt.xlabel("max_iter")

plt.subplot(2,2,2)
plt.plot(para_grid_lr["max_iter"], time_lr[0])
plt.plot(para_grid_lr["max_iter"], time_lr[1])
plt.legend(["Penalty = l1", "Penalty = l2"])
plt.xticks(para_grid_lr["max_iter"])
plt.title("The running time of LR")
plt.ylabel("Time(s)")
plt.xlabel("max_iter")

plt.show()
```



The accuracy of LR

The running time of LR

```
In [23]: ##Random Forest(RF) with 5 fold
from sklearn.ensemble import RandomForestClassifier

RF_score=[]
RF_time=[]
kf = KFold(n_splits=5)
rf_score_sum=0

for train_index, test_index in kf.split(X_train, y_train):
    this_train_x, this_train_y = X_train[train_index], y_train[train_index]
    this_test_x, this_test_y = X_train[test_index], y_train[test_index]
    RFClassifier = RandomForestClassifier(max_depth=11, min_samples_split=2, n_estimators=500, random_state=0)
    RFClassifier.fit(this_train_x, this_train_y)
    prediction = RFClassifier.predict(this_test_x)
    score = accuracy_score(this_test_y, prediction)
    print(score)
    RF_score.append(score)
    start=time.time()
    prediction = RFClassifier.predict(test_norm)
    end=time.time()
    print('Finish classifying, it takes %s' % (end-start)+'s')
    RF_time.append(end-start)
    rf_score_sum = rf_score_sum + score

0.8606060606060606
Finish classifying, it takes 0.252971172352764s
0.8653196531965319
Finish classifying, it takes 0.599959515167263s
0.8589225589225589
Finish classifying, it takes 0.4478431224822998s
0.8533720533720534
Finish classifying, it takes 0.604004146689516s
```

```
In [24]: rf_score_sum/5

Out[24]: 0.8578114781147812
```

# 3.4 Random forest

```
In [25]: # Random Forest ensemble
from sklearn.ensemble import RandomForestClassifier

# set parameter
para_grid_rf = {'max_depth': [7, 9, 11], 'random_state': [0, 12]}

# score data
time_rf = []
acc_rf = []

for max_depth in para_grid_rf["max_depth"]:
    acc_mean = []
    time_list = []
    for random_state in para_grid_rf["random_state"]:
        acc = []
        start = time()
        rf = RandomForestClassifier(max_depth = max_depth, random_state = random_state, n_estimators=500)
        rf.fit(X_train, y_train)
        y_pred = rf.predict(X_valis)
        acc.append(accuracy_score(y_pred, y_valis))
        time_list.append(time()-start)
    print("max_depth = %d, random_state = %d, Accuracy = %.2f%%, running time = %.2fs" % (max_depth, random_state, acc_mean, time_rf.append(time_list)))
    acc_rf.append(acc_mean)

max_depth = 7, random_state = 0, Accuracy = 82.67%, running time = 121.19s
max_depth = 9, random_state = 0, Accuracy = 82.67%, running time = 121.19s
max_depth = 11, random_state = 0, Accuracy = 82.67%, running time = 184.65s
max_depth = 7, random_state = 2, Accuracy = 82.67%, running time = 123.55s
max_depth = 9, random_state = 2, Accuracy = 82.67%, running time = 123.55s
max_depth = 11, random_state = 2, Accuracy = 82.67%, running time = 166.82s
```

```
In [26]: plt.figure(figsize = (15,12))

plt.subplot(2,2,1)
plt.plot(para_grid_rf["max_depth"], acc_rf[0])
plt.plot(para_grid_rf["max_depth"], acc_rf[1])
plt.legend(["random_state = 0", "random_state = 2"])
plt.xticks(para_grid_rf["max_depth"])
plt.title("The accuracy of RF")
plt.ylabel("Accuracy(%)")
plt.xlabel("max_depth")

plt.subplot(2,2,2)
plt.plot(para_grid_rf["max_depth"], time_rf[0])
plt.plot(para_grid_rf["max_depth"], time_rf[1])
plt.legend(["random_state = 0", "random_state = 2"])
plt.xticks(para_grid_rf["max_depth"])
plt.title("The running time of RF")
plt.ylabel("Time(s)")
plt.xlabel("max_depth")

plt.show()
```



The accuracy of RF

Running time of RF

```
In [27]: ##Random Forest(RF) with 5 fold
from sklearn.ensemble import RandomForestClassifier

RF_score=[]
RF_time=[]
kf = KFold(n_splits=5)
rf_score_sum=0

for train_index, test_index in kf.split(X_train, y_train):
    this_train_x, this_train_y = X_train[train_index], y_train[train_index]
    this_test_x, this_test_y = X_train[test_index], y_train[test_index]
    RFClassifier = RandomForestClassifier(max_depth=11, min_samples_split=2, n_estimators=500, random_state=0)
    RFClassifier.fit(this_train_x, this_train_y)
    prediction = RFClassifier.predict(this_test_x)
    score = accuracy_score(this_test_y, prediction)
    print(score)
    RF_score.append(score)
    start=time.time()
    prediction = RFClassifier.predict(test_norm)
    end=time.time()
    print('Finish classifying, it takes %s' % (end-start)+'s')
    RF_time.append(end-start)
    rf_score_sum = rf_score_sum + score

0.8606060606060606
Finish classifying, it takes 0.252971172352764s
0.8653196531965319
Finish classifying, it takes 0.599959515167263s
0.8589225589225589
Finish classifying, it takes 0.4478431224822998s
0.8533720533720534
Finish classifying, it takes 0.604004146689516s
```

```
In [28]: rf_score_sum/5

Out[28]: 0.8578114781147812
```

# 4. output

```
In [29]: # test_input.csv includes 5000 samples used for label prediction. Test samples do not have labels.
data_test_df = pd.read_csv('1./input/test/test_input.csv')
```



```
[27]: data_test_df.head()
```

|   | id | v1 | v2 | v3 | v4 | ... | v780 | v781 | v782 | v783 | v784 |
|---|----|----|----|----|----|-----|------|------|------|------|------|
| 0 | 0  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 1 | 1  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 2 | 2  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 3 | 3  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |
| 4 | 4  | 0  | 0  | 0  | 0  | ... | 0    | 0    | 0    | 0    | 0    |

5 rows x 785 columns

```
In [28]: y_pred_best = RFClassifier.predict(test_norm)
y_pred_best
```

```
Out[28]: array([1, 1, ..., 5, 0, 2], dtype=int64)
```

```
In [29]: # assume output is the predicted labels from classifiers using input as data from test_input.csv
# (5000,)

output_df = pd.DataFrame(y_pred_best, columns = ["label"])
output_df.to_csv('./Output/test_output.csv', sep=",", float_format='%d', index_label="id")
```