

信用卡欺诈判别

本文以kaggle上的Credit Card Fraud Detection为学习对象，利用XGBoost、Logistic Regression、Random Forest方法对信用卡欺诈进行识别，并将这三种模型的结果通过投票的方式进行聚合得到最终的结果。

1. 加载库并读入数据

1.1 加载库

```
1  # -*- coding: utf-8 -*-
2  #!/usr/bin/env python
3
4  import pandas as pd
5  import numpy as np
6  import seaborn as sns
7  import matplotlib.pyplot as plt
8  %matplotlib inline
9  import itertools
10 import xgboost as xgb
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.ensemble import RandomForestClassifier, VotingClassifier
13 from sklearn.model_selection import StratifiedKFold
14 from sklearn.metrics import precision_score, recall_score, confusion_matrix, classification_report,
15   roc_curve, precision_recall_curve
16 from __future__ import division
17 import warnings
18
19 # 忽略告警
20 warnings.filterwarnings('ignore')
```

1.2 读入数据

```
1 data = pd.read_csv('./data/creditcard.csv')
```

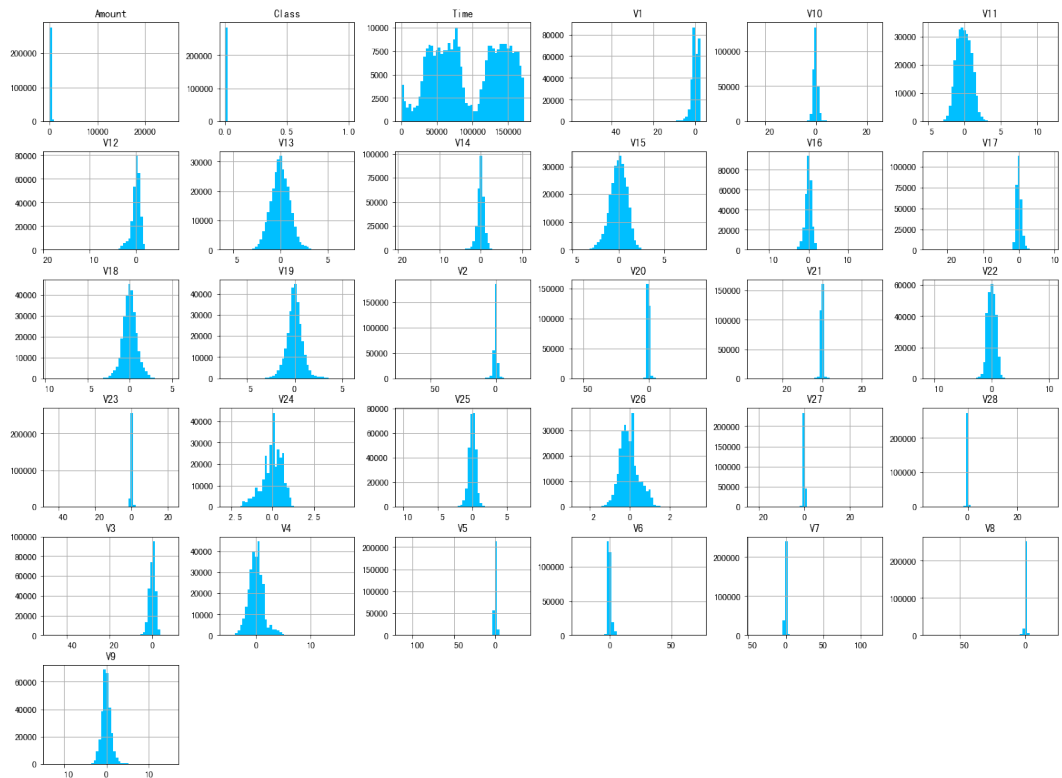
2. 数据探索与分析

数据集包含欧洲持卡人于2013年9月通过信用卡进行的交易。该数据集提供两天内发生的交易，其中在284,807笔交易中有492起欺诈行为。数据集非常不平衡，负面类别（欺诈）占有所有交易的0.172%。

数据经过PCA降维，特征V1,V2,...,V28是PCA获得的主要组件，还有交易时间Time、交易金额Amount。每条记录有一个Class，值为1时表示信用卡诈骗（正例），0表示正常消费（负例）。

2.1 变量分布和描述

```
1 data.hist (bins=50, figsize=(20,15), color = 'deepskyblue')
2
3 plt.show()
```



```
1 data.describe()
```

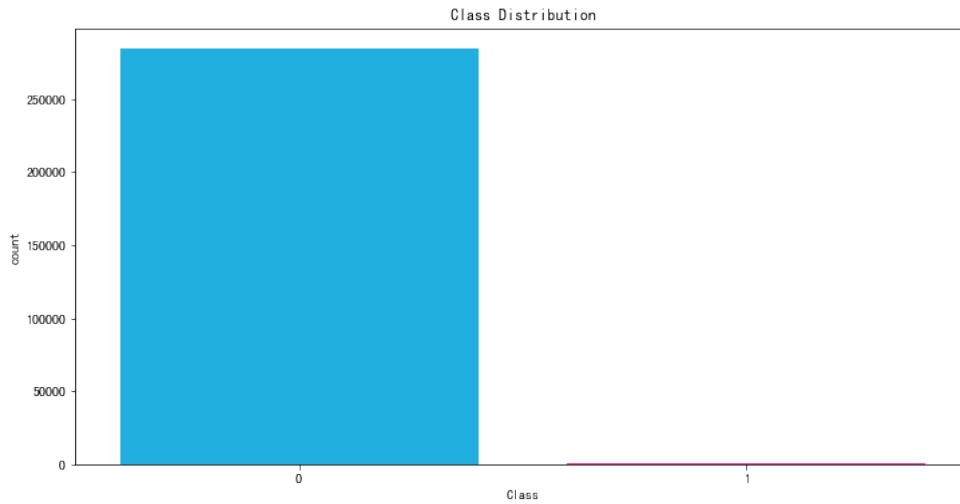
	Time	V1	V2	V3	V4	V5	V6	V7
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	94813.859575	1.165980e-15	3.416908e-16	-1.373150e-15	2.086869e-15	9.604066e-16	1.490107e-15	-5.556467e-16
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+00	1.237094e+00
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00	-1.137433e+02	-2.616051e+01	-4.355724e+01
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-01	-7.682956e-01	-5.540759e-01
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-02	-2.741871e-01	4.010308e-02
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-01	3.985649e-01	5.704361e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+01	7.330163e+01	1.205895e+02

8 rows × 9 columns

2.2 类别分布

数据集非常不平衡，正例（欺诈）数量占全部交易量的0.172%。

```
1 my_pal = {0: 'deepskyblue', 1: 'deeppink'}
2
3 plt.figure(figsize = (12, 6))
4 ax = sns.countplot(x = 'Class', data = data, palette = my_pal)
5 plt.title('Class Distribution')
6 plt.show()
7
8 # 正常交易量和欺诈交易量以及百分比
9 Count_Normal_transacation = len(data[data['Class']==0])
10 Count_Fraud_transacation = len(data[data['Class']==1])
11 Percentage_of_Normal_transacation =
12     float(Count_Normal_transacation)/(Count_Normal_transacation+Count_Fraud_transacation)
13 print('% of normal transacation      :{:0.3f}'.format(Percentage_of_Normal_transacation*100))
14 print('Number of normal transacation  :', Count_Normal_transacation)
15 Percentage_of_Fraud_transacation=
16     float(Count_Fraud_transacation)/(Count_Normal_transacation+Count_Fraud_transacation)
17 print('% of fraud transacation        :{:0.3f}'.format(Percentage_of_Fraud_transacation*100))
18 print('Number of fraud transacation   :', Count_Fraud_transacation)
```



```

1 % of normal transaction      :99.827
2 ('Number of normal transaction  :', 284315)
3 % of fraud transaction      :0.173
4 ('Number of fraud transaction   :', 492)

```

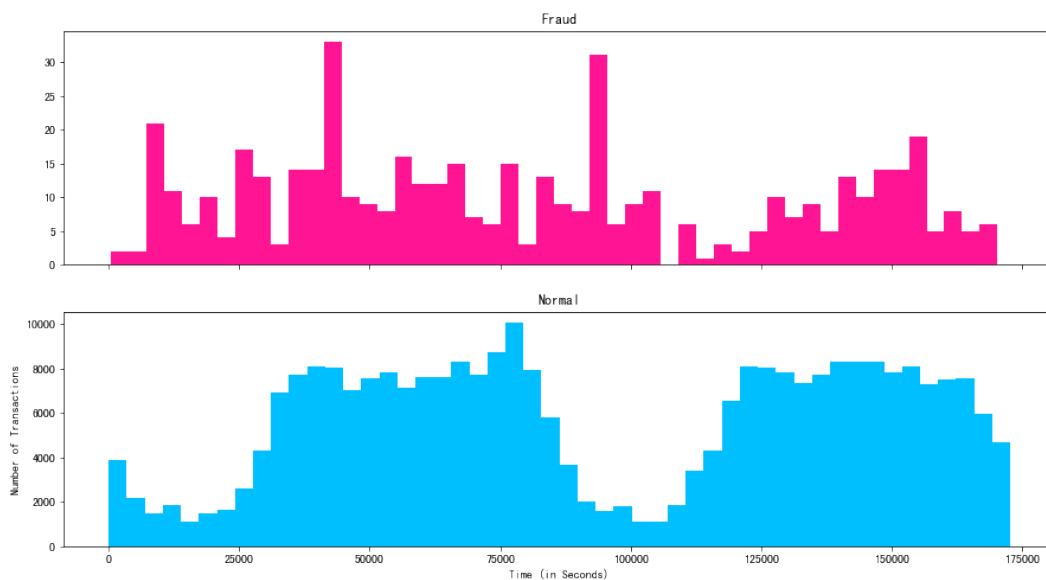
2.3 时间和类别的关系

从时间上看，欺诈在交易活跃和交易少的时段都有发生，在交易很活跃时一些欺诈事件不易被发觉，而交易少的时段都是夜间，此时人们已休息，也易于作案。深夜时正常交易量和欺诈量都较少，此时人们可能都处在睡眠时间。

```

1 f, (ax1, ax2) = plt.subplots(2, 1, sharex=True, figsize=(15,8))
2
3 bins = 50
4
5 ax1.hist(data.Time[data.Class == 1], bins = bins, color = 'deeppink')
6 ax1.set_title('Fraud')
7
8 ax2.hist(data.Time[data.Class == 0], bins = bins, color = 'deepskyblue')
9 ax2.set_title('Normal')
10
11 plt.xlabel('Time (in Seconds)')
12 plt.ylabel('Number of Transactions')
13 plt.show()

```



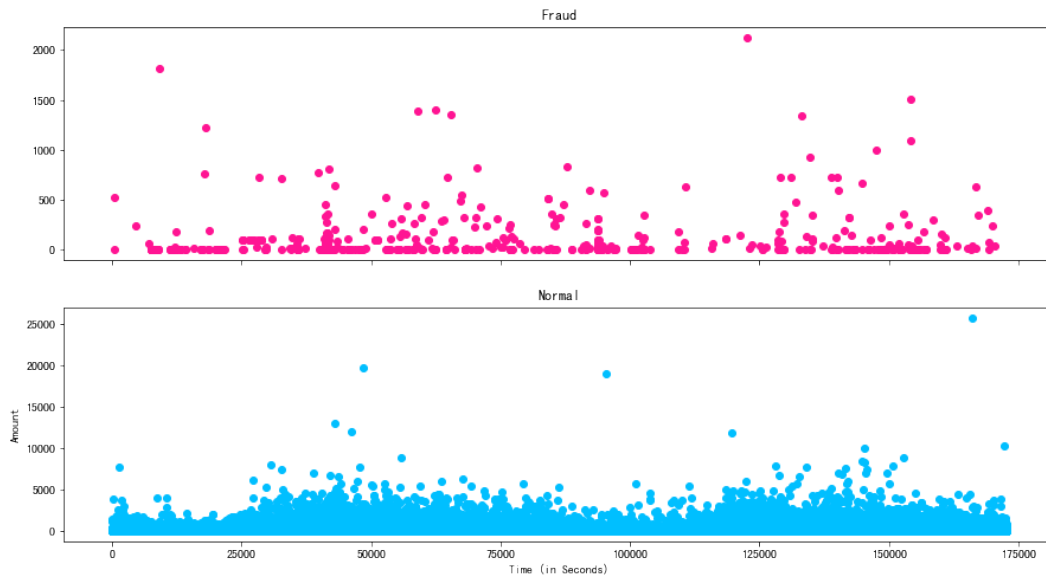
2.4 时间、交易金额和类别的关系

从交易金额来看，大多数欺诈事件的欺诈金额较小，但也有不少欺诈金额在500+以上。

```

1 f, (ax1, ax2) = plt.subplots(2, 1, sharex=True, figsize=(15,8))
2
3 ax1.scatter(data.Time[data.Class == 1], data.Amount[data.Class == 1], color = 'deeppink')
4 ax1.set_title('Fraud')
5
6 ax2.scatter(data.Time[data.Class == 0], data.Amount[data.Class == 0], color = 'deepskyblue')
7 ax2.set_title('Normal')
8
9 plt.xlabel('Time (in Seconds)')
10 plt.ylabel('Amount')
11 plt.show()

```

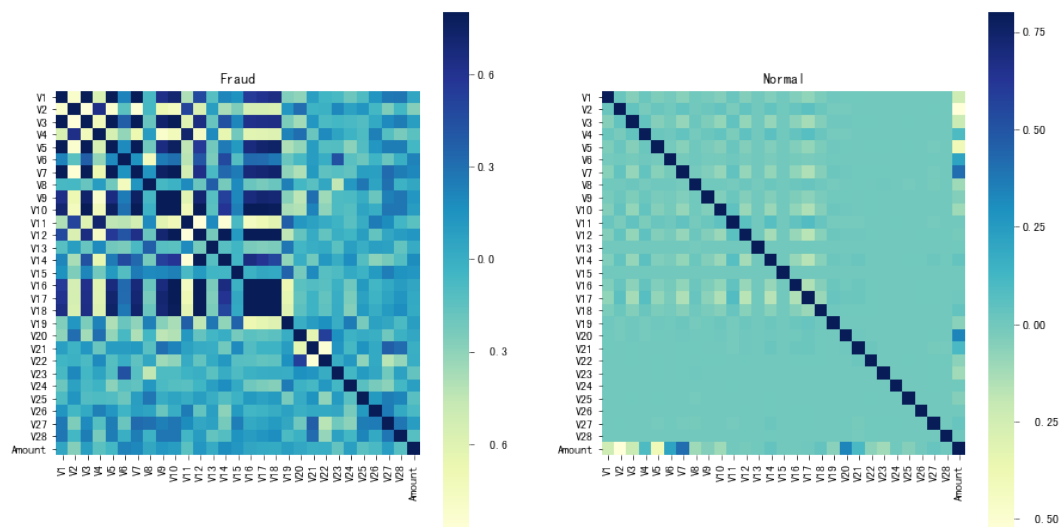


2.5 其它变量与类别的关系

```

1 f, (ax1, ax2) = plt.subplots(1,2,figsize=( 15, 8))
2
3 sns.heatmap(data.query('Class==1').drop(['Class','Time'],1).corr(), vmax = .8, square=True, ax = ax1,
4             cmap = 'YlGnBu')
5 ax1.set_title('Fraud')
6
7 sns.heatmap(data.query('Class==0').drop(['Class','Time'],1).corr(), vmax = .8, square=True, ax = ax2,
8             cmap = 'YlGnBu');
9 ax2.set_title('Normal')
10 plt.show()

```



3. 数据准备

3.1 划分训练集和测试集

训练集用来训练模型和参数调优，占90%。测试集用来评估模型性能，占10%。

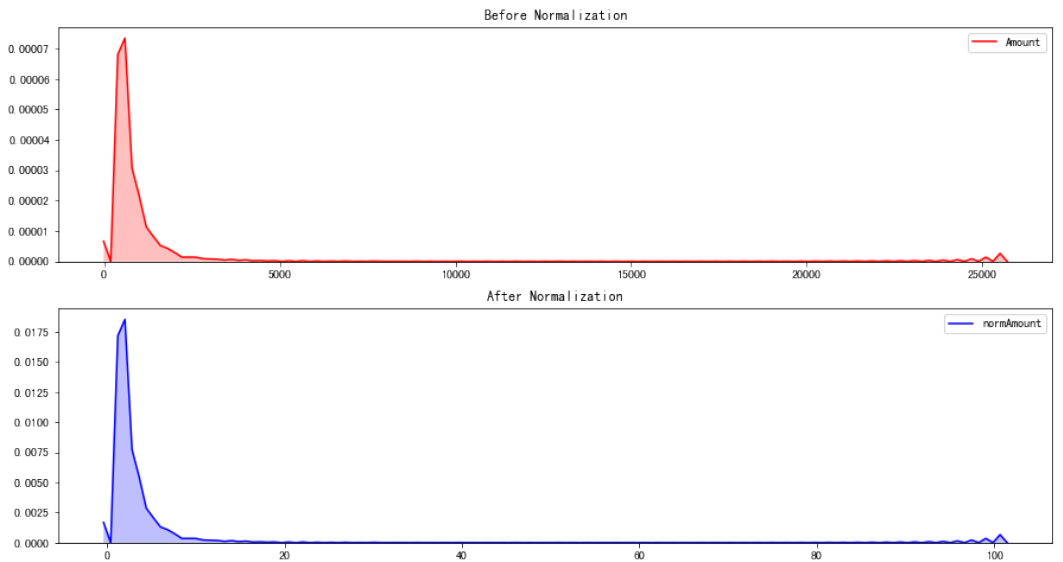
```
1 # 把数据集划分为训练集和测试集，一个用来train(训练) & validation(调参)，另一个是test (评估模型)
2 from sklearn.model_selection import train_test_split
3 train_set, test_set = train_test_split (data, test_size = 0.1, random_state = 42)
4 data = train_set
5 test_data = test_set
```

3.2 交易金额Amount的归一化

```
1 # Amount归一化
2 from sklearn.preprocessing import StandardScaler
3 data['normAmount'] = StandardScaler().fit_transform(data['Amount'].values.reshape(-1,1))
```

查看归一化后的效果：

```
1 f, (ax1, ax2) = plt.subplots(2,1,figsize=( 15, 8))
2
3 sns.kdeplot(data['Amount'],shade=True, ax = ax1, color='red')
4 ax1.set_title('Before Normalization')
5
6 sns.kdeplot(data['normAmount'],shade=True, ax = ax2, color='blue')
7 ax2.set_title('After Normalization')
8
9 plt.show()
```



3.3 删除无用变量

```
1 data = data.drop(['Amount', 'Time'],axis=1)
```

检查数据：

```
1 data.describe()
```

	V1	V2	V3	V4	V5	V6	V7	V8
count	256326.000000	256326.000000	256326.000000	256326.000000	256326.000000	256326.000000	256326.000000	256326.000000
mean	0.000787	0.000279	-0.000545	-0.000420	0.000217	0.000337	0.000506	0.000020
std	1.956494	1.648986	1.514006	1.415594	1.381680	1.333873	1.239372	1.192610
min	-56.407510	-72.715728	-48.325589	-5.683171	-113.743307	-26.160506	-43.557242	-73.216710
25%	-0.920249	-0.598404	-0.890516	-0.848857	-0.690771	-0.768232	-0.552928	-0.208079
50%	0.018564	0.065501	0.179232	-0.019309	-0.054395	-0.274532	0.040901	0.022482
75%	1.315616	0.804021	1.025290	0.743958	0.611975	0.398376	0.569886	0.327616
max	2.454930	22.057729	4.226108	16.875344	34.801666	73.301626	120.589494	20.007200

8 rows × 30 columns

4. 定义评价指标

4.1 混淆矩阵

```
1 # 混淆矩阵绘图
2 def plot_confusion_matrix(cm, classes,
3                             normalize = False,
4                             title = 'Confusion matrix',
5                             cmap = plt.cm.Blues) :
6     plt.imshow(cm, interpolation = 'nearest', cmap = cmap)
7     plt.title(title)
8     plt.colorbar()
9     tick_marks = np.arange(len(classes))
10    plt.xticks(tick_marks, classes, rotation = 0)
11    plt.yticks(tick_marks, classes)
12
13    thresh = cm.max() / 2.
14    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])) :
15        plt.text(j, i, cm[i, j],
16                 horizontalalignment = 'center',
17                 color = 'white' if cm[i, j] > thresh else 'black')
18
19    plt.tight_layout()
20    plt.ylabel('True label')
21    plt.xlabel('Predicted label')
22
```

4.2 Recall, Precision 和 F1_score

```
1 # 显示评估指标
2 def show_metrics():
3     tp = cm[1,1]
4     fn = cm[1,0]
5     fp = cm[0,1]
6     tn = cm[0,0]
7     print('Precision =      {:.3f}'.format(tp/(tp+fp)))
8     print('Recall    =      {:.3f}'.format(tp/(tp+fn)))
9     print('F1_score   =      {:.3f}'.format(2*((tp/(tp+fp))*(tp/(tp+fn)))/
10                                             ((tp/(tp+fp))+(tp/(tp+fn)))))

```

4.3 Precision-Recall曲线

```
1 # 绘制 P-R 曲线
2 def plot_precision_recall():
3     plt.step(recall, precision, color = 'b', alpha = 0.2,
4              where = 'post')
5     plt.fill_between(recall, precision, step = 'post', alpha = 0.2,
6                      color = 'b')
7
8     plt.plot(recall, precision, linewidth=2)
9     plt.xlim([0.0,1])
10    plt.ylim([0.0,1.05])
11    plt.xlabel('Recall')
12    plt.ylabel('Precision')
13    plt.title('Precision Recall Curve')
14    plt.show();

```

4.4 ROC 曲线

```
1 # 绘制 ROC 曲线
2 def plot_roc():
3     plt.plot(fpr, tpr, label = 'ROC curve', linewidth = 2)
4     plt.plot([0,1],[0,1], 'k--', linewidth = 2)
5     plt.xlim([0.0,0.001])
6     plt.ylim([0.0,1.05])
7     plt.xlabel('False Positive Rate')
8     plt.ylabel('True Positive Rate')
9     plt.title('ROC Curve')
10    plt.show();

```

4.5 特征重要性

```
1 predictors = ['V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10',
2              'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19',
3              'V20', 'V21', 'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28',
4              'Normamount']

```

绘制特征重要性的函数：

```
1 # 绘制特征重要性
2 def plot_feature_importance(model):
3     tmp = pd.DataFrame({'Feature': predictors, 'Feature importance': model.feature_importances_})
4     tmp = tmp.sort_values(by='Feature importance', ascending=False)
5     plt.figure(figsize = (15,8))
6     plt.title('Features importance', fontsize=14)
7     s = sns.barplot(x='Feature', y='Feature importance', data=tmp)
8     s.set_xticklabels(s.get_xticklabels(), rotation=90)
9     plt.show()
```

4.6 定义 (X,y) 及交叉验证

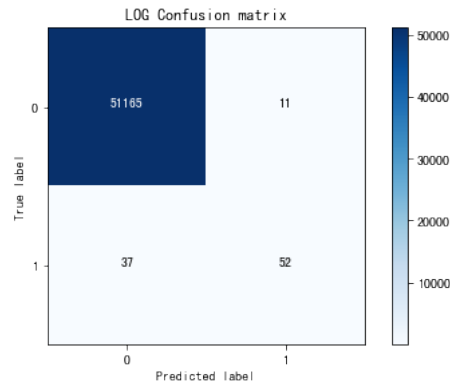
```
1 # 定义 x 和 y
2 y = np.array(data.Class.tolist())
3 data = data.drop('Class', 1)
4 X = np.array(data.as_matrix())
5
6 # K 折交叉验证
7 skf = StratifiedKFold(n_splits = 5, shuffle = True, random_state = 42)
8 for train_index, test_index in skf.split(X, y):
9     X_train, y_train = X[train_index], y[train_index]
10    X_test, y_test = X[test_index], y[test_index]
11
12 # 检查数据
13 print X.shape
14 print X_train.shape
15 print X_test.shape
```

```
1 (256326, 29)
2 (205061, 29)
3 (51265, 29)
```

5. Logistic Regression (LOG)

5.1 LOG - 未调优超参数

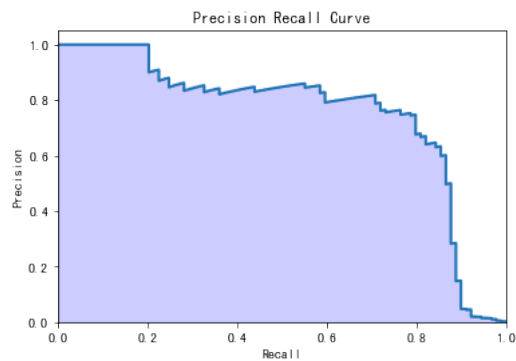
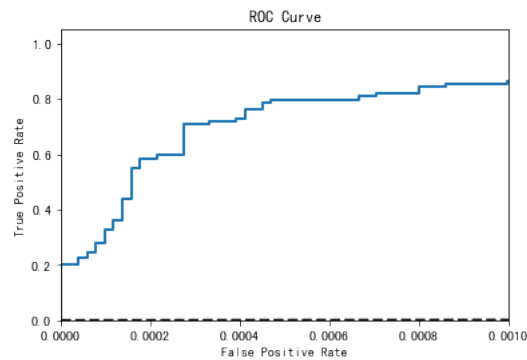
```
1 # 逻辑斯蒂回归
2 log_cfl = LogisticRegression()
3
4 log_cfl.fit(X_train, y_train)
5 y_pred = log_cfl.predict(X_test)
6 y_score = log_cfl.decision_function(X_test)
7
8 # 混淆矩阵 & 评估指标
9 cm = confusion_matrix(y_test, y_pred)
10 class_names = [0,1]
11 plt.figure()
12 plot_confusion_matrix(cm,
13                       classes = class_names,
14                       title = 'LOG Confusion matrix')
15 plt.show()
16
17 cm = cm.astype(np.float64)
18
19 show_metrics()
20
21 # ROC 曲线
22 fpr, tpr, t = roc_curve(y_test, y_score)
23 plot_roc()
24
25 # Precision-recall 曲线
26 precision, recall, thresholds = precision_recall_curve(y_test, y_score)
27 plot_precision_recall()
```



```

1 Precision = 0.825
2 Recall    = 0.584
3 F1_score  = 0.684

```



```

1 # 查看当前参数
2 from pprint import pprint
3 print('Parameters currently in use:\n')
4 pprint(log_cfl.get_params())

```

```

1 Parameters currently in use:
2
3 {'C': 1.0,
4  'class_weight': None,
5  'dual': False,
6  'fit_intercept': True,
7  'intercept_scaling': 1,
8  'max_iter': 100,
9  'multi_class': 'warn',
10 'n_jobs': None,
11 'penalty': 'l2',
12 'random_state': None,
13 'solver': 'warn',
14 'tol': 0.0001,
15 'verbose': 0,
16 'warm_start': False}

```


5.2 LOG - GridSearchCV 搜索超参数

```
1 # 使用GridSearchCV找到最佳的参数组合
2 from sklearn.model_selection import GridSearchCV
3 param_grid = {
4     'penalty' : ['l1','l2'],
5     'class_weight' : ['balanced', None],
6     'C' : [0.1, 1, 10, 100]
7 }
8
9 CV_log_cfl = GridSearchCV(estimator = log_cfl, param_grid = param_grid , scoring = 'recall', verbose
= 1, n_jobs = -1)
10 CV_log_cfl.fit(X_train, y_train)
11
12 best_parameters = CV_log_cfl.best_params_
13 print('The best parameters for using this model is', best_parameters)
```

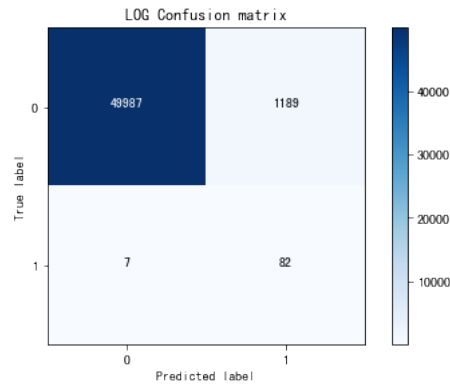
```
1 | Fitting 3 folds for each of 16 candidates, totalling 48 fits
```

```
1 [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
2 [Parallel(n_jobs=-1)]: Done 48 out of 48 | elapsed: 8.9min finished
```

```
1 ('The best parameters for using this model is', {'penalty': 'l1', 'C': 0.1, 'class_weight':
'balanced'})
```

5.3 LOG - 超参数更新

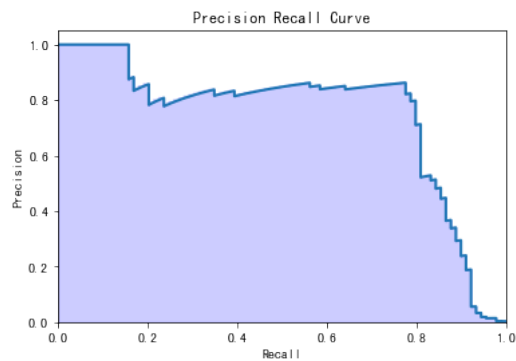
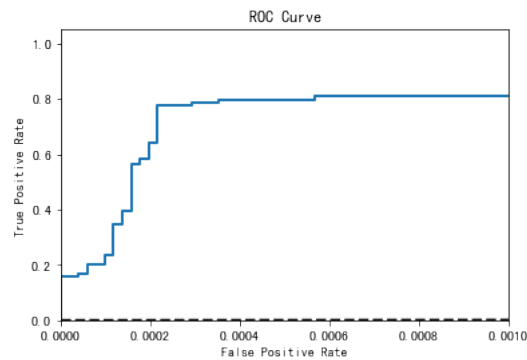
```
1 # 使用搜索的超参数组合重新构建模型, 将结果可视化
2 log_cfl = LogisticRegression(C = 0.1,
3                               penalty = 'l1',
4                               class_weight = 'balanced')
5 log_cfl.fit(X_train, y_train)
6 y_pred = log_cfl.predict(X_test)
7 y_score = log_cfl.decision_function(X_test)
8
9 # 混淆矩阵 & 评估指标
10 cm = confusion_matrix(y_test, y_pred)
11 class_names = [0,1]
12 plt.figure()
13 plot_confusion_matrix(cm,
14                       classes=class_names,
15                       title='LOG Confusion matrix')
16
17 plt.savefig('4.log_cfl_confusion_matrix.png')
18 plt.show()
19
20 cm = cm.astype(np.float64)
21
22 show_metrics()
23
24 # ROC 曲线
25 fpr, tpr, t = roc_curve(y_test, y_score)
26 plot_roc()
27
28 # Precision-recall 曲线
29 precision, recall, thresholds = precision_recall_curve(y_test, y_score)
30 plot_precision_recall()
31
32 fpr_log, tpr_log, t_log = fpr, tpr, t
33 precision_log, recall_log, thresholds_log = precision, recall, thresholds
```



```

1 Precision = 0.065
2 Recall    = 0.921
3 F1_score  = 0.121

```



6. Extreme Gradient Boosting (XGB)

6.1 XGB - 未调优超参数

```

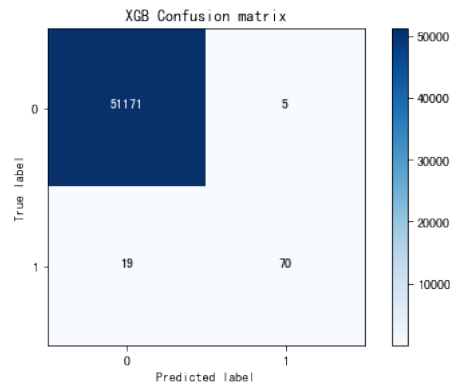
1 # xgb
2 xgb_cfl = xgb.XGBClassifier(n_jobs = -1)
3
4 xgb_cfl.fit(X_train, y_train)
5 y_pred = xgb_cfl.predict(X_test)
6 y_score = xgb_cfl.predict_proba(X_test)[: ,1]
7
8 # 混淆矩阵 & 评估指标
9 cm = confusion_matrix(y_test, y_pred)
10 class_names = [0,1]
11 plt.figure()
12 plot_confusion_matrix(cm,
13                       classes=class_names,
14                       title='XGB Confusion matrix')
15 plt.show()
16
17 cm = cm.astype(np.float64)
18
19 show_metrics()
20

```

```

21 # ROC 曲线
22 fpr, tpr, t = roc_curve(y_test, y_score)
23 plot_roc()
24
25 # Precision-recall 曲线
26 precision, recall, thresholds = precision_recall_curve(y_test, y_score)
27 plot_precision_recall()

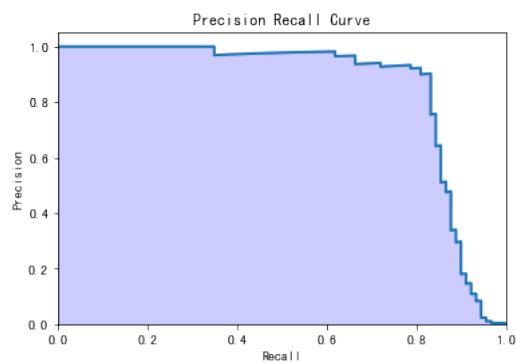
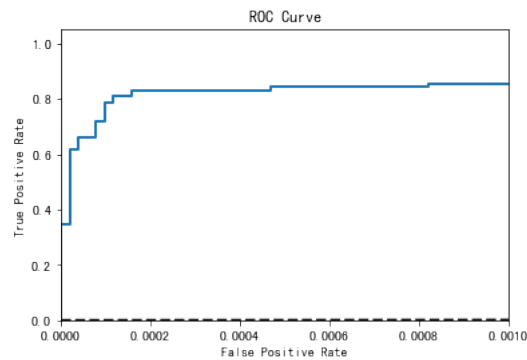
```



```

1 Precision = 0.933
2 Recall    = 0.787
3 F1_score  = 0.854

```



6.2 XGB - GridSearchCV 搜索超参数

```

1 # 使用 GridSearchCV 搜索 XGB 超参数
2 param_grid = {
3     'n_estimators': [100, 200, 300, 400]
4 }
5
6 CV_xgb_cfl = GridSearchCV(estimator = xgb_cfl, param_grid = param_grid, scoring = 'f1', verbose = 2,
7                             n_jobs = -1)
8 CV_xgb_cfl.fit(X_train, y_train)
9
10 best_parameters = CV_xgb_cfl.best_params_
11 print("The best parameters for using this model is", best_parameters)

```

```

1 Fitting 3 folds for each of 4 candidates, totalling 12 fits

```

```

1 | [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
2 | [Parallel(n_jobs=-1)]: Done 12 out of 12 | elapsed: 12.5min remaining: 0.0s
3 | [Parallel(n_jobs=-1)]: Done 12 out of 12 | elapsed: 12.5min finished

```

```

1 | ('The best parameters for using this model is', {'n_estimators': 200})

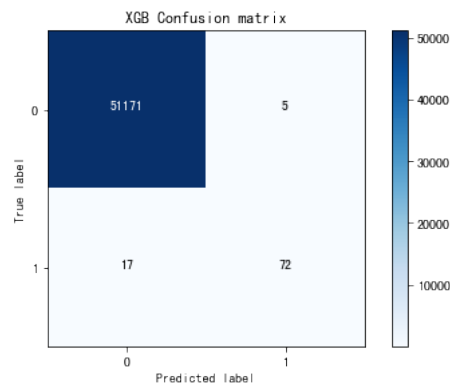
```

6.3 XGB - 超参数更新

```

1 | # xgb
2 | xgb_cfl = xgb.XGBClassifier(n_jobs = -1,
3 |                             n_estimators = 200)
4 |
5 | xgb_cfl.fit(X_train, y_train)
6 | y_pred = xgb_cfl.predict(X_test)
7 | y_score = xgb_cfl.predict_proba(X_test)[:,1]
8 |
9 | # 混淆矩阵 & 评估指标
10 | cm = confusion_matrix(y_test, y_pred)
11 | class_names = [0,1]
12 | plt.figure()
13 | plot_confusion_matrix(cm,
14 |                       classes = class_names,
15 |                       title = 'XGB Confusion matrix')
16 | plt.savefig('2.xgb_cfl_confusion_matrix.png')
17 | plt.show()
18 |
19 | cm = cm.astype(np.float64)
20 |
21 | show_metrics()
22 |
23 | # ROC 曲线
24 | fpr, tpr, t = roc_curve(y_test, y_score)
25 | plot_roc()
26 |
27 | # Precision-recall 曲线
28 | precision, recall, thresholds = precision_recall_curve(y_test, y_score)
29 | plot_precision_recall()
30 |
31 | fpr_xgb, tpr_xgb, t_xgb = fpr, tpr, t
32 | precision_xgb, recall_xgb, thresholds_xgb = precision, recall, thresholds

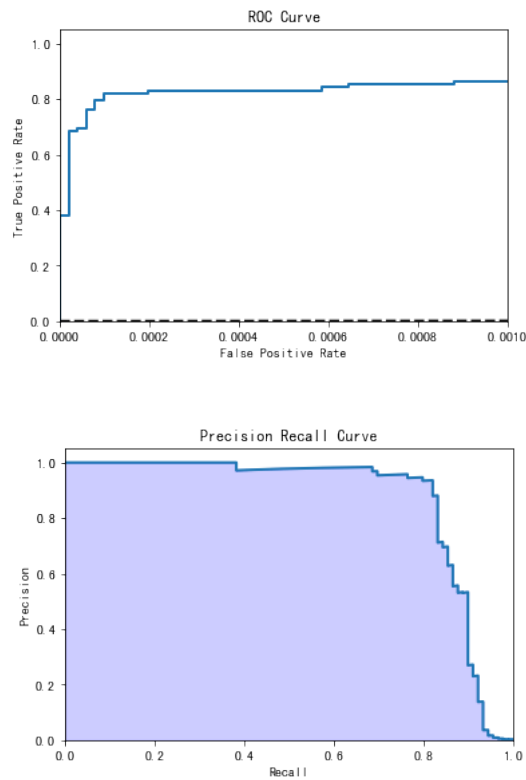
```



```

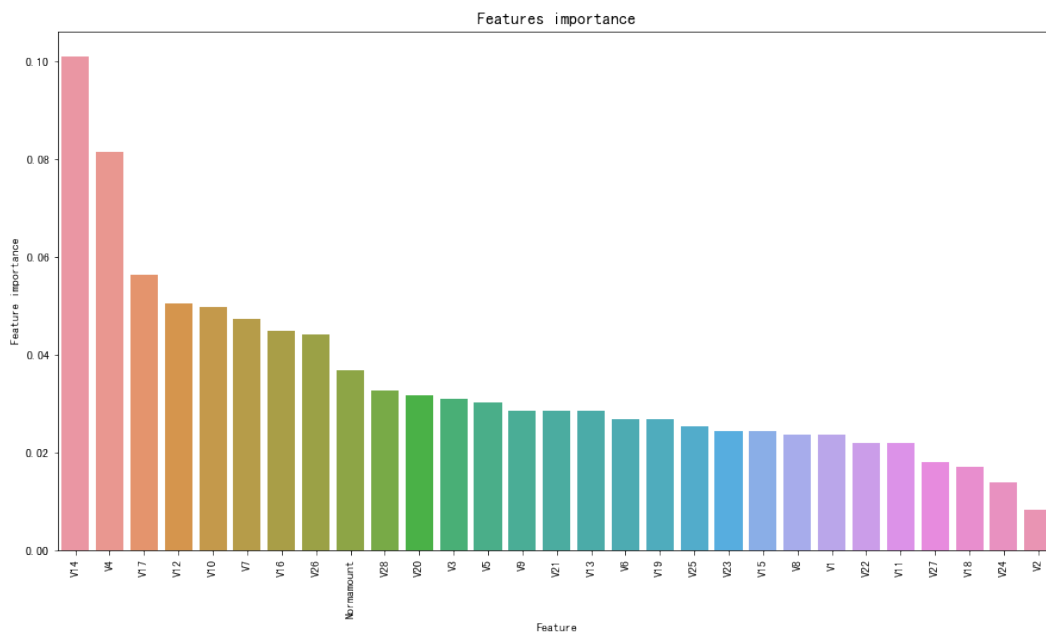
1 | Precision = 0.935
2 | Recall    = 0.809
3 | F1_score  = 0.867

```



对特征重要程度进行可视化：

```
1 plot_feature_importance(xgb_cfl)
```



7. Random Forest (RF)

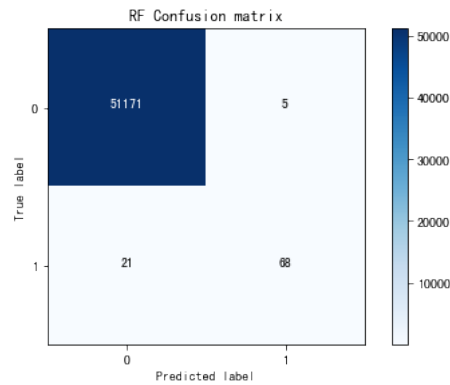
7.1 RF - 未调优超参数

```
1 # 随机森林
2 rf_cfl = RandomForestClassifier(n_jobs = -1,
3                               random_state = 42)
4
5 rf_cfl.fit(X_train, y_train)
6 y_pred = rf_cfl.predict(X_test)
7 y_score = rf_cfl.predict_proba(X_test)[: ,1]
8
9 # 混淆矩阵
10 cm = confusion_matrix(y_test, y_pred)
11 class_names = [0,1]
12 plt.figure()
```

```

13 plot_confusion_matrix(cm,
14                       classes = class_names,
15                       title = 'RF Confusion matrix')
16 plt.show()
17
18 cm = cm.astype(np.float64)
19
20 show_metrics()
21
22 # ROC 曲线
23 fpr, tpr, t = roc_curve(y_test, y_score)
24 plot_roc()
25
26 # Precision-recall 曲线
27 precision, recall, thresholds = precision_recall_curve(y_test, y_score)
28 plot_precision_recall()
29

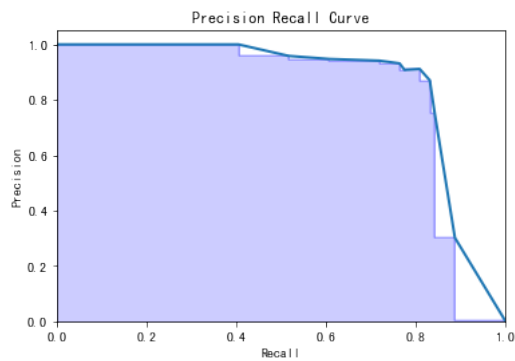
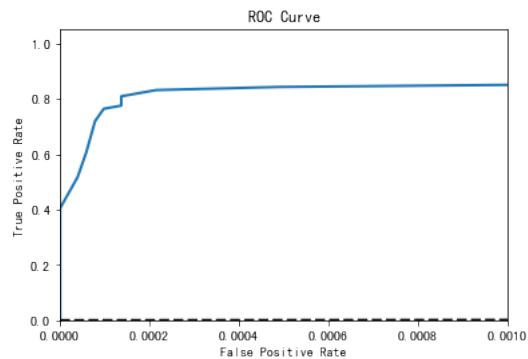
```



```

1 Precision = 0.932
2 Recall = 0.764
3 F1_score = 0.840

```



7.2 RF - GridSearchCV 搜索超参数

```

1 # 使用GridSearchCV 搜索 RF 超参数
2 from sklearn.model_selection import GridSearchCV
3
4 param_grid = {

```

```

5         'n_estimators': [100, 200, 500],
6         'max_features': [2, 3],
7         'min_samples_leaf': [1, 2, 4],
8         'min_samples_split': [2, 5, 10]
9     }
10
11 CV_rf_cfl = GridSearchCV(estimator = rf_cfl, param_grid = param_grid, scoring = 'f1', verbose = 10,
12                          n_jobs = -1)
13 CV_rf_cfl.fit(X_train, y_train)
14
15 best_parameters = CV_rf_cfl.best_params_
16 print("The best parameters for using this model is", best_parameters)

```

```

1 | Fitting 3 folds for each of 54 candidates, totalling 162 fits

```

```

1 [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
2 [Parallel(n_jobs=-1)]: Done   5 tasks      | elapsed:   3.7min
3 [Parallel(n_jobs=-1)]: Done  10 tasks      | elapsed:   8.0min
4 [Parallel(n_jobs=-1)]: Done  17 tasks      | elapsed:  13.9min
5 [Parallel(n_jobs=-1)]: Done  24 tasks      | elapsed:  19.3min
6 [Parallel(n_jobs=-1)]: Done  33 tasks      | elapsed:  26.4min
7 [Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed:  35.0min
8 [Parallel(n_jobs=-1)]: Done  53 tasks      | elapsed:  46.0min
9 [Parallel(n_jobs=-1)]: Done  64 tasks      | elapsed:  57.7min
10 [Parallel(n_jobs=-1)]: Done  77 tasks      | elapsed:  69.9min
11 [Parallel(n_jobs=-1)]: Done  90 tasks      | elapsed:  87.9min
12 [Parallel(n_jobs=-1)]: Done 105 tasks      | elapsed: 106.9min
13 [Parallel(n_jobs=-1)]: Done 120 tasks      | elapsed: 129.5min
14 [Parallel(n_jobs=-1)]: Done 137 tasks      | elapsed: 152.9min
15 [Parallel(n_jobs=-1)]: Done 154 tasks      | elapsed: 176.7min
16 [Parallel(n_jobs=-1)]: Done 162 out of 162 | elapsed: 187.3min finished

```

```

1 | ('The best parameters for using this model is', {'max_features': 3, 'min_samples_split': 2,
2 | 'n_estimators': 200, 'min_samples_leaf': 1})

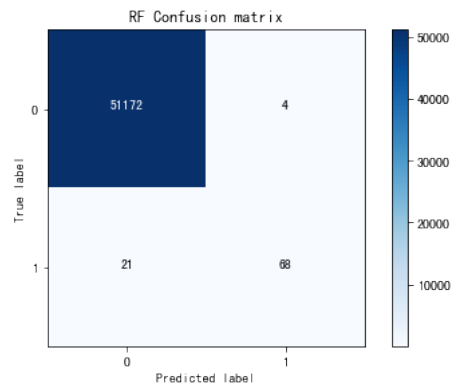
```

7.3 RF - 超参数更新

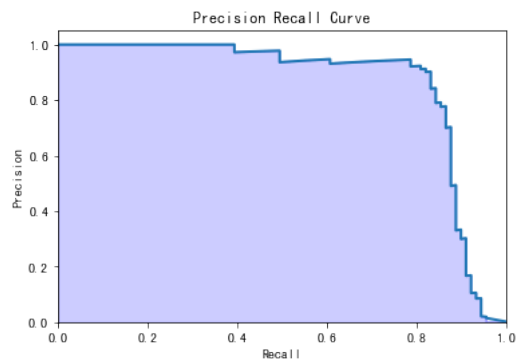
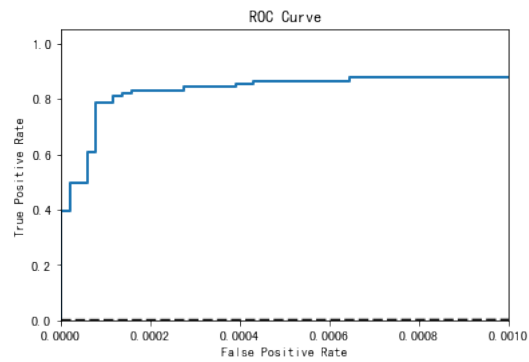
```

1 # 随机深林分类器
2 rf_cfl = RandomForestClassifier(n_estimators = 200,
3                                max_features = 3,
4                                min_samples_leaf = 1,
5                                min_samples_split = 2,
6                                n_jobs = -1,
7                                random_state = 42)
8
9 rf_cfl.fit(X_train, y_train)
10 y_pred = rf_cfl.predict(X_test)
11 y_score = rf_cfl.predict_proba(X_test)[:,:1]
12
13 # 混淆矩阵
14 cm = confusion_matrix(y_test, y_pred)
15 class_names = [0,1]
16 plt.figure()
17 plot_confusion_matrix(cm,
18                       classes = class_names,
19                       title = 'RF Confusion matrix')
20 plt.savefig('3.rf_cfl_confusion_matrix.png')
21 plt.show()
22
23 cm = cm.astype(np.float64)
24
25 show_metrics()
26
27 # ROC 曲线
28 fpr, tpr, t = roc_curve(y_test, y_score)
29 plot_roc()
30
31 # Precision-recall 曲线
32 precision, recall, thresholds = precision_recall_curve(y_test, y_score)
33 plot_precision_recall()
34
35 fpr_rf, tpr_rf, t_rf = fpr, tpr, t
36 precision_rf, recall_rf, thresholds_rf = precision, recall, thresholds

```

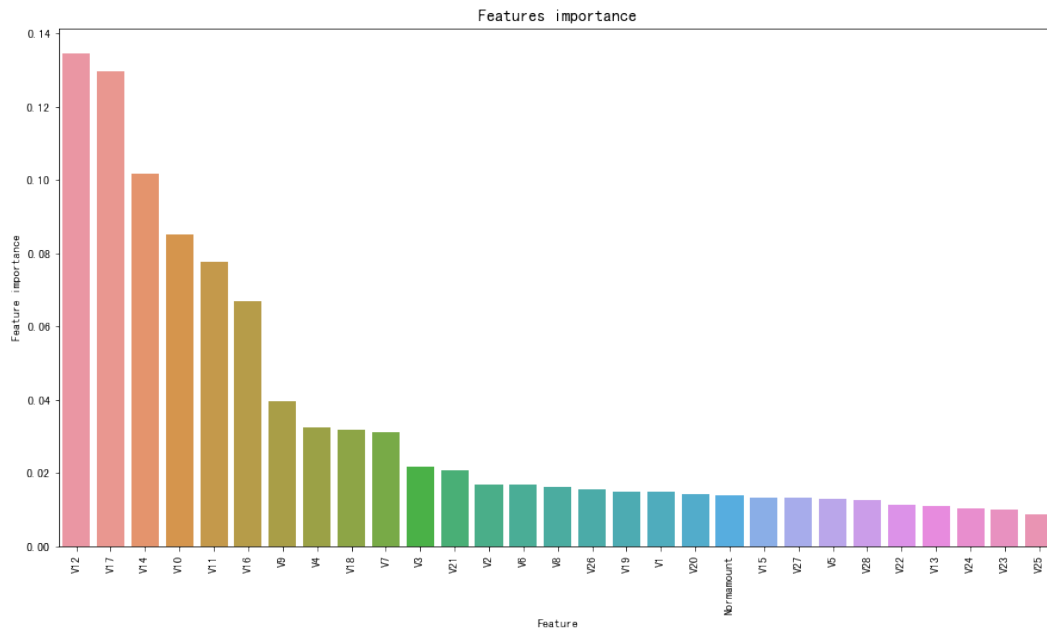


```
1 Precision = 0.944
2 Recall    = 0.764
3 F1_score  = 0.845
```



对特征重要程度可视化:

```
1 plot_feature_importance(rf_cfl)
```

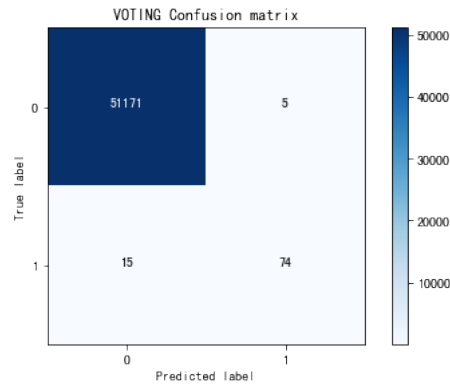
8. VotingClassifier = LOG - XGB - RF (F1=0.881)

8.1 VotingClassifier

```

1  #Voting Classifier
2  voting_cfl = VotingClassifier (
3      estimators = [('xgb', xgb_cfl), ('lt', log_cfl), ('rf', rf_cfl)],
4      voting='soft', weights = [1, 1, 1.33])
5
6  voting_cfl.fit(X_train,y_train)
7
8  y_pred = voting_cfl.predict(X_test)
9  y_score = voting_cfl.predict_proba(X_test)[:,:1]
10
11 # 混淆矩阵
12 cm = confusion_matrix(y_test, y_pred)
13 class_names = [0,1]
14 plt.figure()
15 plot_confusion_matrix(cm,
16                       classes = class_names,
17                       title = 'VOTING Confusion matrix')
18 plt.savefig('1.voting_confusion_matrix.png')
19 plt.show()
20
21 cm = cm.astype(np.float64)
22
23 show_metrics()
24
25 # ROC 曲线
26 fpr, tpr, t = roc_curve(y_test, y_score)
27 plot_roc()
28
29 # Precision-recall 曲线
30 precision, recall, thresholds = precision_recall_curve(y_test, y_score)
31 plot_precision_recall()
32
33 fpr_voting, tpr_voting, t_voting = fpr, tpr, t
34 precision_voting, recall_voting, thresholds_voting = precision, recall, thresholds
35

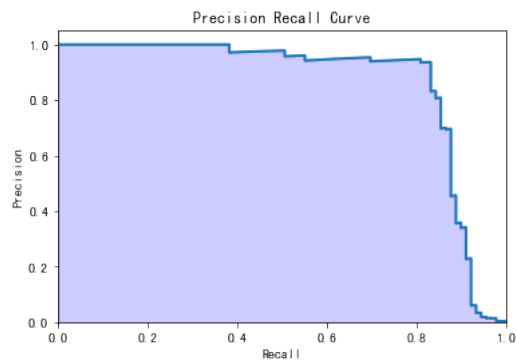
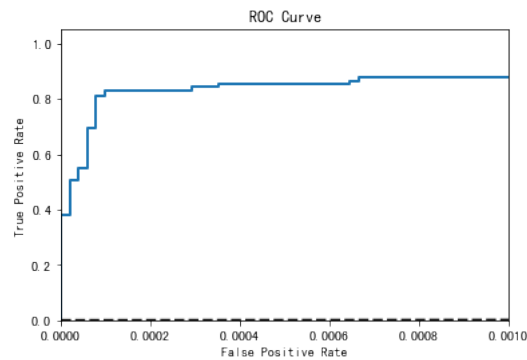
```



```

1 Precision = 0.937
2 Recall    = 0.831
3 F1_score  = 0.881

```



8.2 Precision - Recall - Threshold 曲线

```

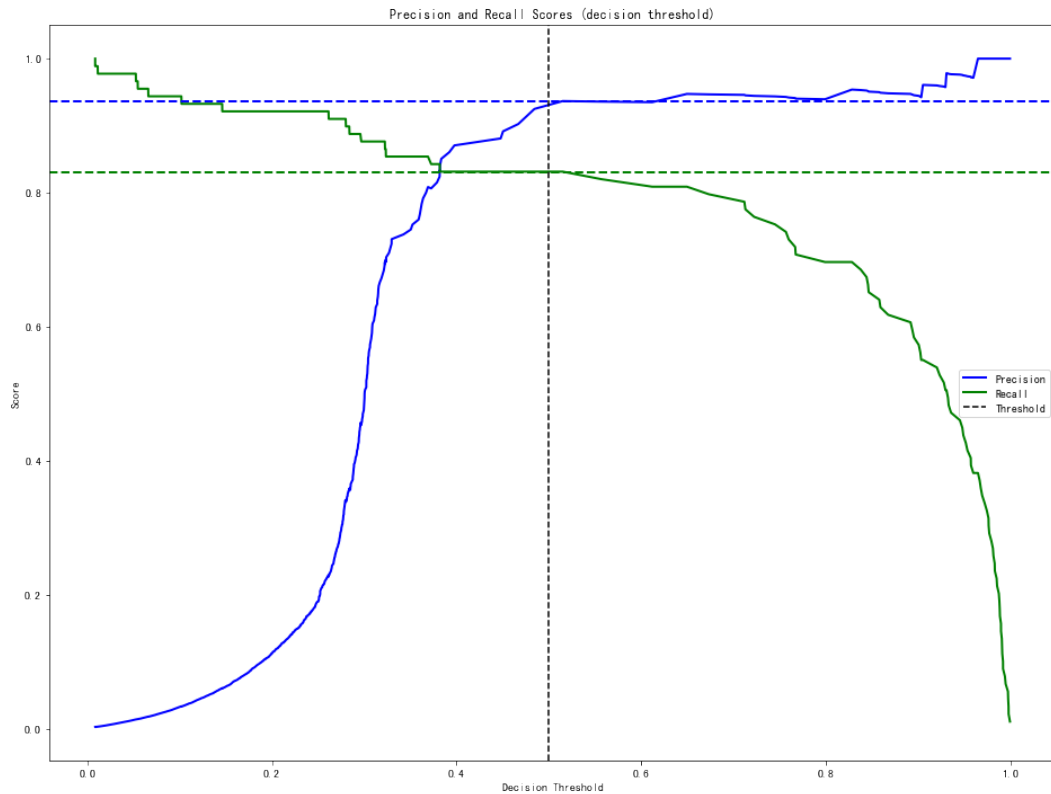
1 pr = 0.937
2 rec = 0.831
3 t = 0.5
4
5 # Precision-recall-threshold 曲线 :
6 def plot_precision_recall_vs_threshold(precisions, recalls, thresholds):
7     plt.figure(figsize=(16, 12))
8     plt.title('Precision and Recall Scores (decision threshold)')
9     plt.plot(thresholds, precisions[:-1], 'b-', linewidth=2, label='Precision')
10    plt.plot(thresholds, recalls[:-1], 'g', linewidth=2, label='Recall')
11    plt.axvline(t, color='k', linestyle='--', label='Threshold')
12    plt.axhline(pr, color='blue', linewidth=2, linestyle='--')
13    plt.axhline(rec, color='green', linewidth=2, linestyle='--')
14    plt.ylabel('Score')
15    plt.xlabel('Decision Threshold')
16    plt.legend(loc='best')
17    plt.savefig('5.prec_recc_threshold.png')
18    plt.show();

```

```

1 plot_precision_recall_vs_threshold(precision, recall, thresholds)

```



8.3 比较 ROC 曲线（所有模型）

```

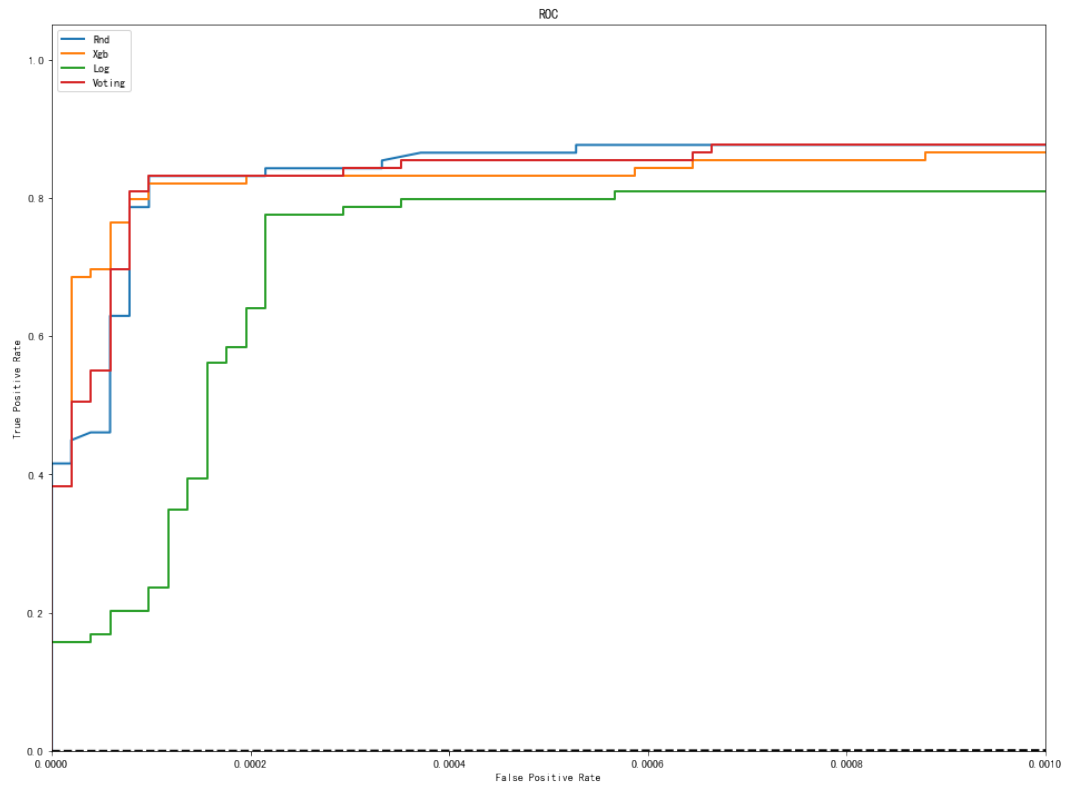
1 def roc_curve_all_models():
2     plt.figure(figsize=(16, 12))
3     plt.plot(fpr_rf, tpr_rf, label = 'ROC curve', linewidth = 2)
4     plt.plot(fpr_xgb, tpr_xgb, label = 'ROC curve', linewidth = 2)
5     plt.plot(fpr_log, tpr_log, label = 'ROC curve', linewidth = 2)
6     plt.plot(fpr_voting, tpr_voting, label = 'ROC curve', linewidth = 2)
7     plt.plot([0,1],[0,1], 'k--', linewidth = 2)
8     plt.xlim([0.0,0.001])
9     plt.ylim([0.0,1.05])
10    plt.xlabel('False Positive Rate')
11    plt.ylabel('True Positive Rate')
12    plt.title('ROC')
13    plt.legend(['Rnd','Xgb', 'Log', 'Voting'], loc='upper left')
14    plt.savefig('6.roc.png')
15    plt.show();

```

```

1 roc_curve_all_models ()

```



8.4 比较 Precision - Recall 曲线（所有模型）

```

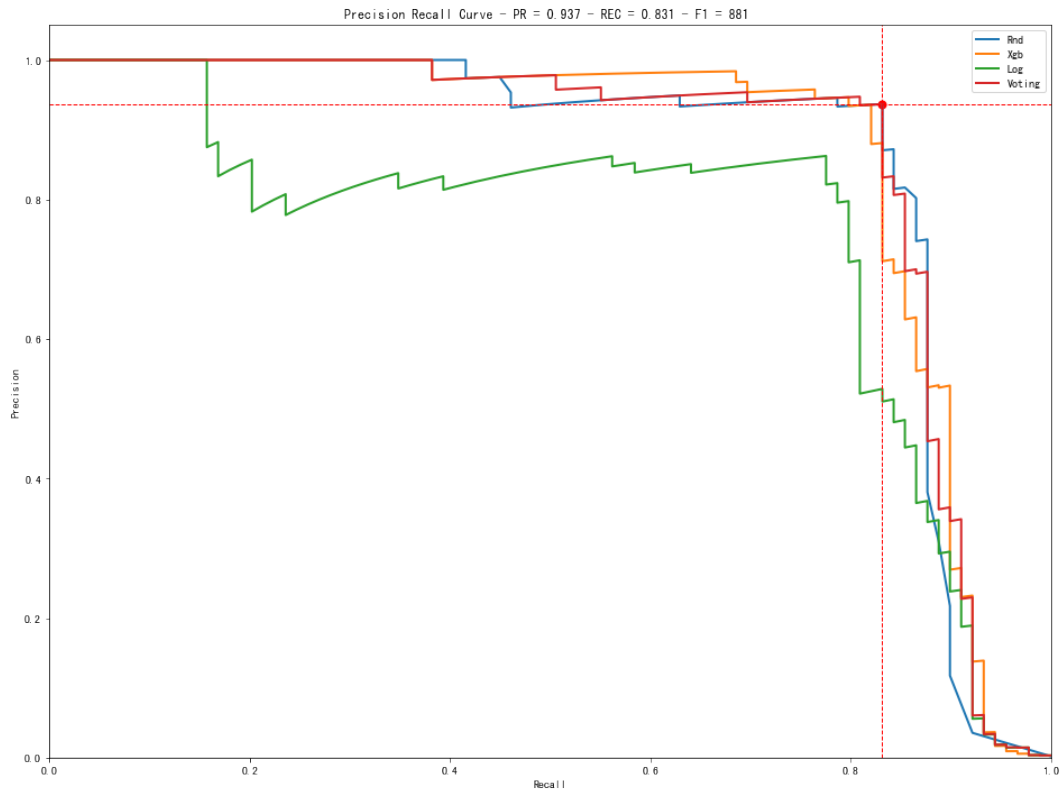
1  def prec_recall_all_models():
2      plt.figure(figsize=(16, 12))
3      plt.plot(recall_rf, precision_rf, linewidth = 2)
4      plt.plot(recall_xgb, precision_xgb, linewidth = 2)
5      plt.plot(recall_log, precision_log, linewidth = 2)
6      plt.plot(recall_voting, precision_voting, linewidth = 2)
7      plt.scatter(rec, pr, linewidth = 2, color = 'red')
8      plt.axvline(rec, color = 'red', linewidth = 1, linestyle='--')
9      plt.axhline(pr, color = 'red', linewidth = 1, linestyle='--')
10     plt.xlim([0.0,1])
11     plt.ylim([0.0,1.05])
12     plt.xlabel('Recall')
13     plt.ylabel('Precision')
14     plt.title('Precision Recall Curve - PR = 0.937 - REC = 0.831 - F1 = 881 ')
15     plt.legend(['Rnd', 'Xgb', 'Log', 'Voting'], loc='upper right')
16     plt.savefig('7.prec_recc.png')
17     plt.show();

```

```

1  prec_recall_all_models ()

```



9. VotingClassifier: Validation (F1 = 0.884)

用测试集数据 test_data 评估模型性能

9.1 Amount归一化、去除无用变量并定义 X,y

```
1 # 归一化 Amount
2 from sklearn.preprocessing import StandardScaler
3 test_data['normAmount'] = StandardScaler().fit_transform(test_data['Amount'].values.reshape(-1,1))
```

```
1 # 删除 test_data 中的time和Amount 字段
2 test_data = test_data.drop(['Amount', 'Time'], axis=1)
```

```
1 # 定义 X & y
2 y = np.array(test_data.Class.tolist())
3 test_data = test_data.drop('Class', 1)
4 X = np.array(test_data.as_matrix())
5
6 print X_test.shape
7 print X.shape
```

```
1 (51265, 29)
2 (28481, 29)
```

9.2 VotingClassifier在测试集上的表现

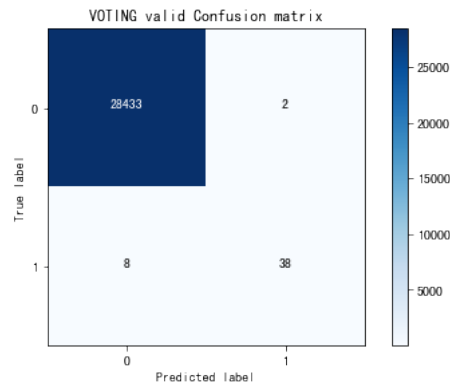
```
1 y_pred = voting_cfl.predict(X)
2 y_score = voting_cfl.predict_proba(X)[:,1]
```

```
1 # 混淆矩阵
2 cm = confusion_matrix(y, y_pred)
3 class_names = [0,1]
4 plt.figure()
5 plot_confusion_matrix(cm,
6                       classes = class_names,
7                       title = 'VOTING valid Confusion matrix')
8 plt.savefig('8.votingvf_cfl_confusion_matrix.png')
9 plt.show()
10
11 show_metrics()
12
13 #ROC
14 fpr, tpr, t = roc_curve(y, y_score)
```

```

15 plot_roc()
16
17 #precision recall
18 precision, recall, thresholds = precision_recall_curve(y, y_score)
19 plot_precision_recall()

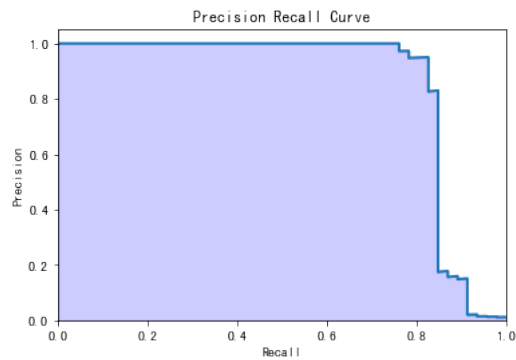
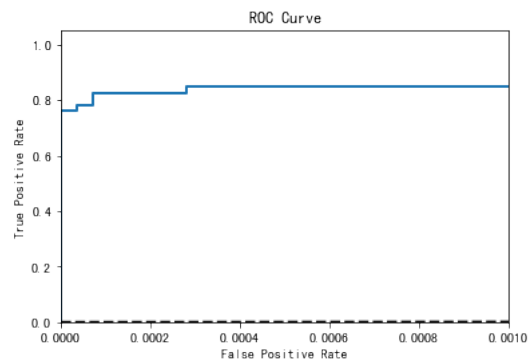
```



```

1 Precision = 0.950
2 Recall    = 0.826
3 F1_score  = 0.884

```



9.3 阈值选择

```

1 thresholds_adj = [0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9]
2
3 plt.figure(figsize = (15,15))
4
5 j = 1
6 for i in thresholds_adj:
7     y_score = voting_cfl.predict_proba(X)[:,-1] > i
8
9     plt.subplot(3,3,j)
10    j += 1
11
12    cm = confusion_matrix(y, y_score)
13
14    tp = cm[1,1]
15    fn = cm[1,0]
16    fp = cm[0,1]

```

```

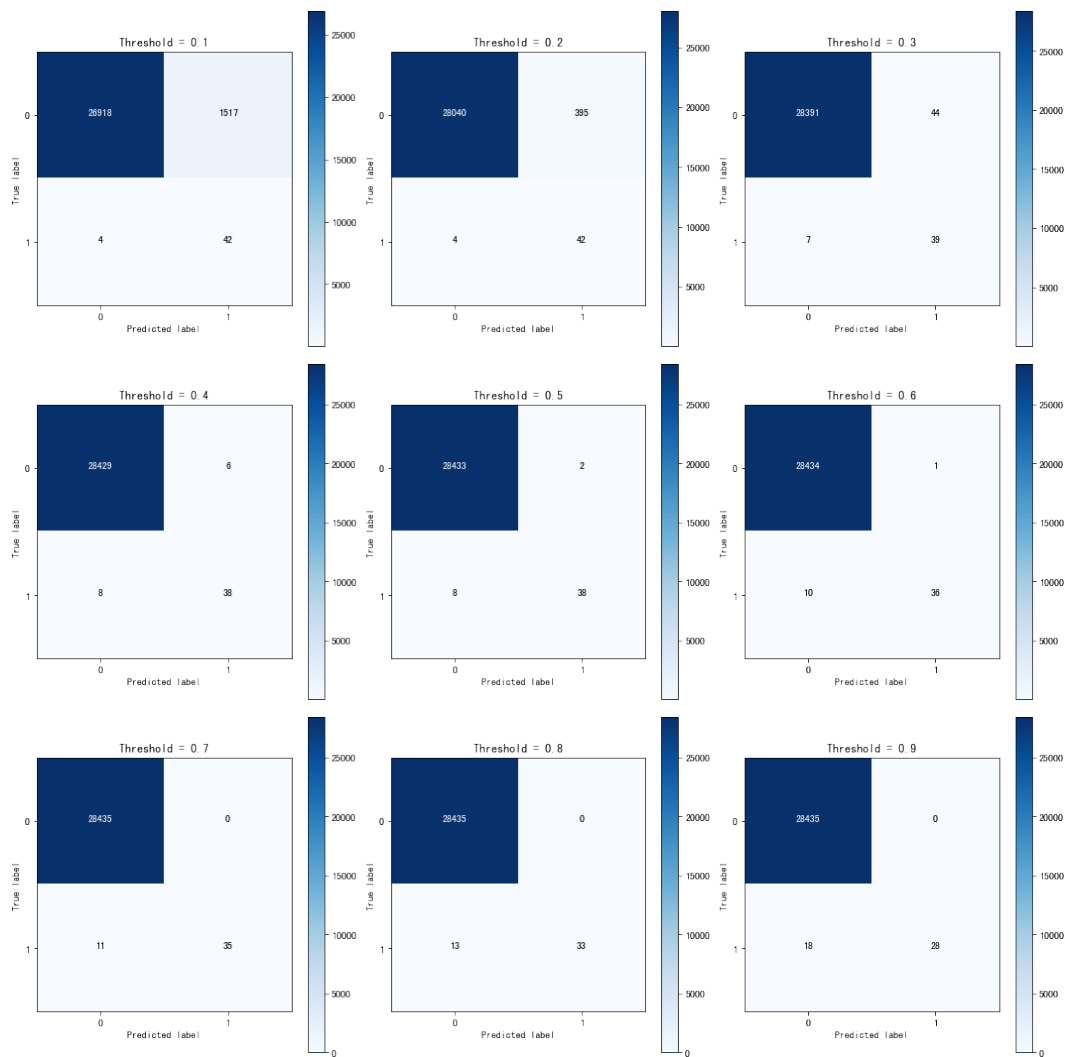
17     tn = cm[0,0]
18
19     print('F1_score w/ threshold = %s :'%i, (2*((tp/(tp+fp))*(tp/(tp+fn)))/
20           ((tp/(tp+fp))+(tp/(tp+fn))))))
21
22     class_names = [0,1]
23     plot_confusion_matrix(cm,
24                           classes=class_names,
25                           title='Threshold = %s'%i)
26
27     plt.savefig('9.confusion_matrix_threshold_select.png')

```

```

1 ('F1_score w/ threshold = 0.1 :', 0.052336448598130844)
2 ('F1_score w/ threshold = 0.2 :', 0.17391304347826086)
3 ('F1_score w/ threshold = 0.3 :', 0.6046511627906976)
4 ('F1_score w/ threshold = 0.4 :', 0.8444444444444444)
5 ('F1_score w/ threshold = 0.5 :', 0.8837209302325583)
6 ('F1_score w/ threshold = 0.6 :', 0.8674698795180723)
7 ('F1_score w/ threshold = 0.7 :', 0.8641975308641976)
8 ('F1_score w/ threshold = 0.8 :', 0.8354430379746834)
9 ('F1_score w/ threshold = 0.9 :', 0.7567567567567568)

```



可见阈值选择0.5时测试集上的 F1 值最大。