# prabhudayala@gmail.com_FB_Models

November 24, 2019

Social network Graph Link Prediction - Facebook Challenge

```
[1]: #Importing Libraries
# please do go through this python notebook:
import warnings
warnings.filterwarnings("ignore")

import csv
import pandas as pd#pandas to create small dataframes
import datetime #Convert to unix time
import time #Convert to unix time
# if numpy is not installed already : pip3 install numpy
import numpy as np#Do aritmetic operations on arrays
# matplotlib: used to plot graphs
import matplotlib
import matplotlib.pylab as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots
from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
import math
import pickle
import os
# to install xgboost: pip3 install xgboost
import xgboost as xgb

import warnings
import networkx as nx
import pdb
import pickle
from pandas import HDFStore,DataFrame
from pandas import read_hdf
from scipy.sparse.linalg import svds, eigs
import gc
from tqdm import tqdm
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
```

```python
[2]: #reading
     from pandas import read_hdf
     df_final_train = read_hdf('data/fea_sample/storage_sample_stage6.h5',␣
      ↪'train_df',mode='r')
     df_final_test = read_hdf('data/fea_sample/storage_sample_stage6.h5',␣
      ↪'test_df',mode='r')
```

```python
[3]: df_final_train.columns
```

```python
[3]: Index(['source_node', 'destination_node', 'indicator_link',
            'jaccard_followers', 'jaccard_followees', 'cosine_followers',
            'cosine_followees', 'num_followers_s', 'num_followees_s',
            'num_followees_d', 'inter_followers', 'inter_followees', 'adar_index',
            'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_out',
            'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
            'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities_s',
            'authorities_d', 'svd_dot_u', 'svd_dot_v', 'preferential_attachment_p',
            'preferential_attachment_s'],
           dtype='object')
```

```python
[4]: y_train = df_final_train.indicator_link
     y_test = df_final_test.indicator_link
```
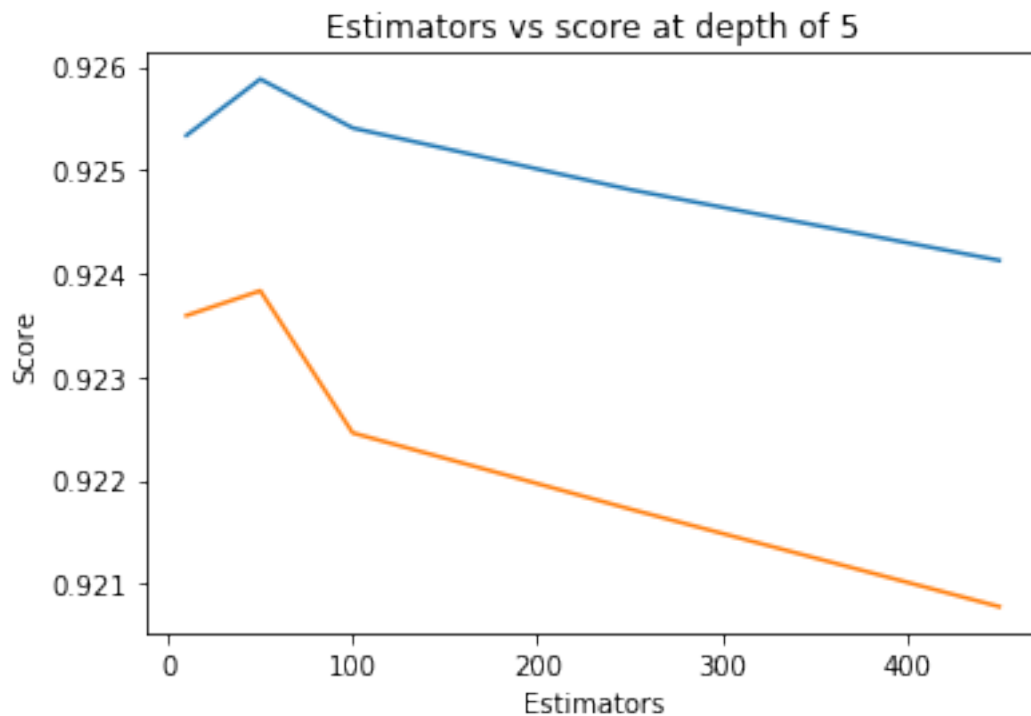
```python
[5]: df_final_train.drop(['source_node',␣
      ↪'destination_node','indicator_link'],axis=1,inplace=True)
     df_final_test.drop(['source_node',␣
      ↪'destination_node','indicator_link'],axis=1,inplace=True)
```

```python
[6]: estimators = [10,50,100,250,450]
     train_scores = []
     test_scores = []
     for i in estimators:
         clf = RandomForestClassifier(bootstrap=True, class_weight=None,␣
      ↪criterion='gini',
                 max_depth=5, max_features='auto', max_leaf_nodes=None,
                 min_impurity_decrease=0.0, min_impurity_split=None,
                 min_samples_leaf=52, min_samples_split=120,
                 min_weight_fraction_leaf=0.0, n_estimators=i,␣
      ↪n_jobs=-1,random_state=25,verbose=0,warm_start=False)
         clf.fit(df_final_train,y_train)
         train_sc = f1_score(y_train,clf.predict(df_final_train))
         test_sc = f1_score(y_test,clf.predict(df_final_test))
         test_scores.append(test_sc)
         train_scores.append(train_sc)
         print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
     plt.plot(estimators,train_scores,label='Train Score')
     plt.plot(estimators,test_scores,label='Test Score')
     plt.xlabel('Estimators')
     plt.ylabel('Score')
```

```
plt.title('Estimators vs score at depth of 5')
```

```
Estimators =   10 Train Score 0.9253332491662195 test Score 0.9235924932975871
Estimators =   50 Train Score 0.9258777360391177 test Score 0.9238324060023605
Estimators =   100 Train Score 0.9254051778572931 test Score 0.9224591584419976
Estimators =   250 Train Score 0.9248103689520584 test Score 0.9217240577504388
Estimators =   450 Train Score 0.9241257602655546 test Score 0.9207803724504877
```

[6]: Text(0.5, 1.0, 'Estimators vs score at depth of 5')



[7]:
```
depths = [3,9,11,15,20,35,50,70,130]
train_scores = []
test_scores = []
for i in depths:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None,
 →criterion='gini',
            max_depth=i, max_features='auto', max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=52, min_samples_split=120,
            min_weight_fraction_leaf=0.0, n_estimators=115,
 →n_jobs=-1,random_state=25,verbose=0,warm_start=False)
    clf.fit(df_final_train,y_train)
    train_sc = f1_score(y_train,clf.predict(df_final_train))
```
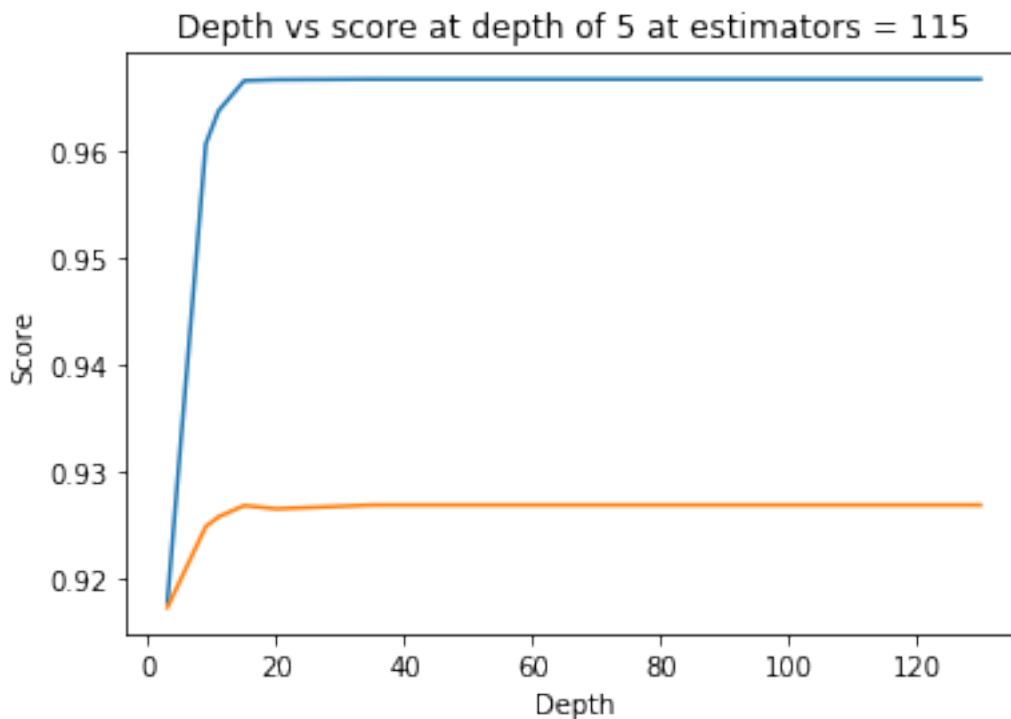
```
    test_sc = f1_score(y_test,clf.predict(df_final_test))
    test_scores.append(test_sc)
    train_scores.append(train_sc)
    print('depth = ',i,'Train Score',train_sc,'test Score',test_sc)
plt.plot(depths,train_scores,label='Train Score')
plt.plot(depths,test_scores,label='Test Score')
plt.xlabel('Depth')
plt.ylabel('Score')
plt.title('Depth vs score at depth of 5 at estimators = 115')
plt.show()
```

```
depth =   3 Train Score 0.917830714226324 test Score 0.9173315510255232
depth =   9 Train Score 0.960628960242471 test Score 0.9248893572181244
depth =  11 Train Score 0.963722253263575 test Score 0.925804278255011
depth =  15 Train Score 0.9665299159185995 test Score 0.9268498407945511
depth =  20 Train Score 0.9666257145468022 test Score 0.9265477218629538
depth =  35 Train Score 0.9666737724721097 test Score 0.9268981208393022
depth =  50 Train Score 0.9666737724721097 test Score 0.9268981208393022
depth =  70 Train Score 0.9666737724721097 test Score 0.9268981208393022
depth =  130 Train Score 0.9666737724721097 test Score 0.9268981208393022
```



[8]:
```
from sklearn.metrics import f1_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
```

```python
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint
from scipy.stats import uniform

param_dist = {"n_estimators":sp_randint(105,125),
              "max_depth": sp_randint(10,15),
              "min_samples_split": sp_randint(110,190),
              "min_samples_leaf": sp_randint(25,65)}

clf = RandomForestClassifier(random_state=25,n_jobs=-1)

rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
  ↳n_iter=5,cv=10,scoring='f1',random_state=25,return_train_score=True)

rf_random.fit(df_final_train,y_train)
print('mean test scores',rf_random.cv_results_['mean_test_score'])
print('mean train scores',rf_random.cv_results_['mean_train_score'])
```

```
mean test scores [0.96454843 0.96396638 0.96254065 0.96390786 0.96553771]
mean train scores [0.96591799 0.96493352 0.96334515 0.96515447 0.96724712]
```

[9]:
```python
print(rf_random.best_estimator_)
```

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                       max_depth=14, max_features='auto', max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=28, min_samples_split=111,
                       min_weight_fraction_leaf=0.0, n_estimators=121,
                       n_jobs=-1, oob_score=False, random_state=25, verbose=0,
                       warm_start=False)
```

[10]:
```python
clf = RandomForestClassifier(bootstrap=True, class_weight=None,
  ↳criterion='gini',
             max_depth=14, max_features='auto', max_leaf_nodes=None,
             min_impurity_decrease=0.0, min_impurity_split=None,
             min_samples_leaf=28, min_samples_split=111,
             min_weight_fraction_leaf=0.0, n_estimators=121, n_jobs=-1,
             oob_score=False, random_state=25, verbose=0, warm_start=False)
```

[11]:
```python
clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)
```

[12]:
```python
from sklearn.metrics import f1_score
print('Train f1 score',f1_score(y_train,y_train_pred))
print('Test f1 score',f1_score(y_test,y_test_pred))
```

```
Train f1 score 0.9679107505070993
Test f1 score 0.9273340358271865
```

[13]:
```python
from sklearn.metrics import confusion_matrix
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)

    A =(((C.T)/(C.sum(axis=1))).T)

    B =(C/C.sum(axis=0))
    plt.figure(figsize=(20,4))

    labels = [0,1]
    # representing A in heatmap format
    cmap=sns.light_palette("blue")
    plt.subplot(1, 3, 1)
    sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels,␣
 ↪yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Confusion matrix")

    plt.subplot(1, 3, 2)
    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels,␣
 ↪yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Precision matrix")

    plt.subplot(1, 3, 3)
    # representing B in heatmap format
    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels,␣
 ↪yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Recall matrix")

    plt.show()
```
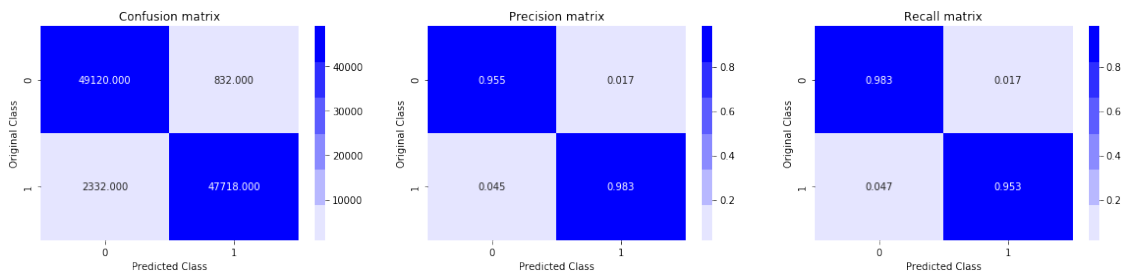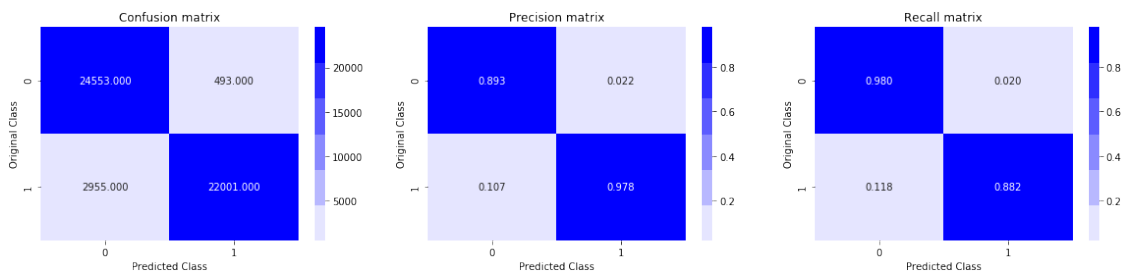
[14]:
```python
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```
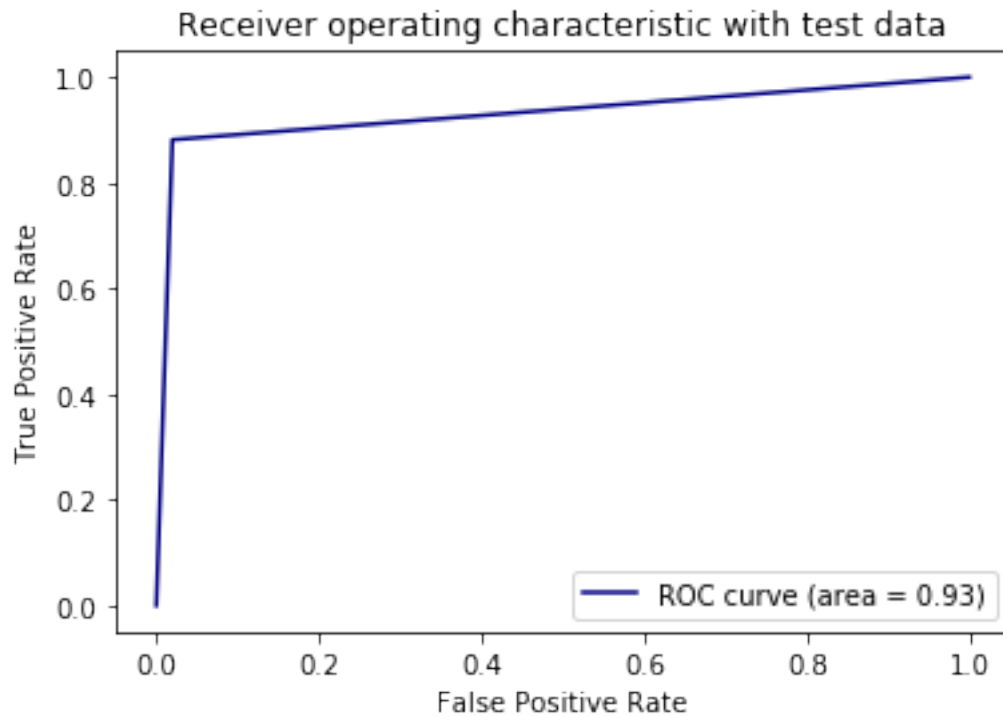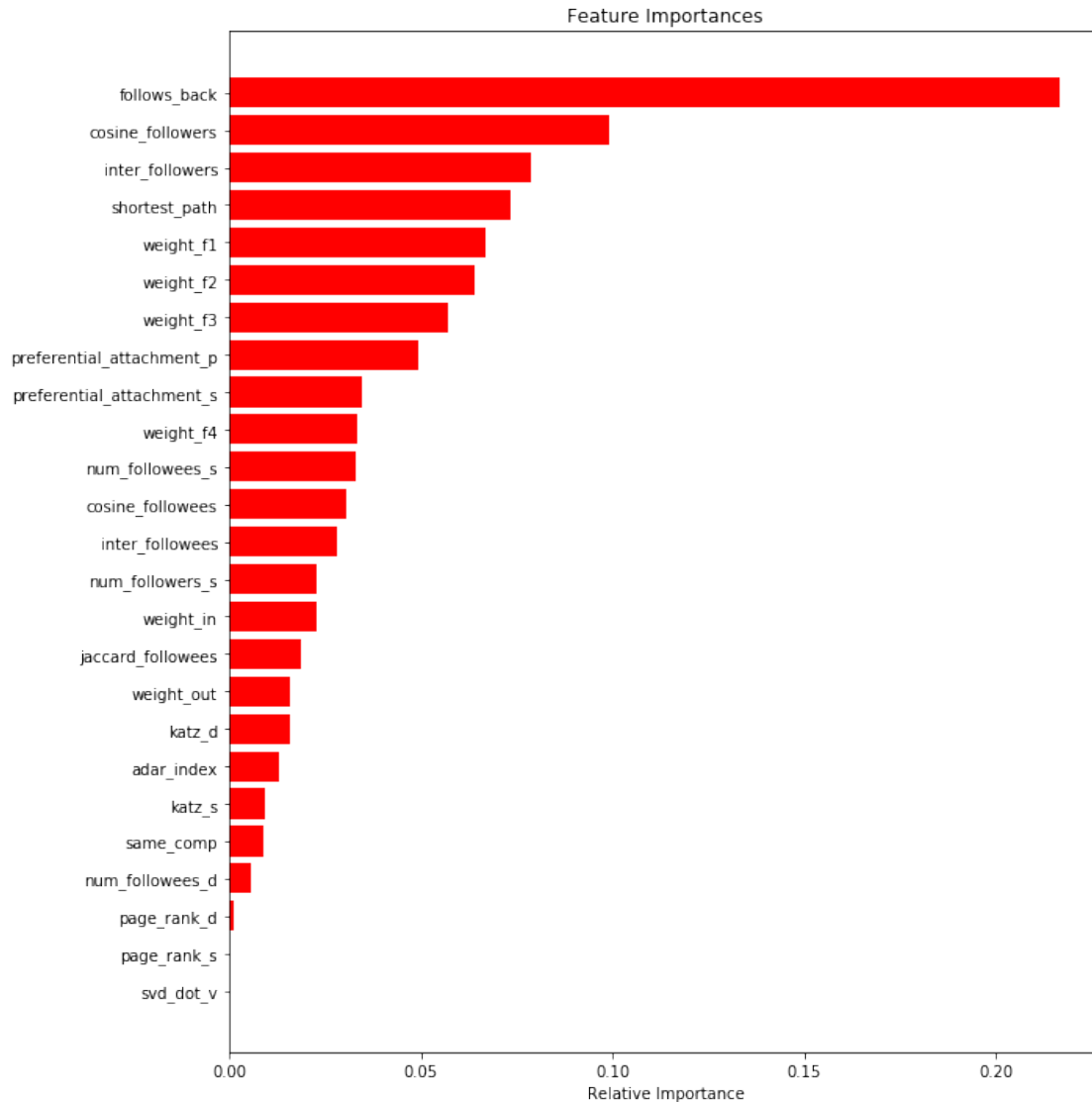
```
Train confusion_matrix
```

Confusion matrix / Precision matrix / Recall matrix (Train)

Test confusion_matrix



Confusion matrix / Precision matrix / Recall matrix (Test)

```
[15]: from sklearn.metrics import roc_curve, auc
fpr,tpr,ths = roc_curve(y_test,y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic with test data')
plt.legend()
plt.show()
```

Receiver operating characteristic with test data

```
features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-25:]
plt.figure(figsize=(10,12))
plt.title('Feature Importances')
plt.barh(range(len(indices)), importances[indices], color='r', align='center')
plt.yticks(range(len(indices)), [features[i] for i in indices])
plt.xlabel('Relative Importance')
plt.show()
```

Feature Importances



# 1 Assignments:

1. Add another feature called Preferential Attachment with followers and followees data of vertex. you can check about Preferential Attachment in below link http://be.amazd.com/link-prediction/
2. Add feature called svd_dot. you can calculate svd_dot as Dot product between sourse node svd and destination node svd features. you can read about this in below pdf https://storage.googleapis.com/kaggle-forum-message-attachments/2594/supervised_link_prediction.pdf
3. Tune hyperparameters for XG boost with all these features and check the error metric.

```python
from sklearn.metrics import f1_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint
from scipy.stats import uniform


param_dist = {
 'max_depth':range(3,10,2),
 'min_child_weight':range(1,6,2),
 'gamma':[i/10.0 for i in range(0,5)],
 'subsample':[i/10.0 for i in range(6,10)],
 'colsample_bytree':[i/10.0 for i in range(6,10)],
 'reg_alpha':[1e-5, 1e-2, 0.1, 1, 100]
}

clf = xgb.XGBClassifier()

rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
                              ⎵
 ↪n_iter=10,cv=10,scoring='f1',random_state=25,return_train_score=True,verbose=1,n_jobs=4)

rf_random.fit(df_final_train,y_train)
print('mean test scores',rf_random.cv_results_['mean_test_score'])
print('mean train scores',rf_random.cv_results_['mean_train_score'])
```

```
Fitting 10 folds for each of 10 candidates, totalling 100 fits

[Parallel(n_jobs=4)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=4)]: Done  42 tasks      | elapsed:  1.5min
[Parallel(n_jobs=4)]: Done 100 out of 100 | elapsed:  4.6min finished

mean test scores [0.9732647  0.97300458 0.97330321 0.97281577 0.97739804
0.97315316
 0.97741073 0.97336856 0.97759844 0.9777435 ]
mean train scores [0.97368159 0.97329602 0.97354947 0.97303076 0.98352133
0.97335821
 0.98397578 0.97349323 0.98407541 0.98462509]
```

```python
print(rf_random.best_estimator_)
```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=0.8, gamma=0.3,
              learning_rate=0.1, max_delta_step=0, max_depth=9,
              min_child_weight=3, missing=None, n_estimators=100, n_jobs=1,
              nthread=None, objective='binary:logistic', random_state=0,
```

```
                reg_alpha=0.01, reg_lambda=1, scale_pos_weight=1, seed=None,
                silent=None, subsample=0.7, verbosity=1)
```

```
[21]: clf = xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                    colsample_bynode=1, colsample_bytree=0.8, gamma=0.3,
                    learning_rate=0.1, max_delta_step=0, max_depth=9,
                    min_child_weight=3, missing=None, n_estimators=100, n_jobs=1,
                    nthread=None, objective='binary:logistic', random_state=0,
                    reg_alpha=0.01, reg_lambda=1, scale_pos_weight=1, seed=None,
                    silent=None, subsample=0.7)
```

```
[22]: clf.fit(df_final_train,y_train)
      y_train_pred = clf.predict(df_final_train)
      y_test_pred = clf.predict(df_final_test)
```

```
[23]: from sklearn.metrics import f1_score
      print('Train f1 score',f1_score(y_train,y_train_pred))
      print('Test f1 score',f1_score(y_test,y_test_pred))
```

```
Train f1 score 0.9842347560668486
Test f1 score 0.9270541672806888
```

```
[24]: from sklearn.metrics import confusion_matrix
      def plot_confusion_matrix(test_y, predict_y):
          C = confusion_matrix(test_y, predict_y)

          A =(((C.T)/(C.sum(axis=1))).T)

          B =(C/C.sum(axis=0))
          plt.figure(figsize=(20,4))

          labels = [0,1]
          # representing A in heatmap format
          cmap=sns.light_palette("blue")
          plt.subplot(1, 3, 1)
          sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels,␣
       ↪yticklabels=labels)
          plt.xlabel('Predicted Class')
          plt.ylabel('Original Class')
          plt.title("Confusion matrix")

          plt.subplot(1, 3, 2)
          sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels,␣
       ↪yticklabels=labels)
          plt.xlabel('Predicted Class')
          plt.ylabel('Original Class')
          plt.title("Precision matrix")
```
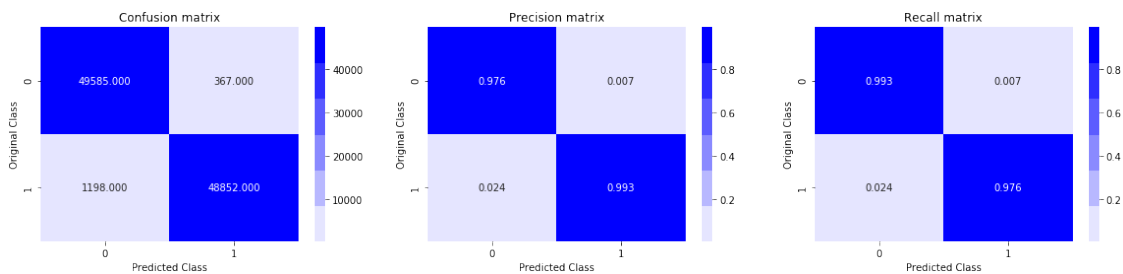
```
    plt.subplot(1, 3, 3)
    # representing B in heatmap format
    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels,␣
 ↪yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Recall matrix")

    plt.show()
```
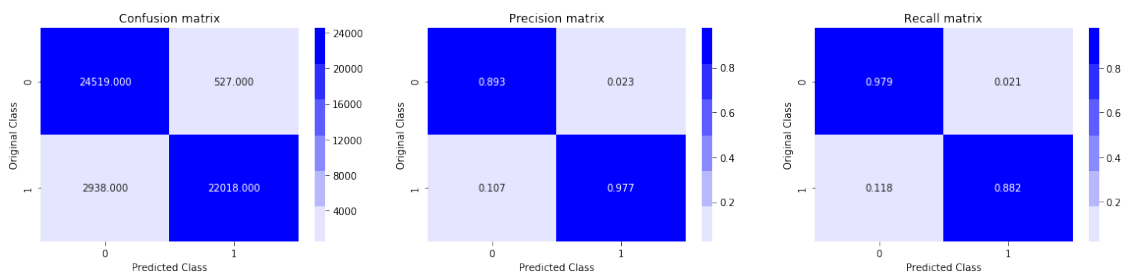
```
[25]: print('Train confusion_matrix')
      plot_confusion_matrix(y_train,y_train_pred)
      print('Test confusion_matrix')
      plot_confusion_matrix(y_test,y_test_pred)
```

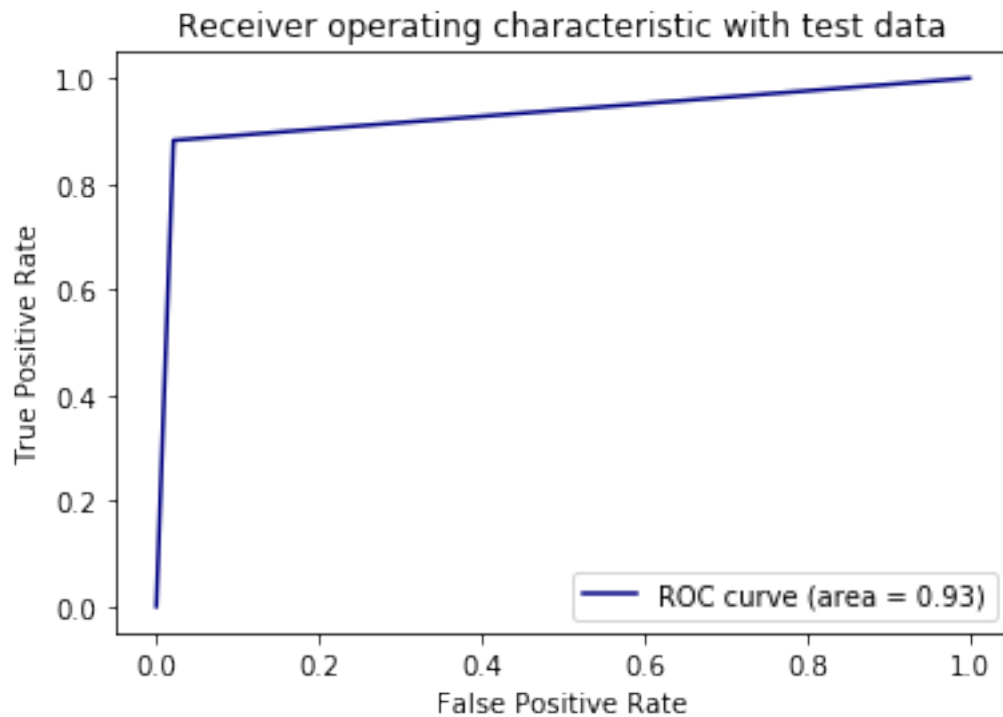Train confusion_matrix



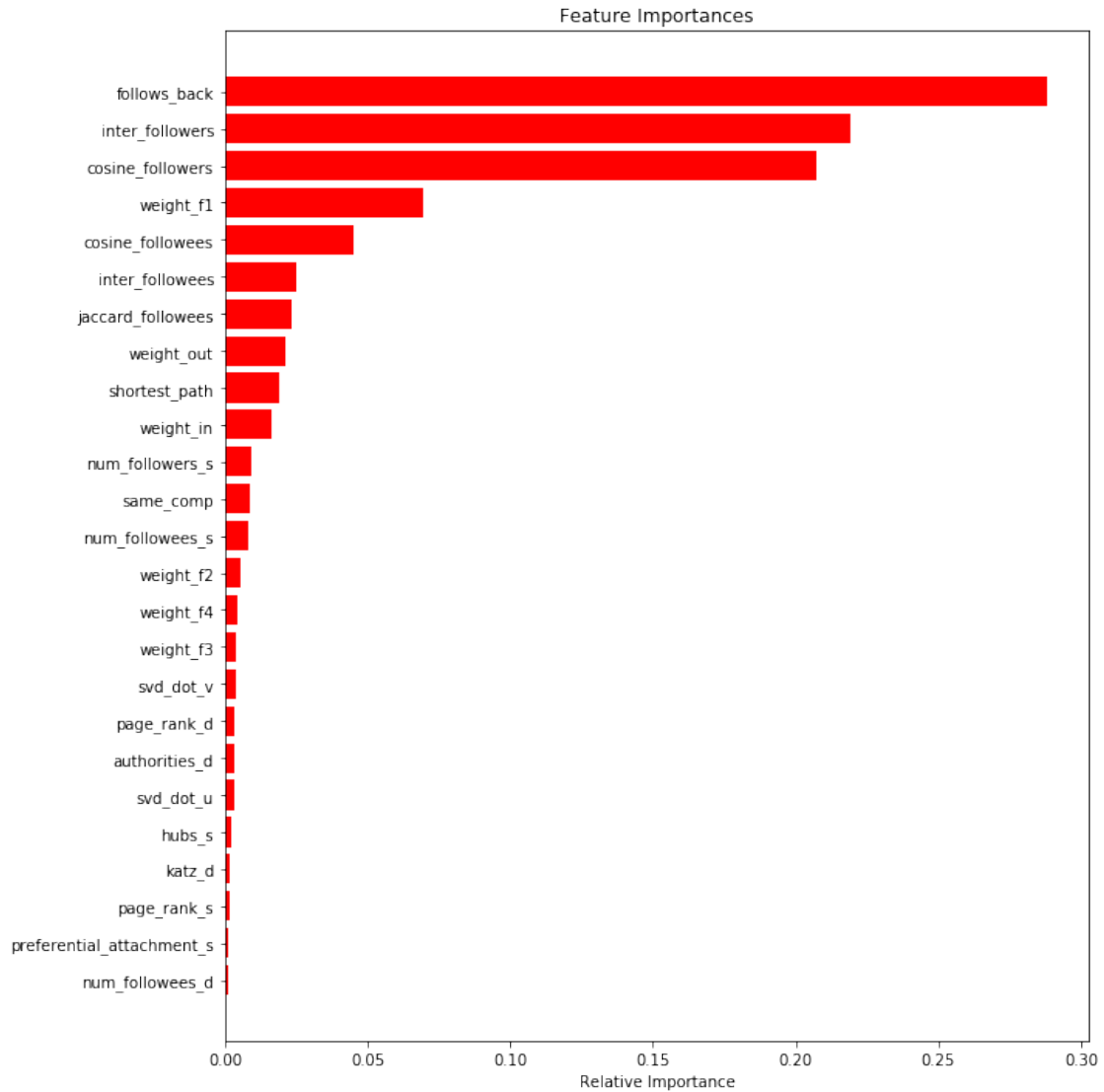Test confusion_matrix



```
[26]: from sklearn.metrics import roc_curve, auc
      fpr,tpr,ths = roc_curve(y_test,y_test_pred)
      auc_sc = auc(fpr, tpr)
      plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
      plt.xlabel('False Positive Rate')
      plt.ylabel('True Positive Rate')
      plt.title('Receiver operating characteristic with test data')
```

```
plt.legend()
plt.show()
```

### Receiver operating characteristic with test data



```
features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-25:]
plt.figure(figsize=(10,12))
plt.title('Feature Importances')
plt.barh(range(len(indices)), importances[indices], color='r', align='center')
plt.yticks(range(len(indices)), [features[i] for i in indices])
plt.xlabel('Relative Importance')
plt.show()
```

Feature Importances

```
[29]: from prettytable import PrettyTable
      x = PrettyTable()
      x.field_names = ["Model", "f1score Train", "f1score Test"]
      x.add_row(["Random forest",  round(0.9679107505070993,3), round(0.
       ↪9273340358271865,3)])
      x.add_row(["XGBoost",  round(0.9842347560668486,3), round(0.
       ↪9270541672806888,3)])
      x.border=True
      print(x)
```

```
+---------------+---------------+--------------+
|     Model     | f1score Train | f1score Test |
+---------------+---------------+--------------+
```

```
| Random forest |      0.968     |     0.927    |
|    XGBoost    |      0.984     |     0.927    |
+---------------+----------------+--------------+
```

   Due to cold start problem the test accuracy is not up to the mark as train. But certainly XGBoost is the best model for this scenario.