

Social network Graph Link Prediction - Facebook Challenge

```
from google.colab import drive
drive.mount('/content/drive')
%cd ./drive/My Drive/Facebook
```

↪ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.m
[Errno 2] No such file or directory: './drive/My Drive/Facebook'
/content/drive/My Drive/Facebook

```
%cd Facebook/
```

↪ [Errno 2] No such file or directory: 'Facebook/'
/content/drive/My Drive/Facebook

```
#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 arithmetic operations on arrays
# matplotlib: used to plot graphs
import matplotlib
import matplotlib.pyplot 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
```

```
if os.path.isfile('data/after_eda/train_pos_after_eda.csv'):
    train_graph=nx.read_edgelist('data/after_eda/train_pos_after_eda.csv',delimiter=',',cr
```

```

print(nx.info(train_graph))
else:
    print("please run the FB_EDA.ipynb or download the files from drive")

```

```

📄 Name:
Type: DiGraph
Number of nodes: 1780722
Number of edges: 7550015
Average in degree: 4.2399
Average out degree: 4.2399

```

```
#reading
```

```
from pandas import read_hdf
```

```
df_final_train = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'train_df',mode='r')
```

```
df_final_test = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'test_df',mode='r')
```

```
#preferential followers&followees
```

```
df_final_train['preferential followers']=df_final_train['num_followers_s']*df_final_train['
```

```
df_final_train['preferential followees']=df_final_train['num_followees_s']*df_final_train['
```

```
df_final_test['preferential followers']=df_final_test['num_followers_s']*df_final_test['nu
```

```
df_final_test['preferential followees']=df_final_test['num_followees_s']*df_final_test['nu
```

```
#svd dot product
```

```
#used code from https://github.com/krpiyush5/Facebook-Friend-Recommendation-using-Graph-Mi
```

```
#for train datasets
```

```
s1,s2,s3,s4,s5,s6=df_final_train['svd_u_s_1'],df_final_train['svd_u_s_2'],df_final_train['
s7,s8,s9,s10,s11,s12=df_final_train['svd_v_s_1'],df_final_train['svd_v_s_2'],df_final_trai
```

```
d1,d2,d3,d4,d5,d6=df_final_train['svd_u_d_1'],df_final_train['svd_u_d_2'],df_final_train['
d7,d8,d9,d10,d11,d12=df_final_train['svd_v_d_1'],df_final_train['svd_v_d_2'],df_final_trai
```

```
svd_dot=[]
```

```
for i in range(len(np.array(s1))):
```

```
    a=[]
```

```
    b=[]
```

```
    a.append(np.array(s1[i]))
```

```
    a.append(np.array(s2[i]))
```

```
    a.append(np.array(s3[i]))
```

```
    a.append(np.array(s4[i]))
```

```
    a.append(np.array(s5[i]))
```

```
    a.append(np.array(s6[i]))
```

```
    a.append(np.array(s7[i]))
```

```
    a.append(np.array(s8[i]))
```

```
    a.append(np.array(s9[i]))
```

```
    a.append(np.array(s10[i]))
```

```
    a.append(np.array(s11[i]))
```

```
    a.append(np.array(s12[i]))
```

```
    b.append(np.array(d1[i]))
```

```
    b.append(np.array(d2[i]))
```

```
    b.append(np.array(d3[i]))
```

```
    b.append(np.array(d4[i]))
```

```
    b.append(np.array(d5[i]))
```

```

    b.append(np.array(u5[i]))
    b.append(np.array(d6[i]))
    b.append(np.array(d7[i]))
    b.append(np.array(d8[i]))
    b.append(np.array(d9[i]))
    b.append(np.array(d10[i]))
    b.append(np.array(d11[i]))
    b.append(np.array(d12[i]))
    svd_dot.append(np.dot(a,b))
df_final_train['svd_dot']=svd_dot

```

```

#for test dataset

```

```

s1,s2,s3,s4,s5,s6=df_final_test['svd_u_s_1'],df_final_test['svd_u_s_2'],df_final_test['svd_u_s_3'],df_final_test['svd_u_s_4'],df_final_test['svd_u_s_5'],df_final_test['svd_u_s_6']
s7,s8,s9,s10,s11,s12=df_final_test['svd_v_s_1'],df_final_test['svd_v_s_2'],df_final_test['svd_v_s_3'],df_final_test['svd_v_s_4'],df_final_test['svd_v_s_5'],df_final_test['svd_v_s_6']

```

```

d1,d2,d3,d4,d5,d6=df_final_test['svd_u_d_1'],df_final_test['svd_u_d_2'],df_final_test['svd_u_d_3'],df_final_test['svd_u_d_4'],df_final_test['svd_u_d_5'],df_final_test['svd_u_d_6']
d7,d8,d9,d10,d11,d12=df_final_test['svd_v_d_1'],df_final_test['svd_v_d_2'],df_final_test['svd_v_d_3'],df_final_test['svd_v_d_4'],df_final_test['svd_v_d_5'],df_final_test['svd_v_d_6']

```

```

svd_dot=[]
for i in range(len(np.array(s1))):
    a=[]
    b=[]
    a.append(np.array(s1[i]))
    a.append(np.array(s2[i]))
    a.append(np.array(s3[i]))
    a.append(np.array(s4[i]))
    a.append(np.array(s5[i]))
    a.append(np.array(s6[i]))
    a.append(np.array(s7[i]))
    a.append(np.array(s8[i]))
    a.append(np.array(s9[i]))
    a.append(np.array(s10[i]))
    a.append(np.array(s11[i]))
    a.append(np.array(s12[i]))
    b.append(np.array(d1[i]))
    b.append(np.array(d2[i]))
    b.append(np.array(d3[i]))
    b.append(np.array(d4[i]))
    b.append(np.array(d5[i]))
    b.append(np.array(d6[i]))
    b.append(np.array(d7[i]))
    b.append(np.array(d8[i]))
    b.append(np.array(d9[i]))
    b.append(np.array(d10[i]))
    b.append(np.array(d11[i]))
    b.append(np.array(d12[i]))
    svd_dot.append(np.dot(a,b))
df_final_test['svd_dot']=svd_dot

```

```

df_final_train.columns

```

```

↳ Index(['source_node', 'destination_node', 'indicator_link',
        'jaccard_followers', 'jaccard_followees', 'cosine_followers',
        'cosine_followees', 'num_followers_s', 'num_followers_d',
        '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_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
        'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4',
        'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3',
        'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1', 'svd_v_d_2',
        'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6',
        'preferential_followers', 'preferential_followees', 'svd_dot'],
        dtype='object')

```

```

y_train = df_final_train.indicator_link
y_test = df_final_test.indicator_link

```

```

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)

```

```

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)
    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')

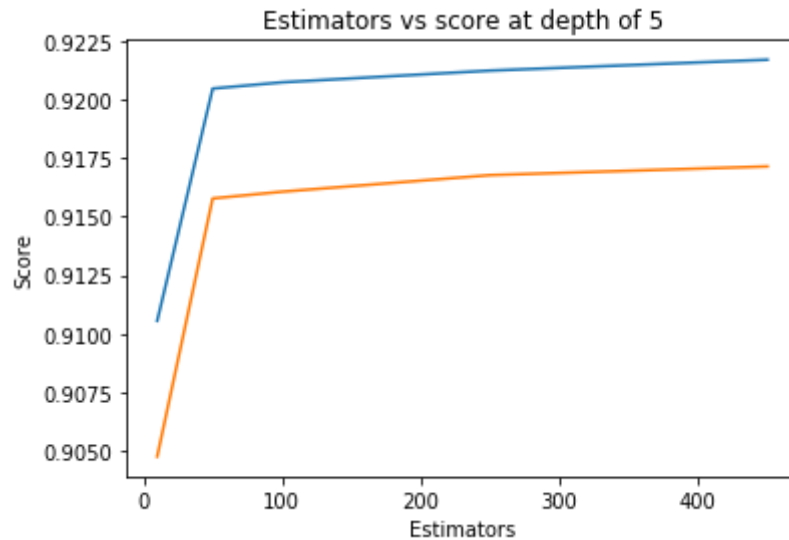
```

↳

```

Estimators = 10 Train Score 0.9105564546246732 test Score 0.9047538772757924
Estimators = 50 Train Score 0.9204793366212765 test Score 0.9157826187425162
Estimators = 100 Train Score 0.9207459689938525 test Score 0.9160706773243584
Estimators = 250 Train Score 0.9212468812439825 test Score 0.9167789190099342
Estimators = 450 Train Score 0.9217123103355483 test Score 0.9171505512076075
Text(0.5, 1.0, 'Estimators vs score at depth of 5')

```



```

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, verb
    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()

```

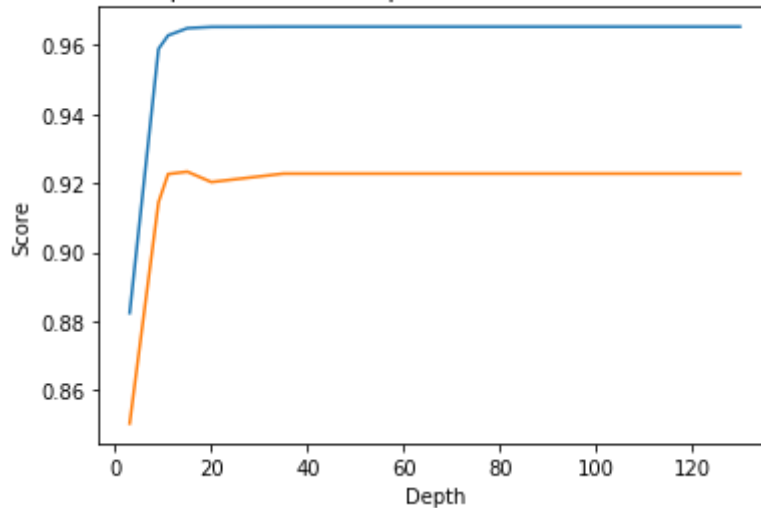


```

depth = 3 Train Score 0.8823547772887599 test Score 0.8503154437818685
depth = 9 Train Score 0.9588946878666504 test Score 0.9144305307096003
depth = 11 Train Score 0.9628045397225725 test Score 0.9226180402775428
depth = 15 Train Score 0.9648937135254453 test Score 0.9232951538119535
depth = 20 Train Score 0.9652764463101808 test Score 0.920290779237342
depth = 35 Train Score 0.965324452517913 test Score 0.922737681282334
depth = 50 Train Score 0.965324452517913 test Score 0.922737681282334
depth = 70 Train Score 0.965324452517913 test Score 0.922737681282334
depth = 130 Train Score 0.965324452517913 test Score 0.922737681282334

```

Depth vs score at depth of 5 at estimators = 115



```

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 = {"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)

rf_random.fit(df_final_train,y_train)

```



```

RandomizedSearchCV(cv=10, error_score='raise-deprecating',
                    estimator=RandomForestClassifier(bootstrap=True,
                                                        class_weight=None,
                                                        criterion='gini',
                                                        max_depth=None,
                                                        max_features='auto',
                                                        max_leaf_nodes=None,
                                                        min_impurity_decrease=0.0,
                                                        min_impurity_split=None,
                                                        min_samples_leaf=1,
                                                        min_samples_split=2,
                                                        min_weight_fraction_leaf=0.0,
                                                        n_estimators='warn',
                                                        n_jobs=-1, oob_sco...
                    'min_samples_leaf': <scipy.stats._distn_infra
                    'min_samples_split': <scipy.stats._distn_infr
                    'n_estimators': <scipy.stats._distn_infrastru
                    pre_dispatch='2*n_jobs', random_state=25, refit=True,
                    return_train_score=False, scoring='f1', verbose=0)

```

```
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)

```

```

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)

```

```

clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)

```

```

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 0.9653662179344864
   Test f1 score 0.9260261233461555

```

```

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))
```

```

D = (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")

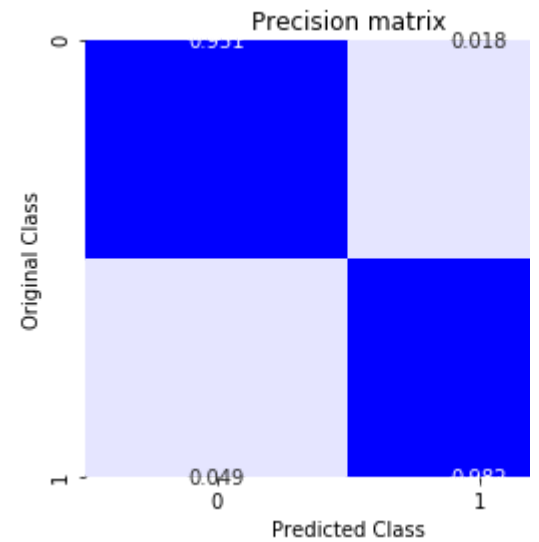
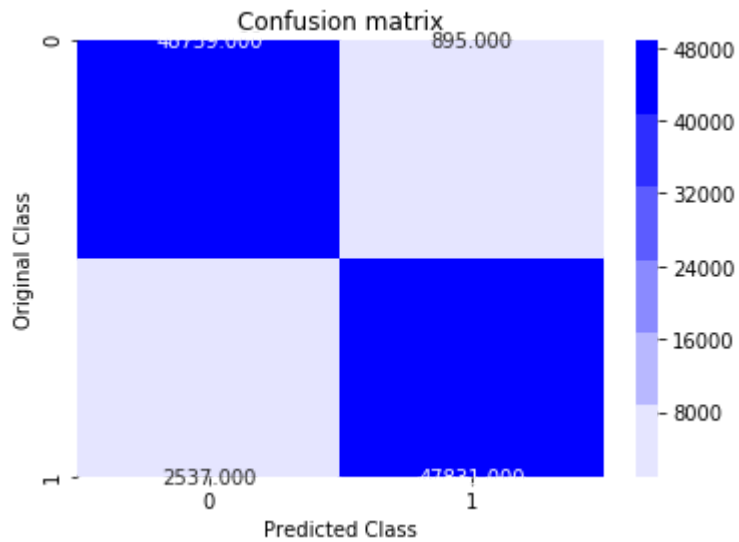
plt.show()

print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)

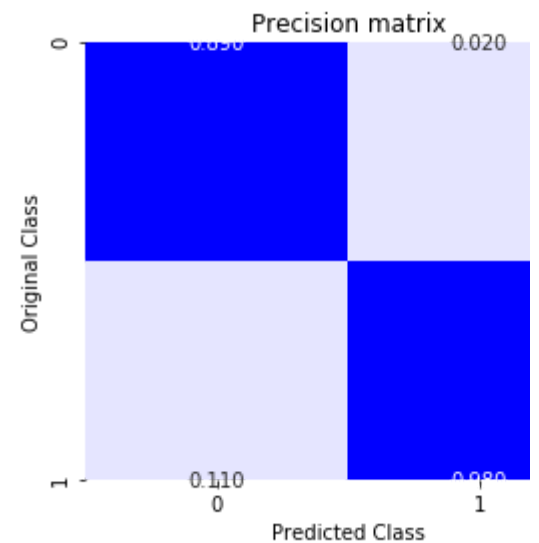
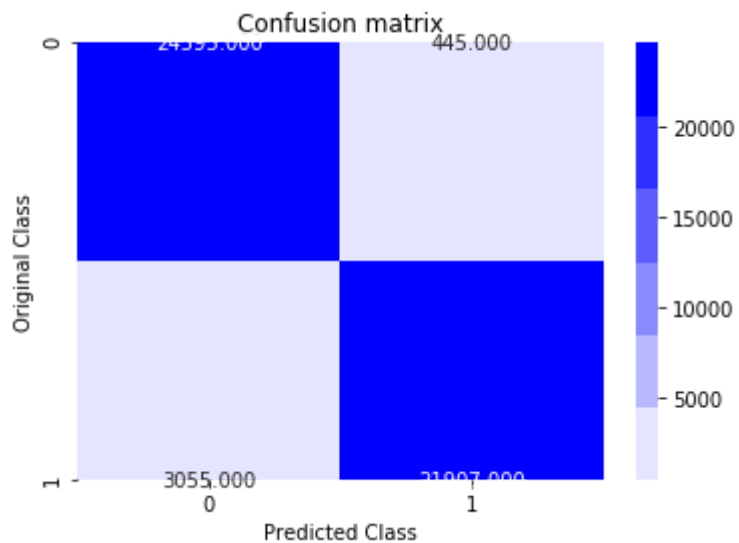
```



Train confusion_matrix



Test confusion_matrix

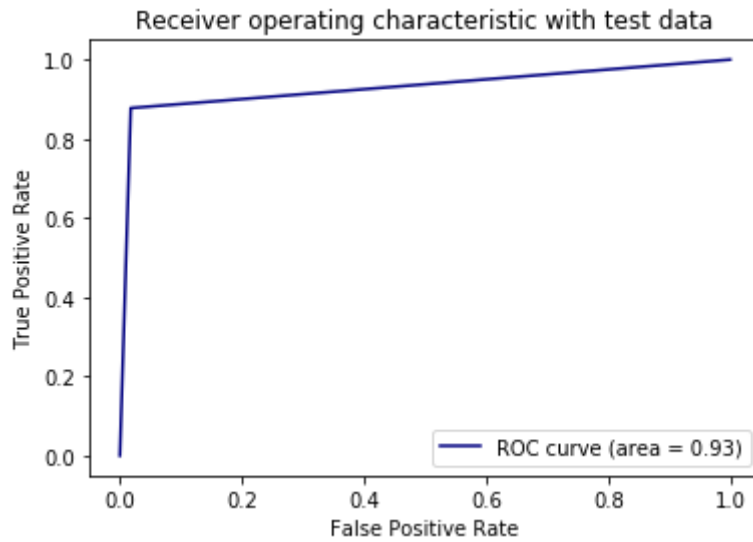


```

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()

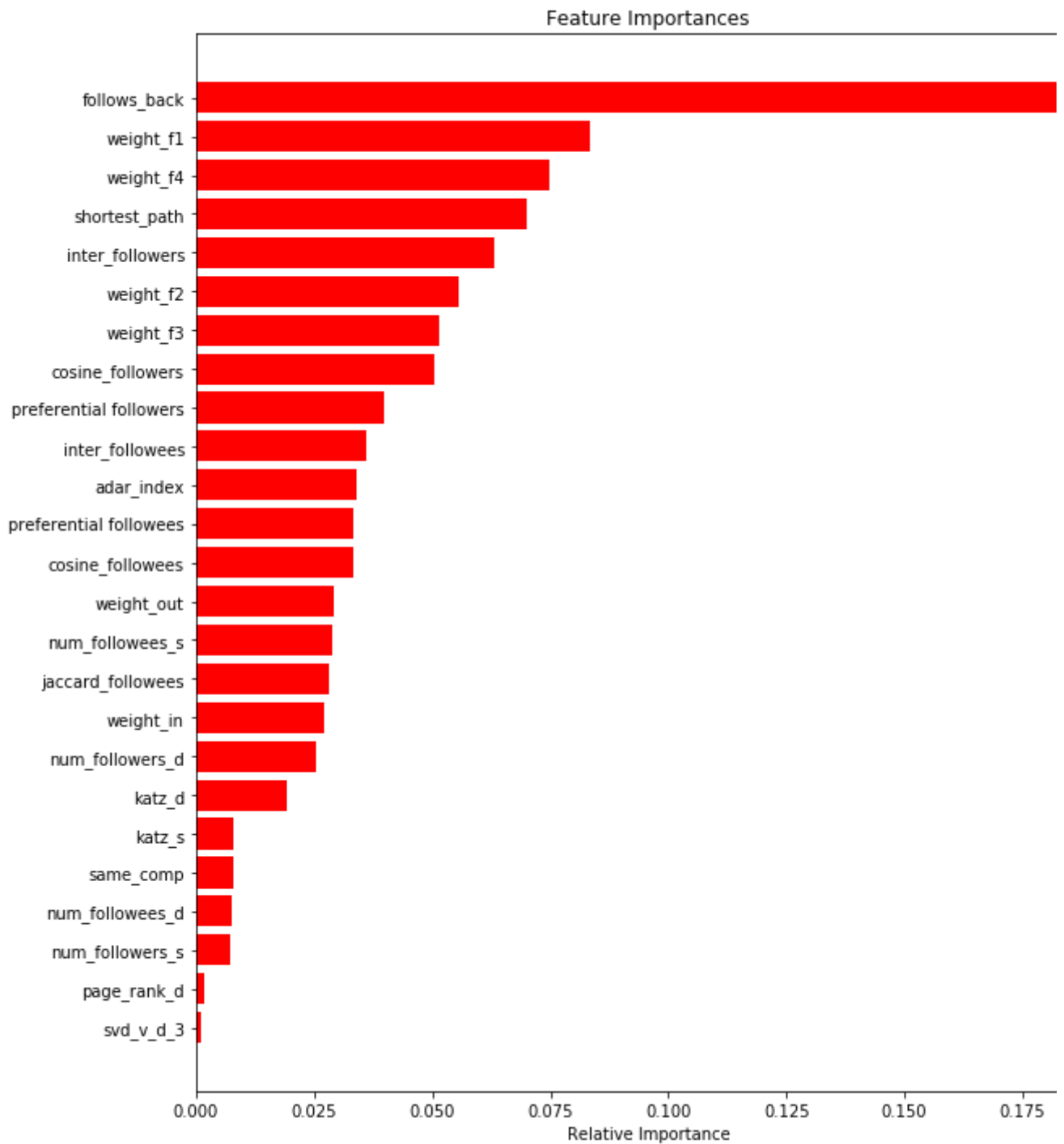
```





```
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()
```





▼ XGBOOST

```
import xgboost as xgb
clf = xgb.XGBClassifier()
param_dist = {"n_estimators": sp_randint(105,125),
              "max_depth": sp_randint(10,15)
              }
model = RandomizedSearchCV(clf, param_distributions=param_dist,
                           n_iter=5, cv=3, scoring='f1', random_state=25)

model.fit(df_final_train, y_train)
```

```

↳ RandomizedSearchCV(cv=3, error_score='raise-deprecating',
                      estimator=XGBClassifier(base_score=0.5, booster='gbtree',
                                             colsample_bylevel=1,
                                             colsample_bynode=1,
                                             colsample_bytree=1, gamma=0,
                                             learning_rate=0.1, max_delta_step=0,
                                             max_depth=3, min_child_weight=1,
                                             missing=None, n_estimators=100,
                                             n_jobs=1, nthread=None,
                                             objective='binary:logistic',
                                             random_state=0, reg_alpha=0...
                                             seed=None, silent=None, subsample=1,
                                             verbosity=1),
                      iid='warn', n_iter=5, n_jobs=None,
                      param_distributions={'max_depth': <scipy.stats._distn_infrastructu
                                             'n_estimators': <scipy.stats._distn_infrastru
                      pre_dispatch='2*n_jobs', random_state=25, refit=True,
                      return_train_score=False, scoring='f1', verbose=0)

```

```
print(model.best_estimator_)
```

```

↳ XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
               colsample_bynode=1, colsample_bytree=1, gamma=0,
               learning_rate=0.1, max_delta_step=0, max_depth=11,
               min_child_weight=1, missing=None, n_estimators=112, n_jobs=1,
               nthread=None, objective='binary:logistic', random_state=0,
               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
               silent=None, subsample=1, verbosity=1)

```

```

clf = xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                       colsample_bynode=1, colsample_bytree=1, gamma=0,
                       learning_rate=0.1, max_delta_step=0, max_depth=11,
                       min_child_weight=1, missing=None, n_estimators=112, n_jobs=1,
                       nthread=None, objective='binary:logistic', random_state=0,
                       reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                       silent=None, subsample=1, verbosity=1)

```

```

clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)

```

```

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 0.9955397144620777
   Test f1 score 0.9278031187016017

```

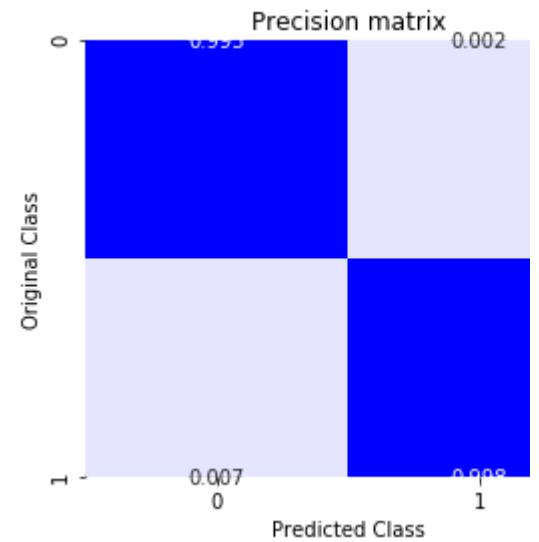
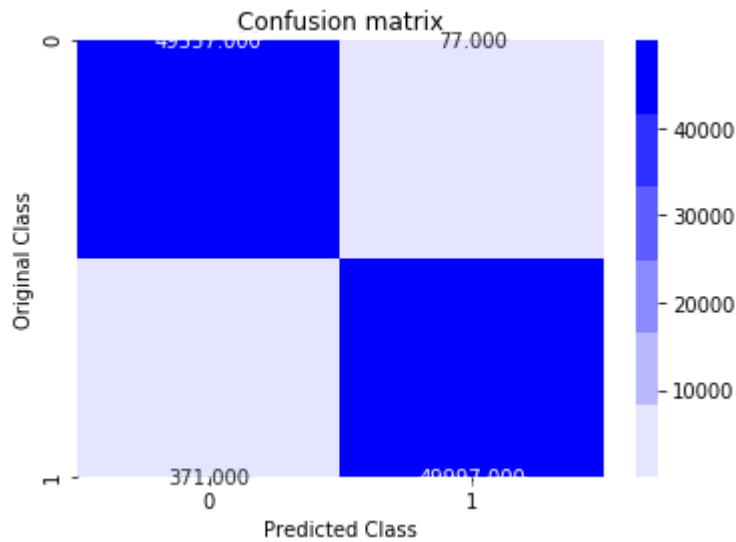
```

print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)

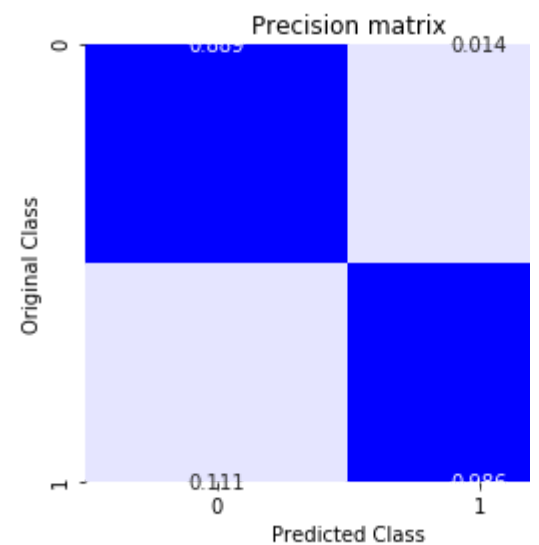
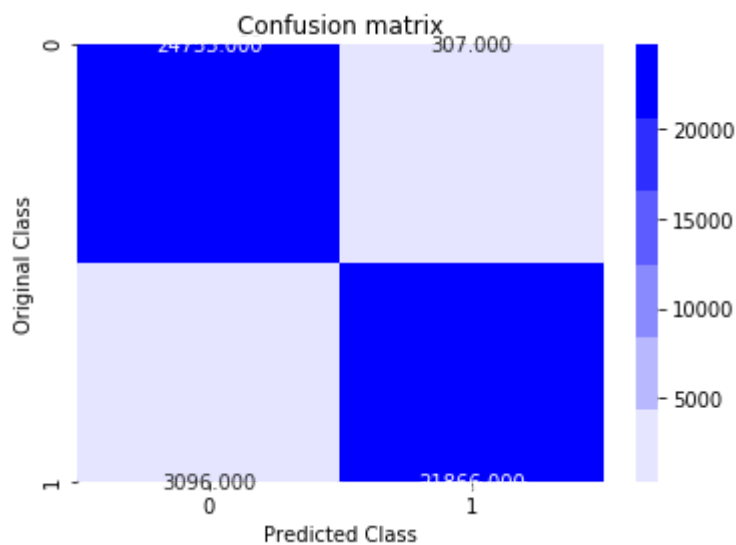
```

```
↳
```

Train confusion_matrix



Test confusion_matrix

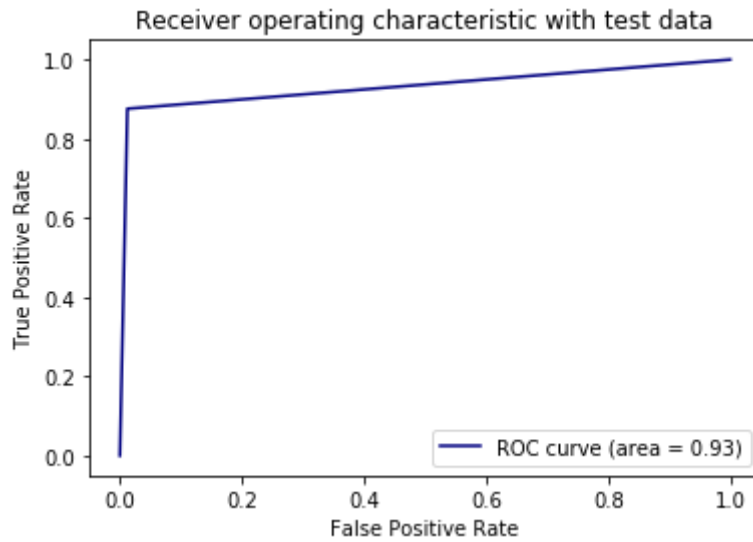


```

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()

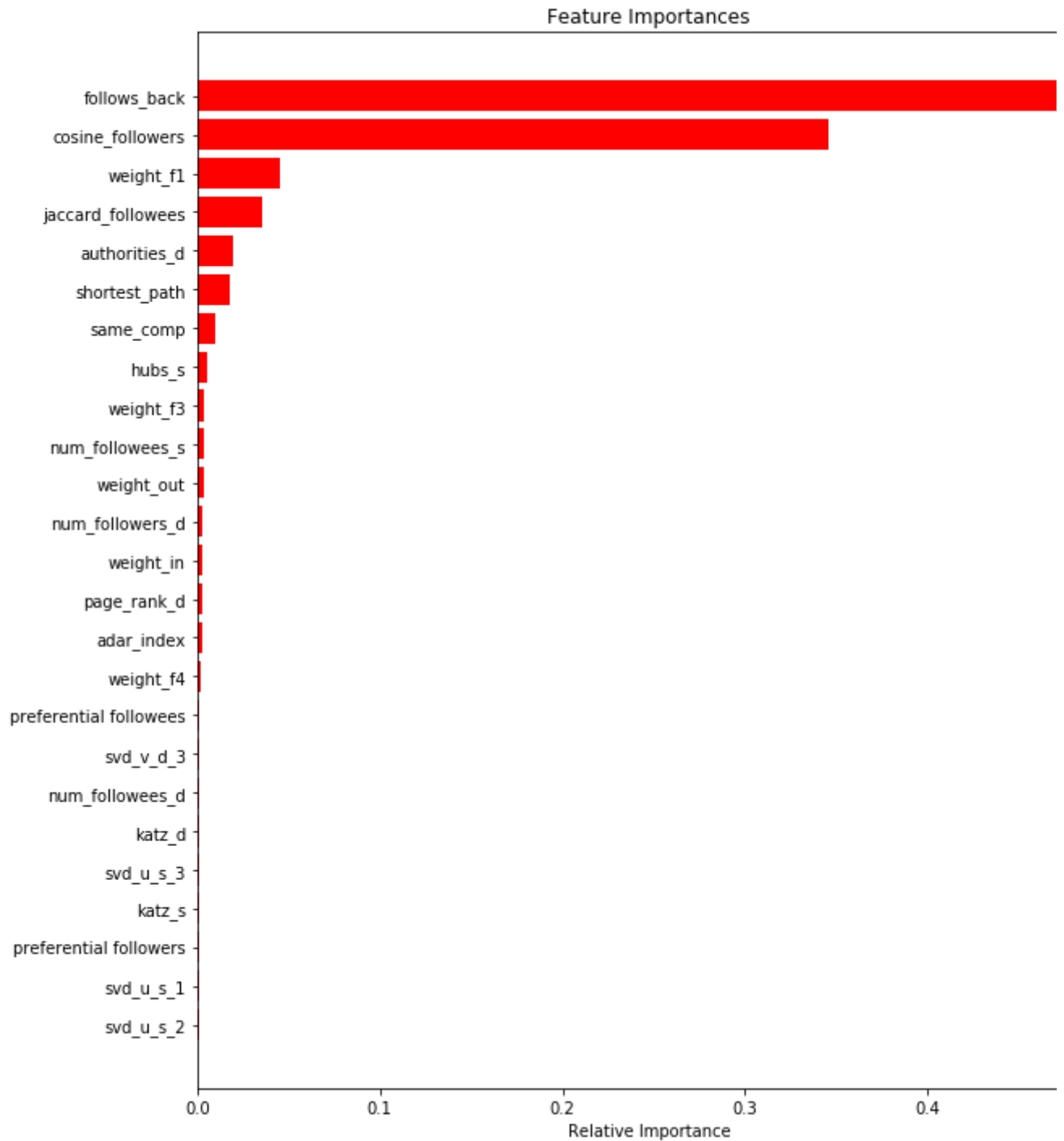
```





```
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()
```





```
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Model", "n_estimators", "max_depth", "Train f1-Score", "Test f1-Score"]
x.add_row(['RFC', '121', '14', '0.965', '0.926'])
x.add_row(['XGB00ST', '112', '11', '0.995', '0.927'])
print(x)
```

Model	n_estimators	max_depth	Train f1-Score	Test f1-Score
RFC	121	14	0.965	0.926
XGBOOST	112	11	0.995	0.927

