

Social network Graph Link Prediction - Facebook Challenge

```
In [3]: #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

In [4]: #reading
from pandas import read_hdf
df_final_train = read_hdf('data/fea_sample/storage_sample_stage14.h5', 'train_df', mode='r')
df_final_test = read_hdf('data/fea_sample/storage_sample_stage14.h5', 'test_df', mode='r')
df_final_train.columns

Out[4]: 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',
'num_followers_d', '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', 'svd_dot_u',
'svd_dot_v', 'Preferential_Attachment_followees',
'Preferential_Attachment_followers'],
dtype='object')

In [5]: y_train = df_final_train.indicator_link
y_test = df_final_test.indicator_link

In [6]: 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)

In [7]: 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))
```

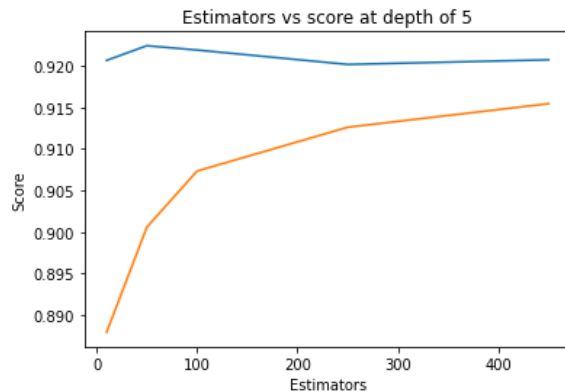
```

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 0.9206418947192234 test Score 0.8879401427431701
 Estimators = 50 Train Score 0.9224029775844764 test Score 0.9005693327931424
 Estimators = 100 Train Score 0.9218875690809748 test Score 0.9073321479976293
 Estimators = 250 Train Score 0.9201547146142588 test Score 0.9126017542752571
 Estimators = 450 Train Score 0.9207118389172028 test Score 0.9154347780503659

Out[7]: Text(0.5, 1.0, 'Estimators vs score at depth of 5')

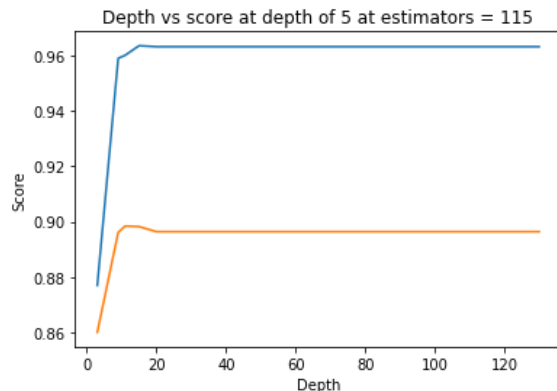


```

In [8]: 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=10, 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 0.8770218558734311 test Score 0.8600957354221062
 depth = 9 Train Score 0.9589063324003083 test Score 0.8961019393443681
 depth = 11 Train Score 0.9600137049801479 test Score 0.8984151394762697
 depth = 15 Train Score 0.9634524013678869 test Score 0.8982105308329399
 depth = 20 Train Score 0.9631118157307456 test Score 0.896384195365235
 depth = 35 Train Score 0.9631118157307456 test Score 0.896384195365235
 depth = 50 Train Score 0.9631118157307456 test Score 0.896384195365235
 depth = 70 Train Score 0.9631118157307456 test Score 0.896384195365235
 depth = 130 Train Score 0.9631118157307456 test Score 0.896384195365235



```
In [9]: 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,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 [0.96176962 0.96168963 0.96028564 0.96142507 0.9630722 ]
mean train scores [0.96281135 0.96218114 0.96054361 0.96222589 0.96414807]
```

```
In [10]: print(rf_random.best_estimator_)

RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                       criterion='gini', max_depth=14, max_features='auto',
                       max_leaf_nodes=None, max_samples=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)
```

```
In [11]: clf = RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                                     criterion='gini', max_depth=14, max_features='auto',
                                     max_leaf_nodes=None, max_samples=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)
```

```
In [12]: clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)
```

```
In [13]: 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 0.964252011648791
Test f1 score 0.9254774789951357

```
In [8]: 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")

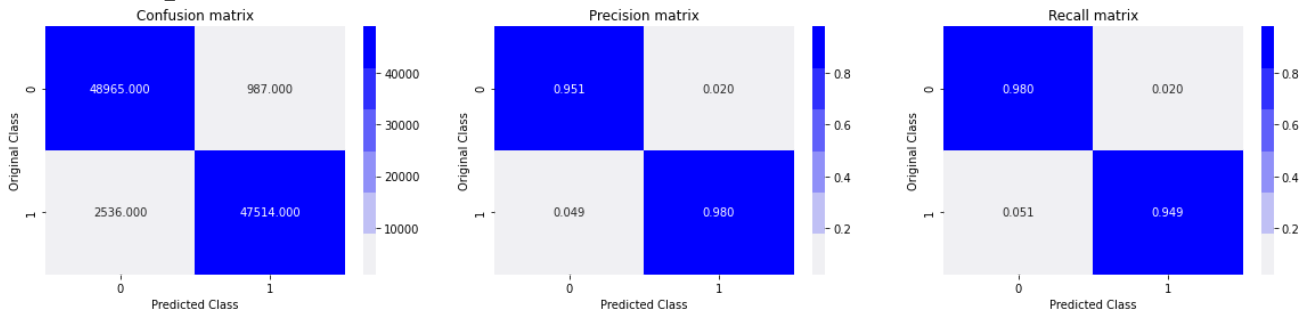
    plt.subplot(1, 3, 2)
    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Precision matrix")

    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")
```

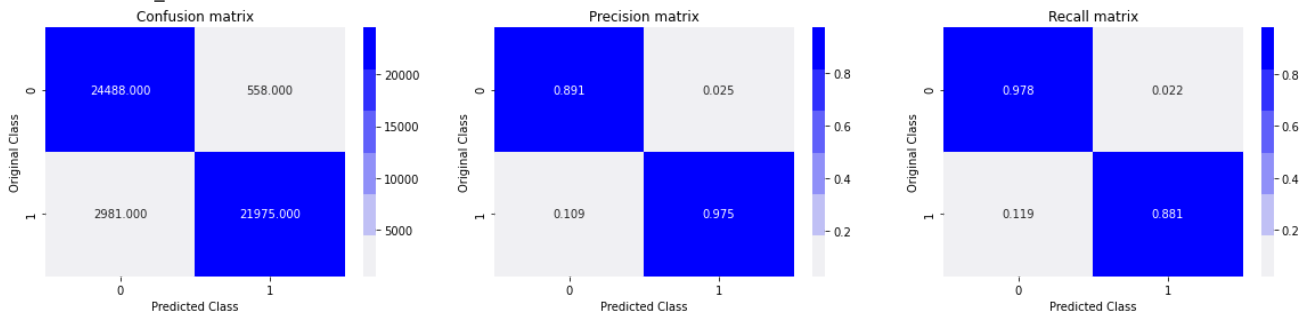
```
plt.show()
```

```
In [15]: print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

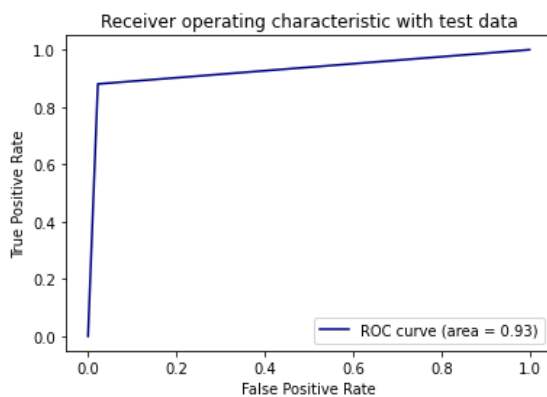
Train confusion_matrix



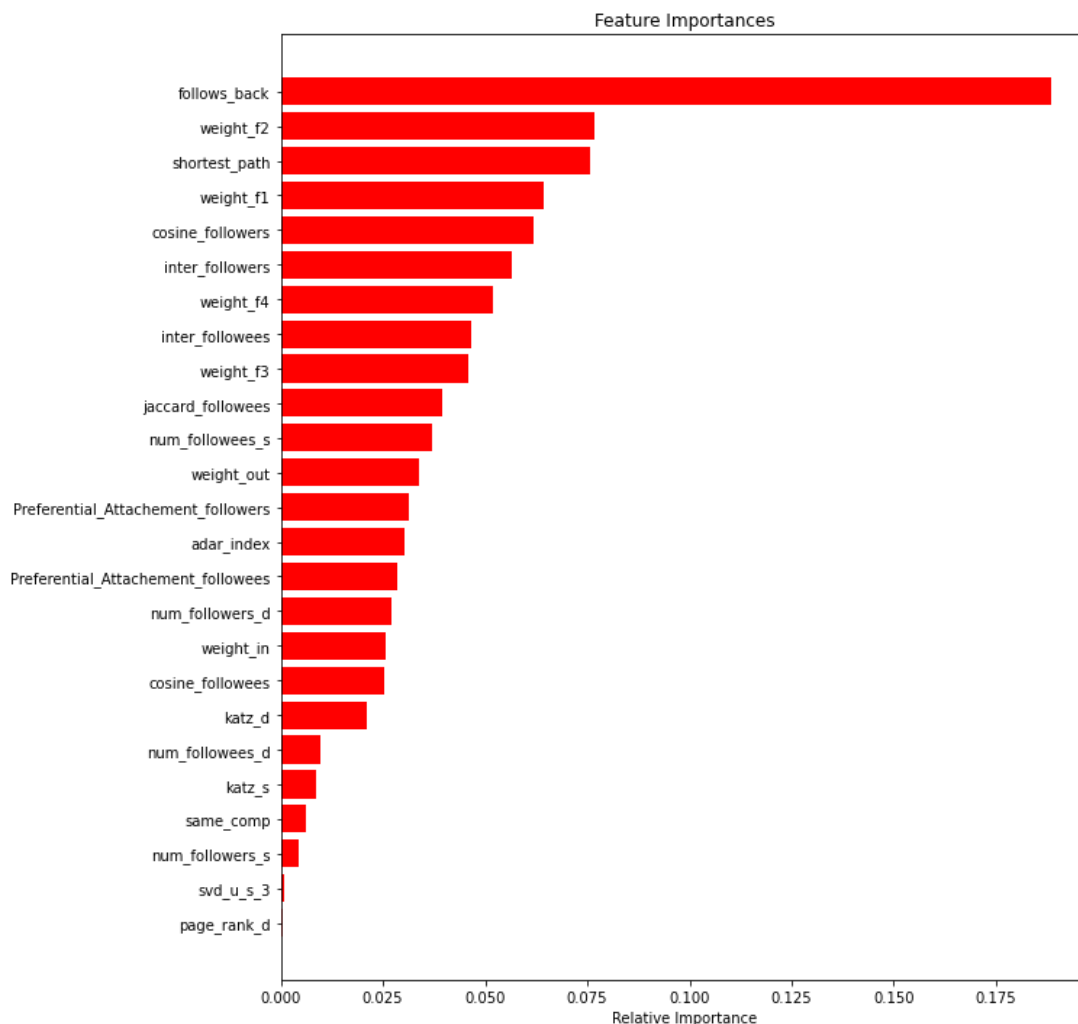
Test confusion_matrix



```
In [16]: 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()
```



```
In [17]: 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()
```



Applying XGBOOST

```
In [10]: import xgboost as xgb
clf_gb = xgb.XGBClassifier(n_jobs=-1)
param_grid = {
    'learning_rate' : [0.0001, 0.001, 0.01, 0.1, 0.2, 0.3],
    'n_estimators' : [5,10,50, 75, 100, 200]}

# Instantiate the grid search model
RandomSearch_xgb = RandomizedSearchCV(clf_gb,param_distributions = param_grid,cv=3,scoring='roc_auc', return_train_score=True)
RandomSearch_xgb.fit(df_final_train,y_train)
```

```
Out[10]: RandomizedSearchCV(cv=3, error_score=nan,
        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,
                                reg_lambda=1, scale_pos_weight=1,
                                seed=None, silent=None, subsample=1,
                                verbosity=1),
        iid='deprecated', n_iter=10, n_jobs=None,
        param_distributions={'learning_rate': [0.0001, 0.001, 0.01,
                                                0.1, 0.2, 0.3],
                             'n_estimators': [5, 10, 50, 75, 100,
                                                200]}},
        pre_dispatch='2*n_jobs', random_state=None, refit=True,
        return_train_score=True, scoring='roc_auc', verbose=0)
```

```
In [11]: f = open('RandomSearch_xgb.pkl', 'wb') # 'wb' instead 'w' for binary file
pickle.dump(RandomSearch_xgb, f, -1) # -1 specifies highest binary protocol
f.close()
```

```
In [12]: f = open('RandomSearch_xgb.pkl', 'rb') # 'rb' for reading binary file
model = pickle.load(f)
f.close()
```

```
In [13]: print('mean test scores',RandomSearch_xgb.cv_results_['mean_test_score'])
print('mean train scores',RandomSearch_xgb.cv_results_['mean_train_score'])

mean test scores [0.9975054  0.92870873 0.9732261  0.93244291 0.99748817 0.99112854
 0.94808495 0.96808614 0.92896493 0.94064139]
mean train scores [0.99776678 0.92887022 0.97340123 0.93283597 0.99771959 0.99116337
 0.94845762 0.96845662 0.92911123 0.94084923]
```

```
In [14]: print(RandomSearch_xgb.best_estimator_)
```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
  colsample_bynode=1, colsample_bytree=1, gamma=0,
  learning_rate=0.2, max_delta_step=0, max_depth=3,
  min_child_weight=1, missing=None, n_estimators=75, 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)
```

```
In [15]: clf_gb = xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
  colsample_bynode=1, colsample_bytree=1, gamma=0,
  learning_rate=0.2, max_delta_step=0, max_depth=3,
  min_child_weight=1, missing=None, n_estimators=75, 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_gb.fit(df_final_train,y_train)
```

```
Out[15]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
  colsample_bynode=1, colsample_bytree=1, gamma=0,
  learning_rate=0.2, max_delta_step=0, max_depth=3,
  min_child_weight=1, missing=None, n_estimators=75, 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)
```

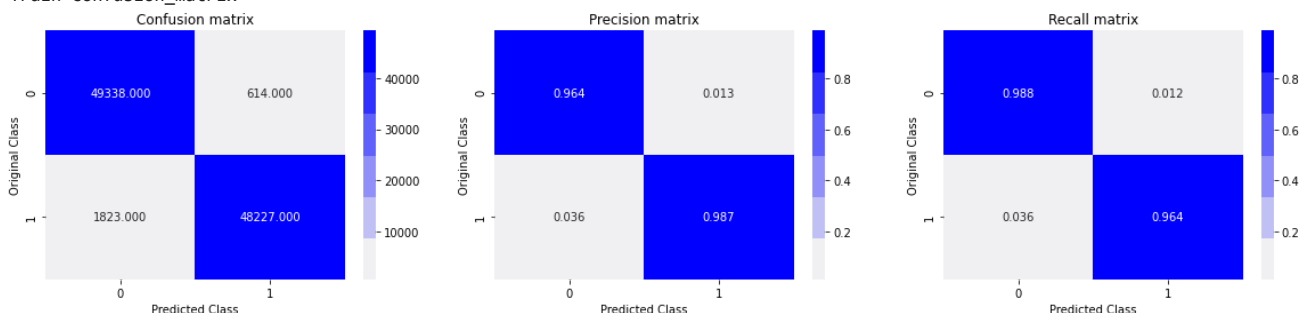
```
In [18]: y_train_pred = clf_gb.predict(df_final_train)
y_test_pred = clf_gb.predict(df_final_test)
```

```
In [19]: 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))
```

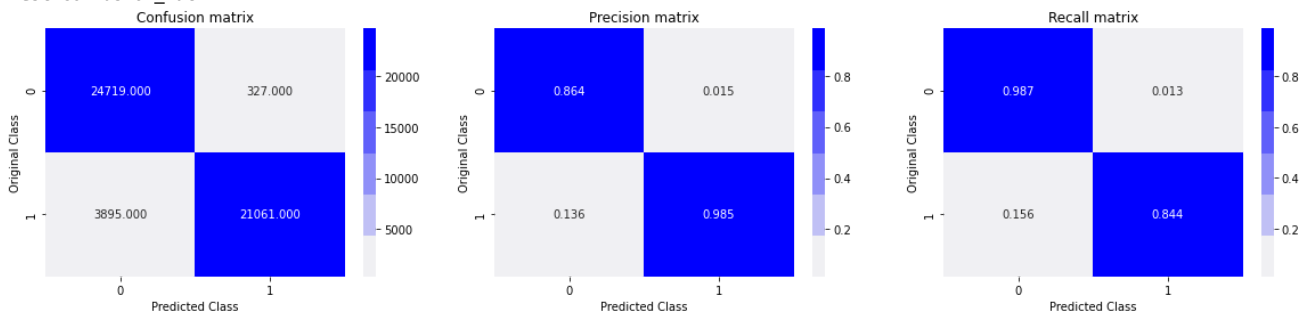
```
Train f1 score 0.975356705868077
Test f1 score 0.9088986708095977
```

```
In [20]: print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

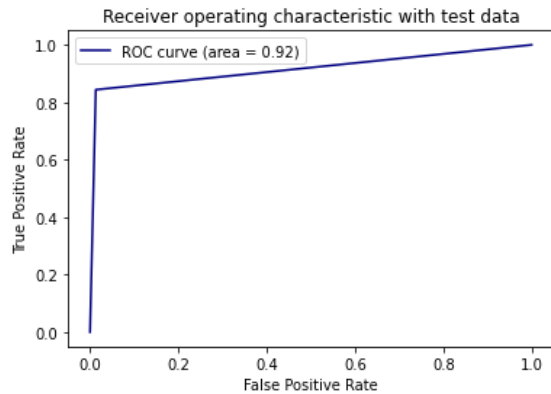
Train confusion_matrix



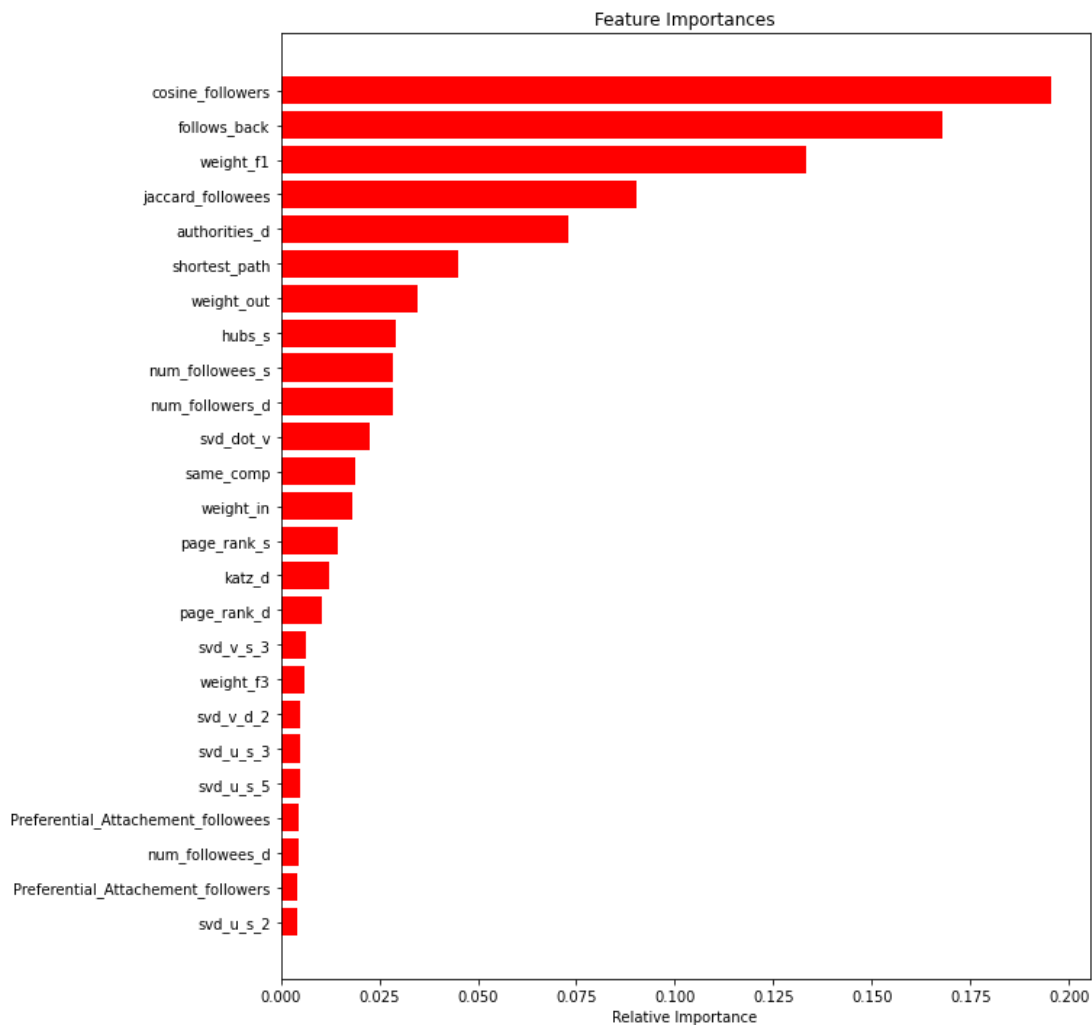
Test confusion_matrix



```
In [21]: 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()
```



```
In [22]: features = df_final_train.columns
importances = clf_gb.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()
```



```
In [ ]: 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
```

```
In [23]: from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Model", "Hyper parameter", "f1-Score"]
x.add_row(['Random Forest', 'max_depth=14, n_estimators = 121', '0.925'])
x.add_row(['XGBOOST', 'learning_rate = 0.2, n_estimators = 75', '0.908'])
print(x)
```

Model	Hyper parameter	f1-Score
Random Forest	max_depth=14, n_estimators = 121	0.925
XGBOOST	learning_rate = 0.2, n_estimators = 75	0.908

1. At first we were given dataset with two column i.e source and destination. Then we perform exploratory data analysis over the dataset.
2. Since we are performing supervised learning , it needs to have at least 2 class label. But there were some edge not present in the graph for classification i.e only edges were present for those node for which there was connection. So, we Generated Bad links from graph which were not in graph and whose shortest path is greater than 2 .
3. Then we split the missing data and given data separately and later we concatenate train and test data separately for feature engineering.
4. Then we performed some feature engineering on dataset like Jaccard Distance followed by Cosine distance , Page Ranking ,Shortest path , Checking for same community , Adamic/Adar Index, Is person was following back,Katz Centrality ,Hits Score, Weight Features, SVD ,Preferential Attachment and svd_dot.
4. After we were done with exploratory data analysis and feature engineering , we perform Random Forest and XGBOOST for which we did RandomizedSearchCV for both of them. We found the best hyperparameter . Then we fit the model with best hyperparameter and found the f1 score for both of them. 6.At last we plotted confusion matrix for train and test data and plotted ROC curve for test data as well as found top 25 features.