# Social network Graph Link Prediction - Facebook Challenge
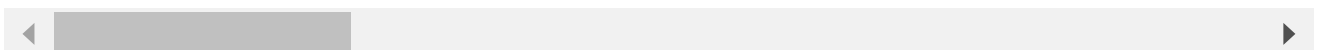
```python
1 #Importing Libraries
2 # please do go through this python notebook:
3 import warnings
4 warnings.filterwarnings("ignore")
5
6 import csv
7 import pandas as pd#pandas to create small dataframes
8 import datetime #Convert to unix time
9 import time #Convert to unix time
10 # if numpy is not installed already : pip3 install numpy
11 import numpy as np#Do aritmetic operations on arrays
12 # matplotlib: used to plot graphs
13 import matplotlib
14 import matplotlib.pylab as plt
15 import seaborn as sns#Plots
16 from matplotlib import rcParams#Size of plots
17 from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
18 import math
19 import pickle
20 import os
21 # to install xgboost: pip3 install xgboost
22 import xgboost as xgb
23 import networkx as nx
24 import pdb
25 import pickle
26 from pandas import HDFStore,DataFrame
27 from pandas import read_hdf
28 from scipy.sparse.linalg import svds, eigs
29 import gc
30 from tqdm import tqdm
31 from sklearn.ensemble import RandomForestClassifier
32 from sklearn.metrics import f1_score
```

```python
1 !wget --header="Host: doc-0o-bk-docs.googleusercontent.com" --header="User-Agent: Mozil
```

```
   --2022-11-28 12:24:19--  https://doc-0o-bk-docs.googleusercontent.com/docs/securesc/
   Resolving doc-0o-bk-docs.googleusercontent.com (doc-0o-bk-docs.googleusercontent.com)
   Connecting to doc-0o-bk-docs.googleusercontent.com (doc-0o-bk-docs.googleusercontent
   HTTP request sent, awaiting response... 403 Forbidden
   2022-11-28 12:24:19 ERROR 403: Forbidden.
```

```python
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
   Mounted at /content/drive
```

```python
1 #reading
```

```
2 from pandas import read_hdf
3 df_final_train = read_hdf('drive/My Drive/Facebook-Recommendation/Data/fea_sample/stora
4 df_final_test = read_hdf('drive/My Drive/Facebook-Recommendation/Data/fea_sample/storag
```

— + Code — — + Text —

```
1 df_final_train.columns
```

```
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_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'],
      dtype='object')
```

```
1 y_train = df_final_train.indicator_link
2 y_test = df_final_test.indicator_link
```

```
1 df_final_train.drop(['source_node', 'destination_node','indicator_link'],axis=1,inplace
2 df_final_test.drop(['source_node', 'destination_node','indicator_link'],axis=1,inplace=
```

```
1 df_final_test.shape
```

```
(50002, 51)
```

```
1 df_final_train.shape
```
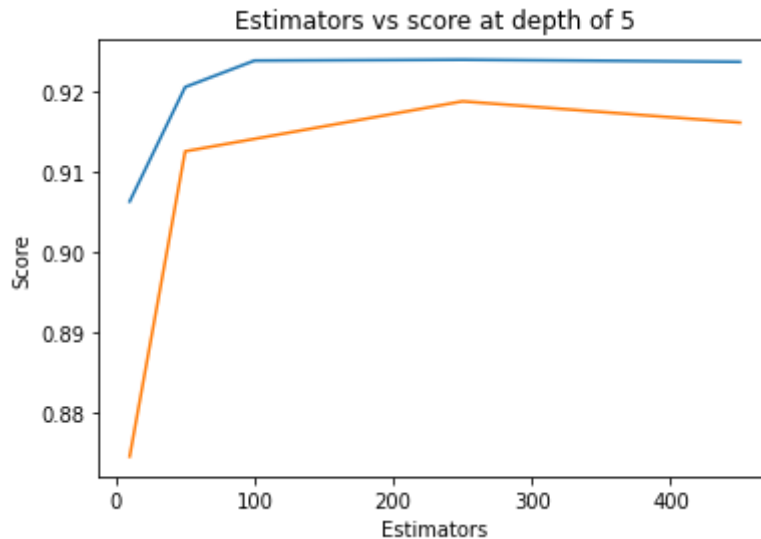
```
(100002, 51)
```

```
 1 estimators = [10,50,100,250,450]
 2 train_scores = []
 3 test_scores = []
 4 for i in estimators:
 5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
 6             max_depth=5, max_features='auto', max_leaf_nodes=None,
 7             min_impurity_decrease=0.0,
 8             min_samples_leaf=52, min_samples_split=120,
 9             min_weight_fraction_leaf=0.0, n_estimators=i, n_jobs=-1,random_state=25,ver
10     clf.fit(df_final_train,y_train)
11     train_sc = f1_score(y_train,clf.predict(df_final_train))
12     test_sc = f1_score(y_test,clf.predict(df_final_test))
13     test_scores.append(test_sc)
14     train_scores.append(train_sc)
15     print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
16 plt.plot(estimators,train_scores,label='Train Score')
17 plt.plot(estimators,test_scores,label='Test Score')
18 plt.xlabel('Estimators')
19 plt.ylabel('Score')
```

```
19 plt.ylabel('Score')
20 plt.title('Estimators vs score at depth of 5')
```

```
Estimators =   10 Train Score 0.9063252121775113 test Score 0.8745605278006858
Estimators =   50 Train Score 0.9205725512208812 test Score 0.9125653355634538
Estimators =  100 Train Score 0.9238690848446947 test Score 0.9141199714153599
Estimators =  250 Train Score 0.9239789348046863 test Score 0.9188007232664732
Estimators =  450 Train Score 0.9237190618658074 test Score 0.9161507685828595
Text(0.5, 1.0, 'Estimators vs score at depth of 5')
```
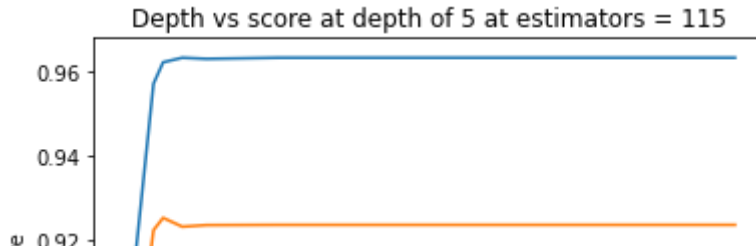


```
 1 depths = [3,9,11,15,20,35,50,70,130]
 2 train_scores = []
 3 test_scores = []
 4 for i in depths:
 5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
 6             max_depth=i, max_features='auto', max_leaf_nodes=None,
 7             min_impurity_decrease=0.0,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

```
depth =   3 Train Score 0.8916120853581238 test Score 0.8687934859875491
depth =   9 Train Score 0.9572226298198419 test Score 0.9222953031452904
depth =  11 Train Score 0.9623451340902863 test Score 0.9252318758281279
depth =  15 Train Score 0.9634267621927706 test Score 0.9231288356496615
depth =  20 Train Score 0.9631629153051491 test Score 0.9235051024711141
depth =  35 Train Score 0.9634333127085721 test Score 0.9235601652753184
depth =  50 Train Score 0.9634333127085721 test Score 0.9235601652753184
depth =  70 Train Score 0.9634333127085721 test Score 0.9235601652753184
depth =  130 Train Score 0.9634333127085721 test Score 0.9235601652753184
```



Depth vs score at depth of 5 at estimators = 115

```
1 from sklearn.metrics import f1_score
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.metrics import f1_score
4 from sklearn.model_selection import RandomizedSearchCV
5 from scipy.stats import randint as sp_randint
6 from scipy.stats import uniform
7
8 param_dist = {"n_estimators":sp_randint(105,125),
9               "max_depth": sp_randint(10,15),
10              "min_samples_split": sp_randint(110,190),
11              "min_samples_leaf": sp_randint(25,65)}
12
13 clf = RandomForestClassifier(random_state=25,n_jobs=-1)
14
15 rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,return_train_score=T
16                                   n_iter=5,cv=10,scoring='f1',random_state=25)
17
18 rf_random.fit(df_final_train,y_train)
19 print('mean test scores',rf_random.cv_results_['mean_test_score'])
20 print('mean train scores',rf_random.cv_results_['mean_train_score'])
```

```
mean test scores [0.96225042 0.96215492 0.9605708  0.96194014 0.96330005]
mean train scores [0.96294922 0.96266735 0.96115674 0.96263457 0.96430539]
```

```
1 print(rf_random.best_estimator_)
```

```
RandomForestClassifier(max_depth=14, min_samples_leaf=28, min_samples_split=111,
                       n_estimators=121, n_jobs=-1, random_state=25)
```

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

```
1 clf.fit(df_final_train,y_train)
```

```
2 y_train_pred = clf.predict(df_final_train)
3 y_test_pred = clf.predict(df_final_test)
```

```
1 from sklearn.metrics import f1_score
2 print('Train f1 score',f1_score(y_train,y_train_pred))
3 print('Test f1 score',f1_score(y_test,y_test_pred))
```
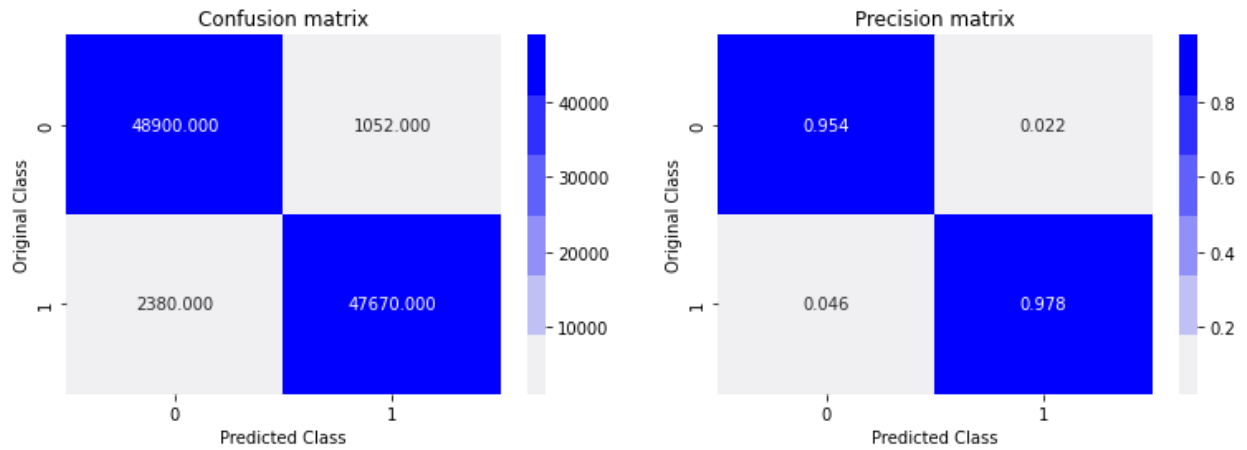
```
    Train f1 score 0.9652533106548414
    Test f1 score 0.9241678239279553
```
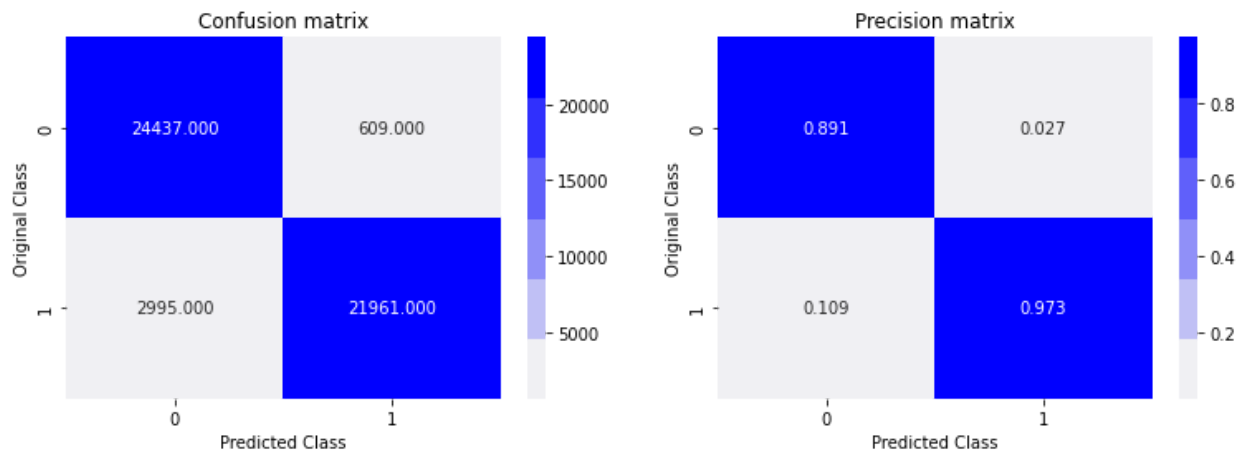
```
 1 from sklearn.metrics import confusion_matrix
 2 def plot_confusion_matrix(test_y, predict_y):
 3     C = confusion_matrix(test_y, predict_y)
 4
 5     A =(((C.T)/(C.sum(axis=1))).T)
 6
 7     B =(C/C.sum(axis=0))
 8     plt.figure(figsize=(20,4))
 9
10     labels = [0,1]
11     # representing A in heatmap format
12     cmap=sns.light_palette("blue")
13     plt.subplot(1, 3, 1)
14     sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
15     plt.xlabel('Predicted Class')
16     plt.ylabel('Original Class')
17     plt.title("Confusion matrix")
18
19     plt.subplot(1, 3, 2)
20     sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
21     plt.xlabel('Predicted Class')
22     plt.ylabel('Original Class')
23     plt.title("Precision matrix")
24
25     plt.subplot(1, 3, 3)
26     # representing B in heatmap format
27     sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
28     plt.xlabel('Predicted Class')
29     plt.ylabel('Original Class')
30     plt.title("Recall matrix")
31
32     plt.show()
```

```
1 print('Train confusion_matrix')
2 plot_confusion_matrix(y_train,y_train_pred)
3 print('Test confusion_matrix')
4 plot_confusion_matrix(y_test,y_test_pred)
```
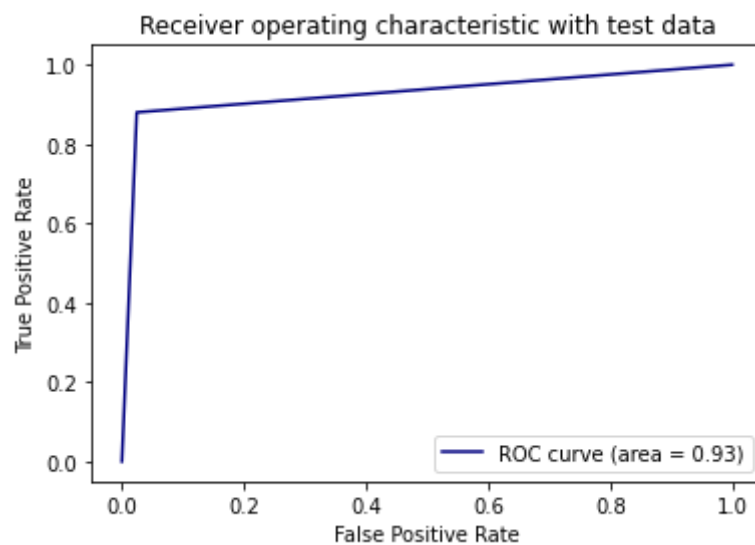
Train confusion_matrix



Test confusion_matrix



```
1 from sklearn.metrics import roc_curve, auc
2 fpr,tpr,ths = roc_curve(y_test,y_test_pred)
3 auc_sc = auc(fpr, tpr)
4 plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('Receiver operating characteristic with test data')
8 plt.legend()
9 plt.show()
```
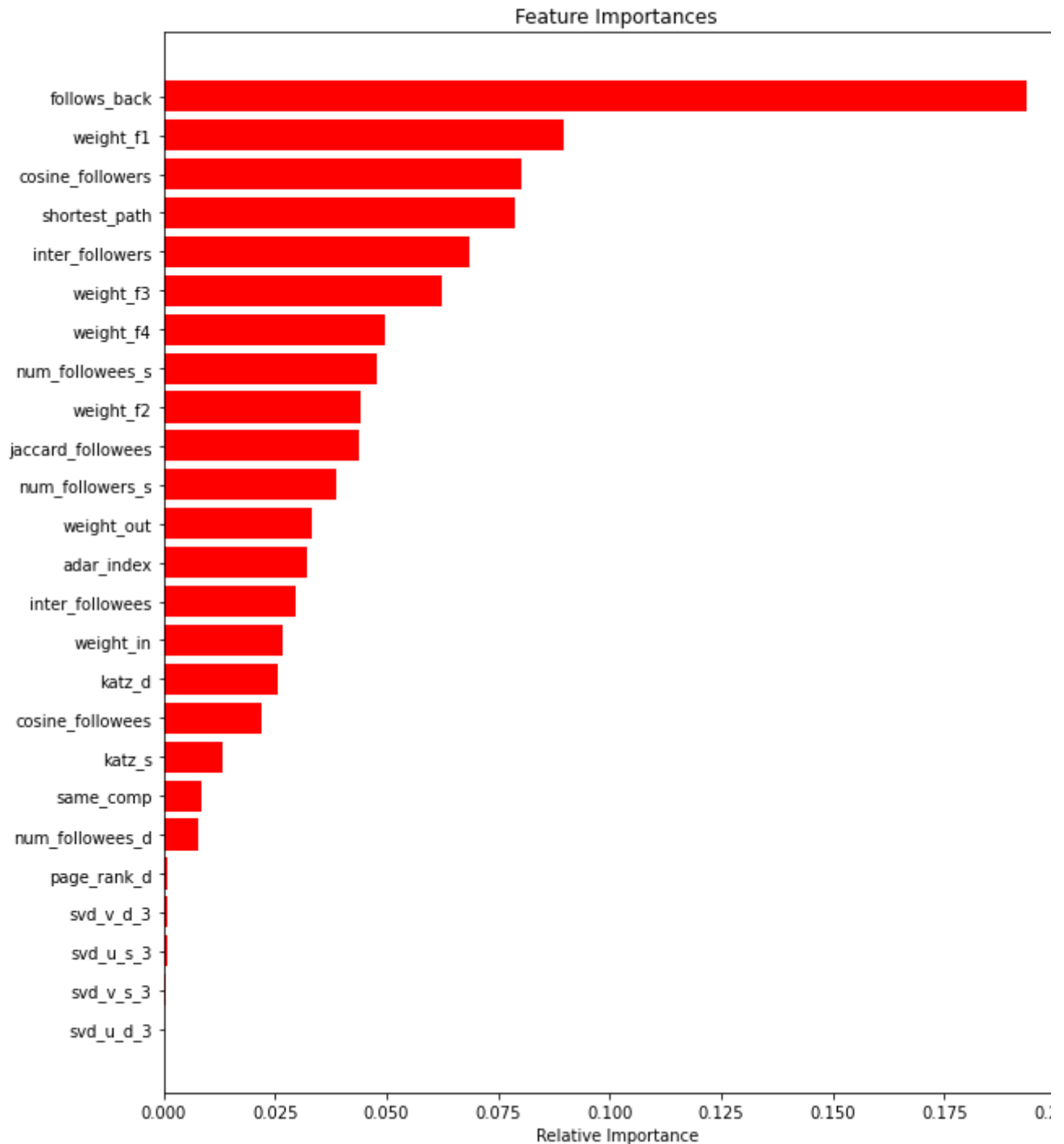


```
1 features = df_final_train.columns
2 importances = clf.feature_importances_
```

```
3 indices = (np.argsort(importances))[-25:]
4 plt.figure(figsize=(10,12))
5 plt.title('Feature Importances')
6 plt.barh(range(len(indices)), importances[indices], color='r', align='center')
7 plt.yticks(range(len(indices)), [features[i] for i in indices])
8 plt.xlabel('Relative Importance')
9 plt.show()
```



## 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.

## 1. Preferential Attachment

https://neo4j.com/docs/graph-data-science/current/alpha-algorithms/preferential-attachment/#:~:text=Preferential%20attachment%20means%20that%20the,of%20nodes%20adjacent%20to%20u%20.

```
Preferential Attachement = |X| * |Y|
```

```
1 #reading
2 from pandas import read_hdf
3 df_final_train = read_hdf('drive/My Drive/Facebook-Recommendation/Data/fea_sample/stora
4 df_final_test = read_hdf('drive/My Drive/Facebook-Recommendation/Data/fea_sample/storag
```

```
1 y_train = df_final_train.indicator_link
2 y_test = df_final_test.indicator_link
```

```
1 y_train.shape
2 #y_test.shape
```

```
    (100002,)
```

```
1 #for followees
2 def pref_att_for_followees(a,b):
3     try:
4         return len(set(train_graph.successors(a))) * len(set(train_graph.successors(b))
5     except:
6         return 0
```

```
1 #for followers
2 def pref_att_for_followers(a,b):
3     try:
4         return len(set(train_graph.predecessors(a)))* len(set(train_graph.predecessors(
5     except:
6         return 0
```

```
1 if not os.path.isfile('data/fea_sample/storage_sample_stage5.h5'):
2     #mapping Preferential Attachment followers to train and test data
```

```
 3     df_final_train['preferential_attachment_followers'] = df_final_train.apply(lambda r
 4                                         pref_att_for_followers(row['source_node'],r
 5     df_final_test['preferential_attachment_followers'] = df_final_test.apply(lambda row
 6                                         pref_att_for_followers(row['source_node'],r
 7
 8     #mapping Preferential Attachment followees to train and test data
 9     df_final_train['preferential_attachment_followees'] = df_final_train.apply(lambda r
10                                         pref_att_for_followees(row['source_node'],r
11     df_final_test['preferential_attachment_followees'] = df_final_test.apply(lambda row
12                                         pref_att_for_followees(row['source_node'],r
13
14     hdf = HDFStore('drive/My Drive/Facebook-Recommendation/Data/fea_sample/storage_samp
15     hdf.put('train_df',df_final_train, format='table', data_columns=True)
16     hdf.put('test_df',df_final_test, format='table', data_columns=True)
17     hdf.close()
18 else:
19     df_final_train = read_hdf('drive/My Drive/Facebook-Recommendation/Data/fea_sample/s
20     df_final_test = read_hdf('drive/My Drive/Facebook-Recommendation/Data/fea_sample/st
```

```
 1 df_final_test.shape
```

```
   (50002, 56)
```

```
 1 df_final_train.shape
```

```
   (100002, 56)
```

```
 1 #for train dataset
 2 us1,us2,us3,us4,us5,us6 = df_final_train['svd_u_s_1'],df_final_train['svd_u_s_2'],df_fi
 3 vs1,vs2,vs3,vs4,vs5,vs6 = df_final_train['svd_v_s_1'],df_final_train['svd_v_s_2'],df_fi
 4
 5 ud1,ud2,ud3,ud4,ud5,ud6 = df_final_train['svd_u_d_1'],df_final_train['svd_u_d_2'],df_fi
 6 vd1,vd2,vd3,vd4,vd5,vd6 = df_final_train['svd_v_d_1'],df_final_train['svd_v_d_2'],df_fi
 7
```

```
 1 svd_dot=[]
 2 for i in range(len(np.array(us1))):
 3     usd=[]
 4     vsd=[]
 5     usd.append(np.array(us1[i]))
 6     usd.append(np.array(us2[i]))
 7     usd.append(np.array(us3[i]))
 8     usd.append(np.array(us4[i]))
 9     usd.append(np.array(us5[i]))
10     usd.append(np.array(us6[i]))
11     usd.append(np.array(vs1[i]))
12     usd.append(np.array(vs2[i]))
13     usd.append(np.array(vs3[i]))
14     usd.append(np.array(vs4[i]))
15     usd.append(np.array(vs5[i]))
16     usd.append(np.array(vs6[i]))
17     vsd.append(np.array(ud1[i]))
18     vsd.append(np.array(ud2[i]))
```

```
19     vsd.append(np.array(ud3[i]))
20     vsd.append(np.array(ud4[i]))
21     vsd.append(np.array(ud5[i]))
22     vsd.append[np.array(ud6[i]))
23     vsd.append(np.array(vd1[i]))
24     vsd.append(np.array(vd2[i]))
25     vsd.append(np.array(vd3[i]))
26     vsd.append(np.array(vd4[i]))
27     vsd.append(np.array(vd5[i]))
28     vsd.append(np.array(vd6[i]))
29     svd_dot.append(np.dot(usd,vsd))
30 df_final_train['svd_dot']=svd_dot
```

```
1 #for test dataset
2 us1,us2,us3,us4,us5,us6 =df_final_test['svd_u_s_1'],df_final_test['svd_u_s_2'],df_final
3 vs1,vs2,vs3,vs4,vs5,vs6 =df_final_test['svd_v_s_1'],df_final_test['svd_v_s_2'],df_final
4
5 ud1,ud2,ud3,ud4,ud5,ud6 =df_final_test['svd_u_d_1'],df_final_test['svd_u_d_2'],df_final
6 vd1,vd2,vd3,vd4,vd5,vd6 =df_final_test['svd_v_d_1'],df_final_test['svd_v_d_2'],df_final
7
```

```
 1 svd_dot=[]
 2 for i in range(len(np.array(us1))):
 3     usd=[]
 4     vsd=[]
 5     usd.append(np.array(us1[i]))
 6     usd.append(np.array(us2[i]))
 7     usd.append(np.array(us3[i]))
 8     usd.append(np.array(us4[i]))
 9     usd.append(np.array(us5[i]))
10     usd.append(np.array(us6[i]))
11     usd.append(np.array(vs1[i]))
12     usd.append(np.array(vs2[i]))
13     usd.append(np.array(vs3[i]))
14     usd.append(np.array(vs4[i]))
15     usd.append(np.array(vs5[i]))
16     usd.append(np.array(vs6[i]))
17     vsd.append(np.array(ud1[i]))
18     vsd.append(np.array(ud2[i]))
19     vsd.append(np.array(ud3[i]))
20     vsd.append(np.array(ud4[i]))
21     vsd.append(np.array(ud5[i]))
22     vsd.append(np.array(ud6[i]))
23     vsd.append(np.array(vd1[i]))
24     vsd.append(np.array(vd2[i]))
25     vsd.append(np.array(vd3[i]))
26     vsd.append(np.array(vd4[i]))
27     vsd.append(np.array(vd5[i]))
28     vsd.append(np.array(vd6[i]))
29     svd_dot.append(np.dot(usd,vsd))
30 df_final_test['svd_dot']=svd_dot
```

```
1 df_final_train.head(2)
```

|   | source_node | destination_node | indicator_link | jaccard_followers | jaccard_followe |
|---|---|---|---|---|---|
| **0** | 273084 | 1505602 | 1 | 0 | 0.0000 |
| **1** | 832016 | 1543415 | 1 | 0 | 0.1871 |

2 rows × 57 columns

```
1 df_final_test.head(2)
```

|   | source_node | destination_node | indicator_link | jaccard_followers | jaccard_followe |
|---|---|---|---|---|---|
| **0** | 848424 | 784690 | 1 | 0 | 0 |
| **1** | 483294 | 1255532 | 1 | 0 | 0 |

2 rows × 57 columns

```
1 df_final_test.columns
```

```
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_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_attachment_followers',
       'preferential_attachment_followees', 'svd_dot'],
      dtype='object')
```

```
1 # Random Forest model after add all features
2 param_dist = {"n_estimators":sp_randint(105,300),
3               "max_depth": sp_randint(10,20),
4               "min_samples_split": sp_randint(100,200),
5               "min_samples_leaf": sp_randint(25,70)}
6
7 clf_rf = RandomForestClassifier(random_state=25,n_jobs=-1)
```

```
 8
 9 rf_random_search = RandomizedSearchCV(clf_rf, param_distributions=param_dist,return_tra
10                                       n_iter=5,cv=10,scoring='f1',random_state=25)
11
12 rf_random_search.fit(df_final_train,y_train)
13 print('mean test scores',rf_random_search.cv_results_['mean_test_score'])
14 print('mean train scores',rf_random_search.cv_results_['mean_train_score'])
15
16 best_params_rdsearch = rf_random_search.best_params_
17 best_score_rf = rf_random_search.best_score_
18
19 print("Best Params from RandomSearchCV for Random Forest Classifier ", best_params_rdse
20 print("Best score",best_score_rf)
```

```
mean test scores [1. 1. 1. 1. 1.]
mean train scores [1. 1. 1. 1. 1.]
Best Params from RandomSearchCV for Random Forest Classifier  {'max_depth': 14, 'min_
Best score 1.0
```
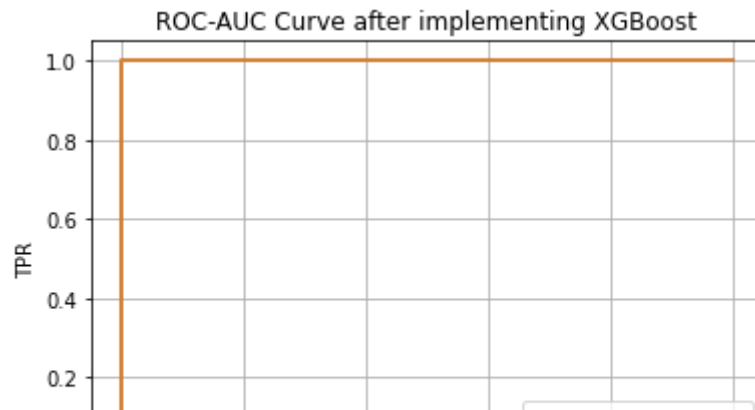
```
 1 from sklearn.metrics import accuracy_score
 2
 3 rdf_clf = RandomForestClassifier(max_depth = 14, min_samples_leaf = 51, min_samples_spl
 4 rdf_clf.fit(df_final_train,y_train)
 5
 6 y_train_predicted_rdf = rdf_clf.predict(df_final_train)
 7 y_test_predicted_rdf = rdf_clf.predict(df_final_test)
 8
 9 rdf_train_fpr, rdf_train_tpr, rdf_train_threshold = roc_curve(y_train, y_train_predicte
10 rdf_test_fpr, rdf_test_tpr, rdf_test_threshold = roc_curve(y_test, y_test_predicted_rdf
11
12 Accuracy_of_Model_rdfc = accuracy_score(y_test,y_test_predicted_rdf)
13 print("Accuracy",Accuracy_of_Model_rdfc)
14
15 # calculate scores
16 train_auc_rdfc = str(auc(rdf_train_fpr, rdf_train_tpr))
17 test_auc_rdfc = str(auc(rdf_test_fpr, rdf_test_tpr))
18
19
20 plt.plot(rdf_train_fpr, rdf_train_tpr, label="Train AUC = "+train_auc_rdfc)
21 plt.plot(rdf_test_fpr, rdf_test_tpr, label="Test AUC = "+test_auc_rdfc)
22
23
24 plt.legend()
25 plt.xlabel('FPR')
26 plt.ylabel('TPR')
27 plt.grid()
28 plt.title('ROC-AUC Curve after implementing XGBoost')
29 plt.show()
```

Accuracy 1.0



```
1 train_f1_score_rf = f1_score(y_train,y_train_predicted_rdf)
2 test_f1_score_rf = f1_score(y_test,y_test_predicted_rdf)
3 print('Train f1 score',train_f1_score_rf)
4 print('Test f1 score',test_f1_score_rf)
```

Train f1 score 1.0
Test f1 score 1.0

```
1 print('Train confusion_matrix')
2 plot_confusion_matrix(y_train,y_train_predicted_rdf)
3 print('Test confusion_matrix')
4 plot_confusion_matrix(y_test,y_test_predicted_rdf)
```

Train confusion_matrix
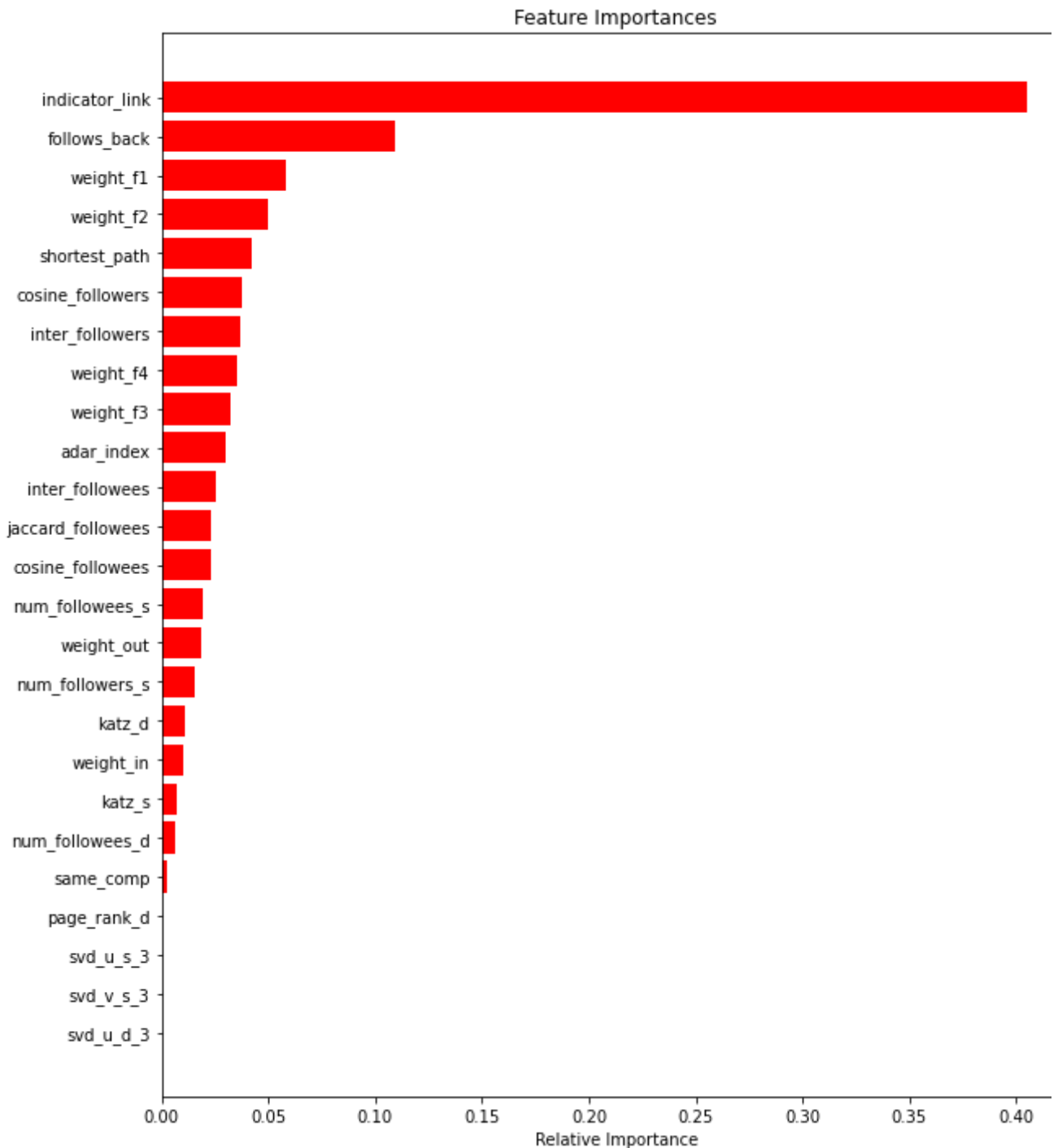


```
1 features = df_final_train.columns
2 importances = rdf_clf.feature_importances_
3 indices = (np.argsort(importances))[-25:]
4 plt.figure(figsize=(10,12))
5 plt.title('Feature Importances')
6 plt.barh(range(len(indices)), importances[indices], color='r', align='center')
7 plt.yticks(range(len(indices)), [features[i] for i in indices])
8 plt.xlabel('Relative Importance')
9 plt.show()
```



```
1 from xgboost.sklearn import XGBClassifier
2
```

```
 3 xgb_clf = XGBClassifier()
 4
 5 params = {"learning_rate" : uniform(0.001,0.3),
 6              "n_estimators" : sp_randint(10,600),
 7              "max_depth"    : sp_randint(5,20)}
 8
 9 rs_xgb = RandomizedSearchCV(xgb_clf, param_distributions=params,return_train_score=True
10                                      n_iter=5,cv=10,scoring='f1',random_state=25)
11
12 rs_xgb.fit(df_final_train,y_train)
13
14 best_params_rndsearch_xgb = rs_xgb.best_params_
15 best_score_xgb = rs_xgb.best_score_
16
17 print('mean test scores',rs_xgb.cv_results_['mean_test_score'])
18 print('mean train scores',rs_xgb.cv_results_['mean_train_score'])
19
20 print("Best Params from RandomSearchCV with XGB ", best_params_rndsearch_xgb)
21 print("Best score",best_score_xgb)
22
23 #Accuracy_of_Model_xgb = accuracy_score(y_test,y_train_pred)
```

```
    mean test scores [1. 1. 1. 1. 1.]
    mean train scores [1. 1. 1. 1. 1.]
    Best Params from RandomSearchCV with XGB  {'learning_rate': 0.26203724098816356, 'max
    Best score 1.0
```
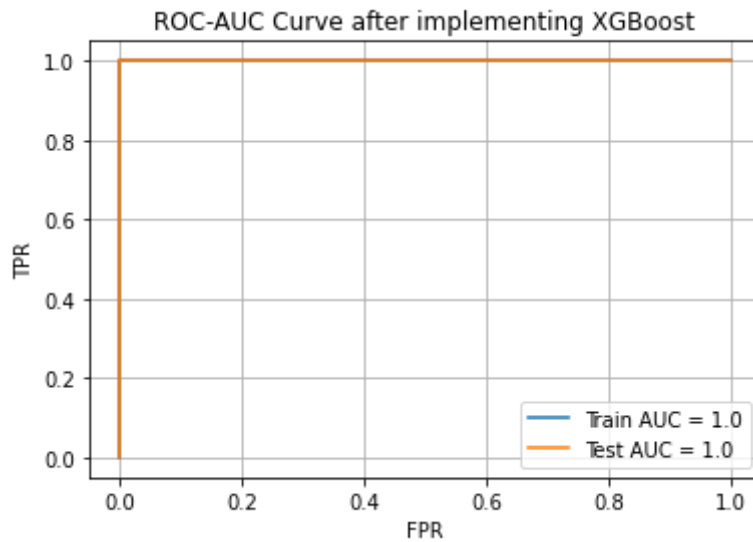
```
 1 print(rs_xgb.best_estimator_)
```

```
    XGBClassifier(learning_rate=0.26203724098816356, max_depth=15, n_estimators=153)
```

```
 1 xgb_clf = XGBClassifier(learning_rate = 0.26203724098816356, max_depth = 15, n_estimato
 2 xgb_clf.fit(df_final_train, y_train)
 3
 4 y_train_predicted_xgb = xgb_clf.predict(df_final_train)
 5 y_test_predicted_xgb = xgb_clf.predict(df_final_test)
 6
 7 xgb_train_fpr, xgb_train_tpr, xgb_train_threshold = roc_curve(y_train, y_train_predicte
 8 xgb_test_fpr, xgb_test_tpr, xgb_test_threshold = roc_curve(y_test, y_test_predicted_xgb
 9
10 Accuracy_of_Model_xgb = accuracy_score(y_test,y_test_predicted_xgb)
11 print("Accuracy",Accuracy_of_Model_xgb)
12
13 # calculate scores
14 train_auc_xgb = str(auc(xgb_train_fpr, xgb_train_tpr))
15 test_auc_xgb = str(auc(xgb_test_fpr, xgb_test_tpr))
16
17
18 plt.plot(xgb_train_fpr, xgb_train_tpr, label="Train AUC = "+train_auc_xgb)
19 plt.plot(xgb_test_fpr, xgb_test_tpr, label="Test AUC = "+test_auc_xgb)
20
21
```

```
22 plt.legend()
23 plt.xlabel('FPR')
24 plt.ylabel('TPR')
25 plt.grid()
26 plt.title('ROC-AUC Curve after implementing XGBoost')
27 plt.show()
```
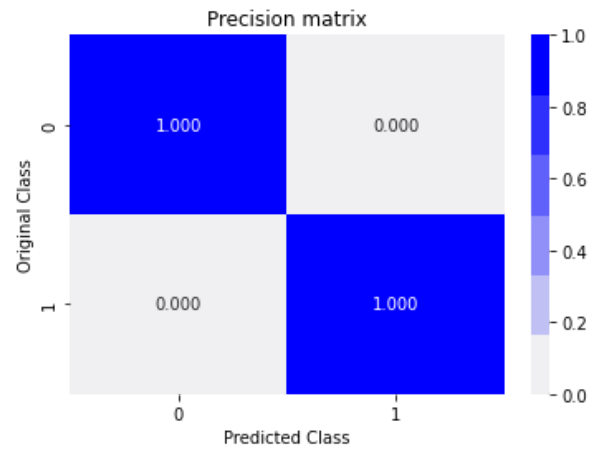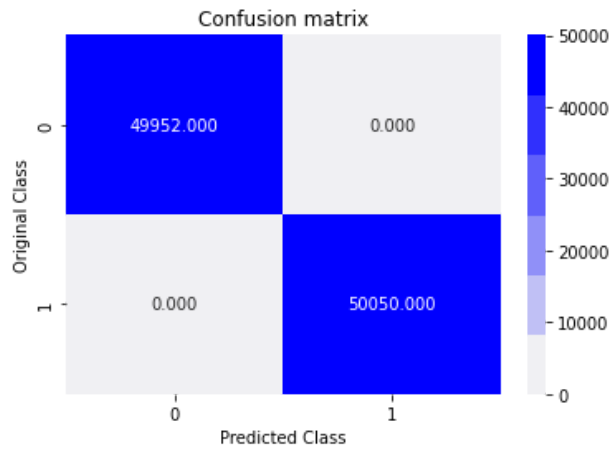
Accuracy 1.0



```
1 train_f1_score_xgb = f1_score(y_train,y_train_predicted_xgb)
2 test_f1_score_xgb = f1_score(y_test,y_test_predicted_xgb)
3 print('Train f1 score',train_f1_score_xgb)
4 print('Test f1 score',test_f1_score_xgb)
```

Train f1 score 1.0
Test f1 score 1.0

```
1 print('Train confusion_matrix')
2 plot_confusion_matrix(y_train,y_train_predicted_xgb)
3 print('Test confusion_matrix')
4 plot_confusion_matrix(y_test,y_test_predicted_xgb)
```
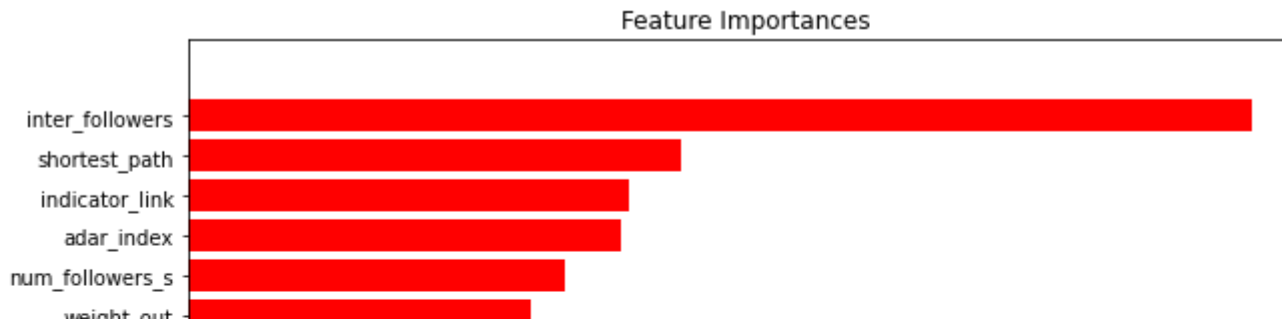
Train confusion_matrix



Test confusion_matrix



```
1 features = df_final_train.columns
2 importances = clf.feature_importances_
3 indices = (np.argsort(importances))[-30:]
4 plt.figure(figsize=(10,12))
5 plt.title('Feature Importances')
6 plt.barh(range(len(indices)), importances[indices], color='r', align='center')
7 plt.yticks(range(len(indices)), [features[i] for i in indices])
8 plt.xlabel('Relative Importance')
9 plt.show()
```

### Feature Importances



```
1 pretty_table = pd.DataFrame(columns = ['Model','Best-Hyper-parameter','Train_f1_score',
2 pretty_table['Model'] = ["Random Forest","XGB"]
3 pretty_table['Best-Hyper-parameter'] = [best_params_rdsearch,best_params_rndsearch_xgb]
4 pretty_table['Train_f1_score'] = [train_f1_score_rf,train_f1_score_xgb]
5 pretty_table['Test_f1_score'] = [test_f1_score_rf,test_f1_score_xgb]
6 pretty_table['Train-AUC'] = [train_auc_rdfc,train_auc_xgb]
7 pretty_table['Test-AUC'] = [test_auc_rdfc,test_auc_xgb]
8 pretty_table['Accuracy'] = [Accuracy_of_Model_rdfc,Accuracy_of_Model_xgb]
9 pretty_table
```

| | Model | Best-Hyper-parameter | Train_f1_score | Test_f1 |
|---|---|---|---|---|
| 0 | Random Forest | {'max_depth': 14, 'min_samples_leaf': 51, 'min... | 1.0 | |
| 1 | XGB | {'learning_rate': 0.26203724098816356, 'max_de... | 1.0 | |

**Observations**

Feature engineering on the dataset like finding the shortest path, Katz centrality, Jaccard distances, page rank, preferential attachments,etc were helped to improve machine learning model training, leading to better performance and greater accuracy.It shows that the available data source was also good enough.

Applying SVD Matrix Factorization helps to incorporate implicit feedback information that is not directly given but can be derived from analyzing user behavior.This will helps our model robust for both seen and unseen data also.

At the end we have plotted the confusion matrix and pretty-table for both XGBoost and Random Forest algorithms found the best hyperparameters we got best accuracy and F1 scores for both training and testing data.

Colab paid products  -  Cancel contracts here