

In []:

```
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 arithmetic operations on arrays
12 # matplotlib: used to plot graphs
13 import matplotlib
14 import matplotlib.pyplot 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 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
24
25 import warnings
26 import networkx as nx
27 import pdb
28 import pickle
29 from pandas import HDFStore, DataFrame
30 from pandas import read_hdf
31 from scipy.sparse.linalg import svds, eigs
32 import gc
33 from tqdm import tqdm
34 from sklearn.ensemble import RandomForestClassifier
35 from sklearn.metrics import f1_score
```

In []:

```
1 #Getting after_eda data
2 !wget --header="Host: doc-0s-80-drive-data-export.googleusercontent.com" --header="User-Agent: curl/7.28.0"
```

```
--2020-09-02 06:55:47-- https://doc-0s-80-drive-data-export.googleusercontent.com/download/sghr8pnsjffiha0mq8imgl110723ufn6/dk0plg57034lp2684omnrvf1ubebje3k/1599025500000/51fbc4b6-29b4-4fdb-b0b9-717a0a8c3a63/100968300222125116581/ADt3v-NWmVn6oYzxH1_stsqDLgfb8XWEyQ8q9lRpqBnNhva--067JxM8CHaiBvcHZx71smjTSa4ItbtEmCkmFL3btGACtiYDR8r2TSEHI2YVKN4XEWVeciepyGzQWRN_YYkwtMQiE2rm2dQqez1GYTsYk40ShWVjtyVeNWKHmNRBlbgbZv_9a33ayKvsevQeehaFuP3swBS7vtKjUohOFq1Y7_C0kY1_C9TAebJpHUjFpsYLAjqwYY6UEzGxPn51HI3oS04vFvMRwQQ1SHG5GfI_0tGDYRk8EbMh8HZYUyQ19irmTZ-BvjRVzfQBwcawQnKfvdFWQrd?authuser=0&nonce=b34jo1tr99bjq&user=100968300222125116581&hash=4mvd5ln84an7rk81gn71nl3fp7biloin (https://doc-0s-80-drive-data-export.googleusercontent.com/download/sghr8pnsjffiha0mq8imgl110723ufn6/dk0plg57034lp2684omnrvf1ubebje3k/1599025500000/51fbc4b6-29b4-4fdb-b0b9-717a0a8c3a63/100968300222125116581/ADt3v-NWmVn6oYzxH1_stsqDLgfb8XWEyQ8q9lRpqBnNhva--067JxM8CHaiBvcHZx71smjTSa4ItbtEmCkmFL3btGACtiYDR8r2TSEHI2YVKN4XEWVeciepyGzQWRN_YYkwtMQiE2rm2dQqez1GYTsYk40ShWVjtyVeNWKHmNRBlbgbZv_9a33ayKvsevQeehaFuP3swBS7vtKjUohOFq1Y7_C0kY1_C9TAebJpHUjFpsYLAjqwYY6UEzGxPn51HI3oS04vFvMRwQQ1SHG5GfI_0tGDYRk8EbMh8HZYUyQ19irmTZ-BvjRVzfQBwcawQnKfvdFWQrd?authuser=0&nonce=b34jo1tr99bjq&user=100968300222125116581&hash=4mvd5ln84an7rk81gn71nl3fp7biloin)
Resolving doc-0s-80-drive-data-export.googleusercontent.com (doc-0s-80-drive-data-export.googleusercontent.com)... 172.217.203.132, 2607:f8b0:400c:c07::84
Connecting to doc-0s-80-drive-data-export.googleusercontent.com (doc-0s-80-drive-data-export.googleusercontent.com)|172.217.203.132|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 416199167 (397M) [application/octet-stream]
Saving to: 'drive-download-20200902T064408Z-001.zip'

drive-download-2020 100%[=====>] 396.92M  125MB/s   in 3.2s

2020-09-02 06:55:52 (125 MB/s) - 'drive-download-20200902T064408Z-001.zip' saved [416199167/416199167]
```

In []:

```
1 !unzip 'drive-download-20200902T064408Z-001.zip'
```

```
Archive: drive-download-20200902T064408Z-001.zip
  inflating: test_neg_after_eda.csv
  inflating: test_pos_after_eda.csv
  inflating: train_pos_after_eda.csv
  inflating: train_woheader.csv
  inflating: train_neg_after_eda.csv
  inflating: missing_edges_final.p
  inflating: train_after_eda.csv
  inflating: test_after_eda.csv
```

In []:

```
1 #Getting fea_sample data
2 !wget --header="Host: doc-3s-20-drive-data-export.googleusercontent.com" --header="User-Agent: curl/7.28.0" https://doc-3s-20-drive-data-export.googleusercontent.com/download/sghr8pnsjffiha0mq8imgl11o723ufn6/hr17r51t9k3qdbe5h7gs8h09gp6lep10/1599025500000/919e3811-e419-4176-824e-e545b09a3b53/100968300222125116581/ADt3v-Pw9RoFjhn6TfXjSqSOQDZ0-j1Wz0g4RbiO-PNfsajPiPhuQ7g8apoM9UEVsVfuyW507azeiJ_Qf1AP60lBIqY-pHLL5vI2ato01eoBysWGUpFMbR0xtityAiPJ-rYDEmNDdDW0Q1A4UPKw3D_Z9WVDuYFmsuDObOYy7r-BjZu9QdWTNtrOLNfoXL04yC8cufxzkeg98RB0W8pNvweOhtKFNkx-T14veyrN0v9fPnkXA0KhG2v6EOHI99qrPS054I1JcvOZZ4Z3gGuk1nmjtuNVp14p9bml0rBDV5X60PzCG6bHFnOFE_3WqTQeRfH3eKNwaJ2h?authuser=0&nonce=p0gl2as7auk2e&user=100968300222125116581&hash=n4spur6l1sn86ghmi6tc96a8mc81h60 (https://doc-3s-20-drive-data-export.googleusercontent.com/download/sghr8pnsjffiha0mq8imgl11o723ufn6/hr17r51t9k3qdbe5h7gs8h09gp6lep10/1599025500000/919e3811-e419-4176-824e-e545b09a3b53/100968300222125116581/ADt3v-Pw9RoFjhn6TfXjSqSOQDZ0-j1Wz0g4RbiO-PNfsajPiPhuQ7g8apoM9UEVsVfuyW507azeiJ_Qf1AP60lBIqY-pHLL5vI2ato01eoBysWGUpFMbR0xtityAiPJ-rYDEmNDdDW0Q1A4UPKw3D_Z9WVDuYFmsuDObOYy7r-BjZu9QdWTNtrOLNfoXL04yC8cufxzkeg98RB0W8pNvweOhtKFNkx-T14veyrN0v9fPnkXA0KhG2v6EOHI99qrPS054I1JcvOZZ4Z3gGuk1nmjtuNVp14p9bml0rBDV5X60PzCG6bHFnOFE_3WqTQeRfH3eKNwaJ2h?authuser=0&nonce=p0gl2as7auk2e&user=100968300222125116581&hash=n4spur6l1sn86ghmi6tc96a8mc81h60)
```

```
--2020-09-02 06:42:49-- https://doc-3s-20-drive-data-export.googleusercontent.com/download/sghr8pnsjffiha0mq8imgl11o723ufn6/hr17r51t9k3qdbe5h7gs8h09gp6lep10/1599025500000/919e3811-e419-4176-824e-e545b09a3b53/100968300222125116581/ADt3v-Pw9RoFjhn6TfXjSqSOQDZ0-j1Wz0g4RbiO-PNfsajPiPhuQ7g8apoM9UEVsVfuyW507azeiJ_Qf1AP60lBIqY-pHLL5vI2ato01eoBysWGUpFMbR0xtityAiPJ-rYDEmNDdDW0Q1A4UPKw3D_Z9WVDuYFmsuDObOYy7r-BjZu9QdWTNtrOLNfoXL04yC8cufxzkeg98RB0W8pNvweOhtKFNkx-T14veyrN0v9fPnkXA0KhG2v6EOHI99qrPS054I1JcvOZZ4Z3gGuk1nmjtuNVp14p9bml0rBDV5X60PzCG6bHFnOFE_3WqTQeRfH3eKNwaJ2h?authuser=0&nonce=p0gl2as7auk2e&user=100968300222125116581&hash=n4spur6l1sn86ghmi6tc96a8mc81h60 (https://doc-3s-20-drive-data-export.googleusercontent.com/download/sghr8pnsjffiha0mq8imgl11o723ufn6/hr17r51t9k3qdbe5h7gs8h09gp6lep10/1599025500000/919e3811-e419-4176-824e-e545b09a3b53/100968300222125116581/ADt3v-Pw9RoFjhn6TfXjSqSOQDZ0-j1Wz0g4RbiO-PNfsajPiPhuQ7g8apoM9UEVsVfuyW507azeiJ_Qf1AP60lBIqY-pHLL5vI2ato01eoBysWGUpFMbR0xtityAiPJ-rYDEmNDdDW0Q1A4UPKw3D_Z9WVDuYFmsuDObOYy7r-BjZu9QdWTNtrOLNfoXL04yC8cufxzkeg98RB0W8pNvweOhtKFNkx-T14veyrN0v9fPnkXA0KhG2v6EOHI99qrPS054I1JcvOZZ4Z3gGuk1nmjtuNVp14p9bml0rBDV5X60PzCG6bHFnOFE_3WqTQeRfH3eKNwaJ2h?authuser=0&nonce=p0gl2as7auk2e&user=100968300222125116581&hash=n4spur6l1sn86ghmi6tc96a8mc81h60)
```

```
Resolving doc-3s-20-drive-data-export.googleusercontent.com (doc-3s-20-drive-data-export.googleusercontent.com)... 172.217.203.132, 2607:f8b0:400c:c07::84
```

```
Connecting to doc-3s-20-drive-data-export.googleusercontent.com (doc-3s-20-drive-data-export.googleusercontent.com)|172.217.203.132|:443... connected.
```

```
HTTP request sent, awaiting response... 200 OK
```

```
Length: 178254609 (170M) [application/octet-stream]
```

```
Saving to: 'drive-download-20200902T063450Z-001.zip'
```

```
drive-download-2020 100%[=====>] 170.00M 90.6MB/s in 1.9s
```

```
2020-09-02 06:42:51 (90.6 MB/s) - 'drive-download-20200902T063450Z-001.zip' saved [178254609/178254609]
```

In []:

```
1 !unzip 'drive-download-20200902T063450Z-001.zip'
```

```
Archive: drive-download-20200902T063450Z-001.zip
```

```
 inflating: storage_sample_stage2.h5
```

```
 inflating: storage_sample_stage1.h5
```

```
 inflating: storage_sample_stage3.h5
```

```
 inflating: hits.p
```

```
 inflating: storage_sample_stage4.h5
```

```
 inflating: katz.p
```

```
 inflating: page_rank.p
```

In []:

```
1 import pandas as pd
2 train_pos = pd.read_csv('train_pos_after_eda.csv', header=None)
3 train_pos.head()
```

Out[10]:

	0	1
0	273084	1505602
1	912810	1678443
2	365429	1523458
3	527014	1605979
4	1228116	471233

In []:

```
1 len(train_pos)
```

Out[11]:

7550015

In []:

```
1 df_final_train = read_hdf('storage_sample_stage4.h5', 'train_df', mode='r')
2 df_final_test = read_hdf('storage_sample_stage4.h5', 'test_df', mode='r')
```

In []:

```
1 print(len(df_final_train))
2 len(df_final_test)
```

100002

Out[15]:

50002

In []:

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

Out[16]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_fc
0	273084	1505602	1	0	0.000000	C
1	832016	1543415	1	0	0.187135	C
2	1325247	760242	1	0	0.369565	C
3	1368400	1006992	1	0	0.000000	C
4	140165	1708748	1	0	0.000000	C

In []:

```
1 #Creating Graph from data_new
2 train_graph=nx.read_edgelist('train_pos_after_eda.csv',delimiter=',',create_using=nx.G
3 print(nx.info(train_graph))
```

Name:

Type: Graph

Number of nodes: 1780722

Number of edges: 5457004

Average degree: 6.1290

###Preferential Attachment

In []:

```
1 #Preferential Attachment Function
2 def preferential_attachment_(a,b):
3     if train_graph.has_node(a) and train_graph.has_node(b):
4         preds = nx.preferential_attachment(train_graph, [(a,b)])
5     else:
6         preds = [(a,b,-1)]
7     for u, v, p in preds:
8         sc = p
9     # print('%d, %d) -> %d' % (u, v, p))
10    return sc
```

In []:

```
1 preferential_attachment_(273084,15) #Node not present in Graph
```

Out[30]:

-1

In []:

```
1 df_final_train['preferential_attachment'] = df_final_train.apply(lambda row:
2     preferential_attachment_(row['source_node']
3 df_final_test['preferential_attachment'] = df_final_test.apply(lambda row:
4     preferential_attachment_(row['source_node']
```

###Svd Dot between svd's of Source and Dest.

In []:

```
1 #SVD_DOT Function
2 def svd_dot_(df_svd):
3     svd_dot_u = 0
4     pt = 0
5     for i in range(6):
6         u_s = df_svd[pt] #Values of svd_u_s
7         u_d = df_svd[pt+6] #Values of svd_u_d
8         svd_dot_u = svd_dot_u + (u_s+u_d) #Elementwise Mul and Sum : (svd_u_s_i*svd_u_d_i
9         pt = pt+1
10    #print(svd_dot_u)
11    return svd_dot_u
```

In []:

```
1 df_final_train['svd_dot_u'] = df_final_train.apply(lambda row:
2     svd_dot_(row[['svd_u_s_1','svd_u_s_2','svd_
3     'svd_u_d_1','svd_u_d_2','svd_
```

In []:

```
1 df_final_train['svd_dot_v'] = df_final_train.apply(lambda row:
2     svd_dot_(row[['svd_v_s_1','svd_v_s_2','svd_
3     'svd_v_d_1','svd_v_d_2','svd_
```

In []:

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

Out[68]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_fr
0	273084	1505602	1	0	0.000000	C
1	832016	1543415	1	0	0.187135	C
2	1325247	760242	1	0	0.369565	C
3	1368400	1006992	1	0	0.000000	C
4	140165	1708748	1	0	0.000000	C

In []:

```
1 df_final_test['svd_dot_u'] = tqdm(df_final_test.apply(lambda row:
2                                     svd_dot_(row[['svd_u_s_1','svd_u_s_2','svd_
3                                     'svd_u_d_1','svd_u_d_2','svd_
4 df_final_test['svd_dot_v'] = tqdm(df_final_test.apply(lambda row:
5                                     svd_dot_(row[['svd_v_s_1','svd_v_s_2','svd_
6                                     'svd_v_d_1','svd_v_d_2','svd_
```

```
100%|██████████| 50002/50002 [00:00<00:00, 2101291.38it/s]
100%|██████████| 50002/50002 [00:00<00:00, 2191308.77it/s]
```

In []:

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

Out[70]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_fr
0	848424	784690	1	0	0.0	C
1	483294	1255532	1	0	0.0	C
2	626190	1729265	1	0	0.0	C
3	947219	425228	1	0	0.0	C
4	991374	975044	1	0	0.2	C

###Implementation for XGBoost Starts Here

In []:

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

In []:

```
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=
```

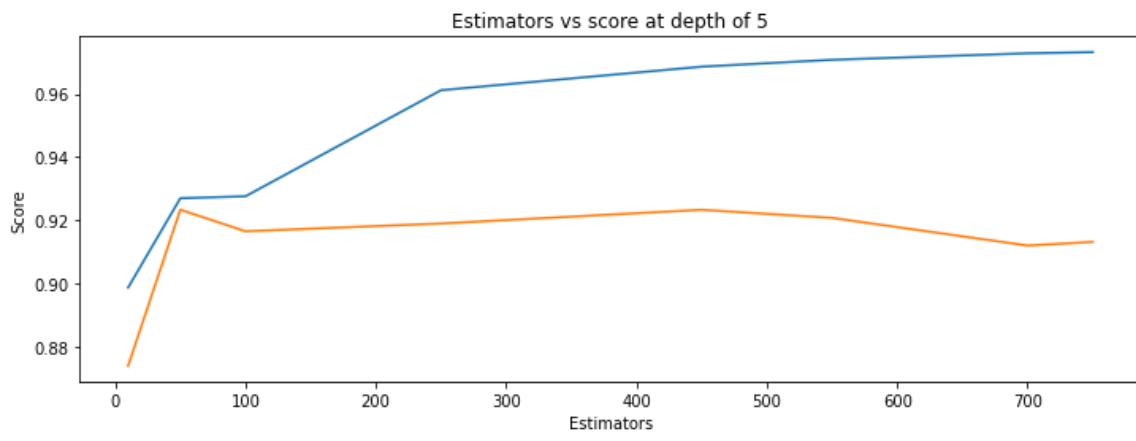
In []:

```
1 #Implementing XGBoost Here
2 estimators = [10,50,100,250,450,550,700,750]
3 train_scores = []
4 test_scores = []
5 for i in estimators:
6     clf = XGBClassifier(silent=False,
7                         scale_pos_weight=1,
8                         learning_rate=0.01,
9                         colsample_bytree = 0.4,
10                        subsample = 0.8,
11                        objective='binary:logistic',
12                        n_estimators=i,
13                        reg_alpha = 0.3,
14                        max_depth=4,
15                        gamma=10)
16     clf.fit(df_final_train,y_train)
17     train_sc = f1_score(y_train,clf.predict(df_final_train))
18     test_sc = f1_score(y_test,clf.predict(df_final_test))
19     test_scores.append(test_sc)
20     train_scores.append(train_sc)
21     print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
22 plt.plot(estimators,train_scores,label='Train Score')
23 plt.plot(estimators,test_scores,label='Test Score')
24 plt.xlabel('Estimators')
25 plt.ylabel('Score')
26 plt.title('Estimators vs score at depth of 5')
```

```
Estimators = 10 Train Score 0.8987468187236235 test Score 0.874023694602896
Estimators = 50 Train Score 0.9269906841461872 test Score 0.923373351119016
Estimators = 100 Train Score 0.9276093214818552 test Score 0.91654627347667
22
Estimators = 250 Train Score 0.961156351791531 test Score 0.919004831736882
3
Estimators = 450 Train Score 0.9686075025839531 test Score 0.92333778427756
94
Estimators = 550 Train Score 0.9707796479792714 test Score 0.92078217758580
34
Estimators = 700 Train Score 0.9728346098200104 test Score 0.91203574420552
91
Estimators = 750 Train Score 0.973202720207254 test Score 0.913176930176651
1
```

Out[90]:

Text(0.5, 1.0, 'Estimators vs score at depth of 5')



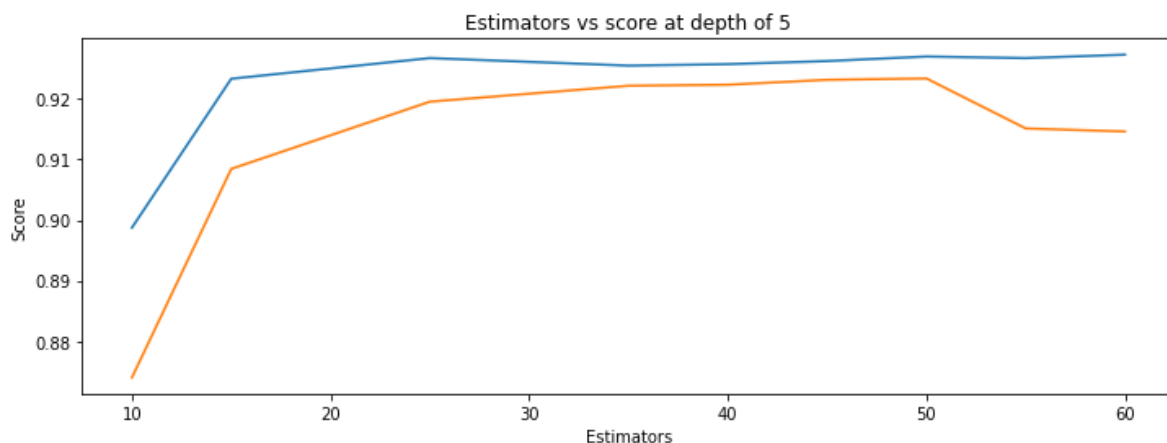
In []:

```
1 estimators = [10,15,25,35,40,45,50,55,60]
2 train_scores = []
3 test_scores = []
4 for i in estimators:
5     clf = XGBClassifier(silent=False,
6                         scale_pos_weight=1,
7                         learning_rate=0.01,
8                         colsample_bytree = 0.4,
9                         subsample = 0.8,
10                        objective='binary:logistic',
11                        n_estimators=i,
12                        reg_alpha = 0.3,
13                        max_depth=4,
14                        gamma=10)
15     clf.fit(df_final_train,y_train)
16     train_sc = f1_score(y_train,clf.predict(df_final_train))
17     test_sc = f1_score(y_test,clf.predict(df_final_test))
18     test_scores.append(test_sc)
19     train_scores.append(train_sc)
20     print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
21 plt.plot(estimators,train_scores,label='Train Score')
22 plt.plot(estimators,test_scores,label='Test Score')
23 plt.xlabel('Estimators')
24 plt.ylabel('Score')
25 plt.title('Estimators vs score at depth of 5')
```

```
Estimators = 10 Train Score 0.8987468187236235 test Score 0.874023694602896
Estimators = 15 Train Score 0.9233213057153978 test Score 0.908443040945654
8
Estimators = 25 Train Score 0.9267309174205398 test Score 0.919538271499765
7
Estimators = 35 Train Score 0.9254898653990071 test Score 0.922161269141973
6
Estimators = 40 Train Score 0.925730673527055 test Score 0.922333799415279
Estimators = 45 Train Score 0.9262387565509896 test Score 0.923142142524057
9
Estimators = 50 Train Score 0.9269906841461872 test Score 0.923373351119016
Estimators = 55 Train Score 0.9267256786022161 test Score 0.915124441608168
5
Estimators = 60 Train Score 0.9272975814931651 test Score 0.914611813806907
7
```

Out[91]:

Text(0.5, 1.0, 'Estimators vs score at depth of 5')



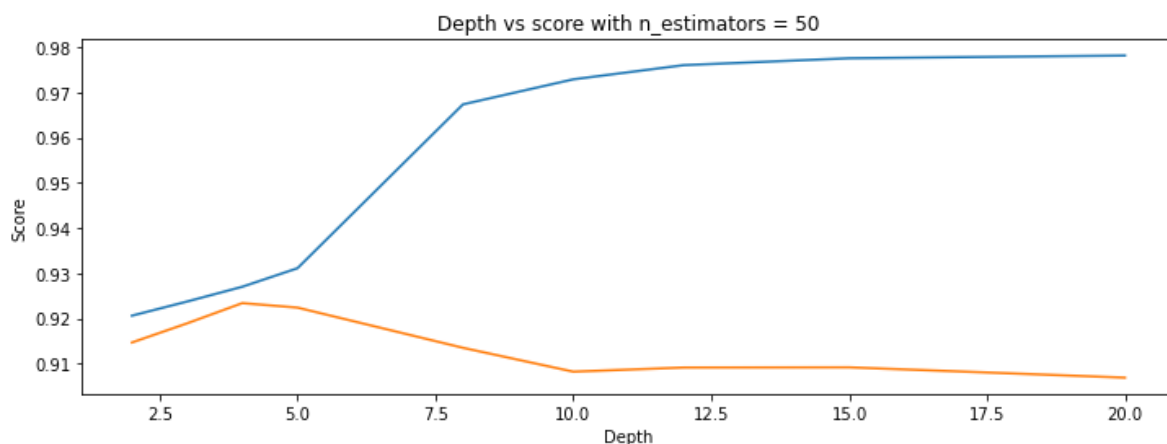
In []:

```
1 depths = [2,3,4,5,8,10,12,15,20]
2 train_scores = []
3 test_scores = []
4 for i in depths:
5     clf = XGBClassifier(silent=False,
6                         scale_pos_weight=1,
7                         learning_rate=0.01,
8                         colsample_bytree = 0.4,
9                         subsample = 0.8,
10                        objective='binary:logistic',
11                        n_estimators=50,
12                        reg_alpha = 0.3,
13                        max_depth=i,
14                        gamma=10)
15     clf.fit(df_final_train,y_train)
16     train_sc = f1_score(y_train,clf.predict(df_final_train))
17     test_sc = f1_score(y_test,clf.predict(df_final_test))
18     test_scores.append(test_sc)
19     train_scores.append(train_sc)
20     print('Depth = ',i,'Train Score',train_sc,'test Score',test_sc)
21 plt.plot(depths,train_scores,label='Train Score')
22 plt.plot(depths,test_scores,label='Test Score')
23 plt.xlabel('Depth')
24 plt.ylabel('Score')
25 plt.title('Depth vs score with n_estimators = 50')
```

```
Depth = 2 Train Score 0.9205948018029403 test Score 0.9146623153701074
Depth = 3 Train Score 0.9237252687497386 test Score 0.9188973400008404
Depth = 4 Train Score 0.9269906841461872 test Score 0.923373351119016
Depth = 5 Train Score 0.9310883818515604 test Score 0.9223829390711366
Depth = 8 Train Score 0.9673295167196002 test Score 0.9134871434353763
Depth = 10 Train Score 0.972870167246966 test Score 0.908207805546046
Depth = 12 Train Score 0.976004016064257 test Score 0.9091258733787739
Depth = 15 Train Score 0.9775280898876405 test Score 0.9091686623332192
Depth = 20 Train Score 0.9781313323836848 test Score 0.9068895317689918
```

Out[94]:

Text(0.5, 1.0, 'Depth vs score with n_estimators = 50')



In []:

```
1 from scipy.stats import randint as sp_randint
2 from scipy.stats import uniform
```

In []:

```
1 param_dist = {"n_estimators": [40, 45, 50, 55, 60],
2               "max_depth": [3, 4, 5, 6, 7],
3               "learning_rate": [0.01, 0.05, 0.1, 0.15, 0.2],
4               "subsample": [0.75, 0.8, 0.85, 0.9, 0.95],
5               "gamma": [0, 1, 5, 8, 10],
6               "colsample_bytree": [0.4, 0.6, 0.8, 0.9, 1],
7               "reg_alpha": [0.2, 0.3, 0.4, 0.5, 0.6]}
9
10 clf = XGBClassifier(random_state=25, n_jobs=-1)
11
12 xgb_random = RandomizedSearchCV(clf, param_distributions=param_dist,
13                                 n_iter=5, cv=10, scoring='f1', random_state=25, return_
14
15 xgb_random.fit(df_final_train, y_train)
16 print('mean test scores', xgb_random.cv_results_['mean_test_score'])
17 print('mean train scores', xgb_random.cv_results_['mean_train_score'])
```

```
mean test scores [0.9741902  0.92802887 0.96636358 0.97428761 0.9773855 ]
mean train scores [0.9749252  0.92821009 0.96688544 0.97487048 0.98057475]
```

In []:

```
1 print(xgb_random.best_estimator_)
```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=0.9, gamma=8,
              learning_rate=0.2, max_delta_step=0, max_depth=7,
              min_child_weight=1, missing=None, n_estimators=50, n_jobs=-1,
              nthread=None, objective='binary:logistic', random_state=25,
              reg_alpha=0.4, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=None, subsample=0.95, verbosity=1)
```

In []:

```
1 clf = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
2                     colsample_bynode=1, colsample_bytree=0.9, gamma=8,
3                     learning_rate=0.2, max_delta_step=0, max_depth=7,
4                     min_child_weight=1, missing=None, n_estimators=50, n_jobs=-1,
5                     nthread=None, objective='binary:logistic', random_state=25,
6                     reg_alpha=0.4, reg_lambda=1, scale_pos_weight=1, seed=None,
7                     silent=None, subsample=0.95, verbosity=1)
```

In []:

```
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)
```

In []:

```
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 0.980836113995383
Test f1 score 0.9314090102273688
```

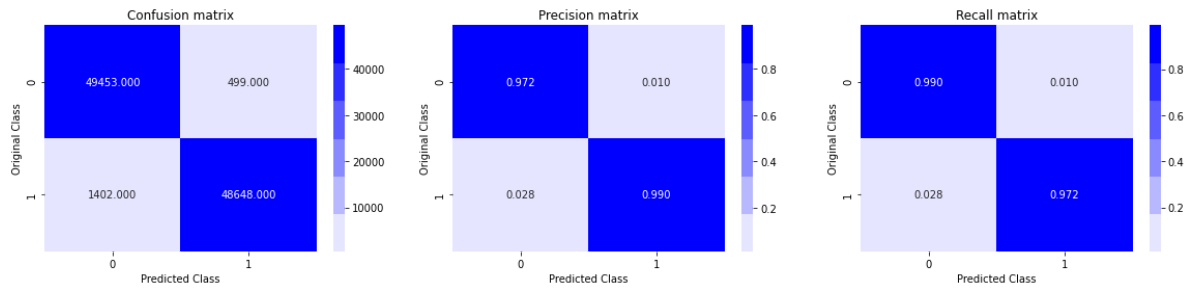
In []:

```
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=labels)
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=labels)
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=labels)
28    plt.xlabel('Predicted Class')
29    plt.ylabel('Original Class')
30    plt.title("Recall matrix")
31
32    plt.show()
```

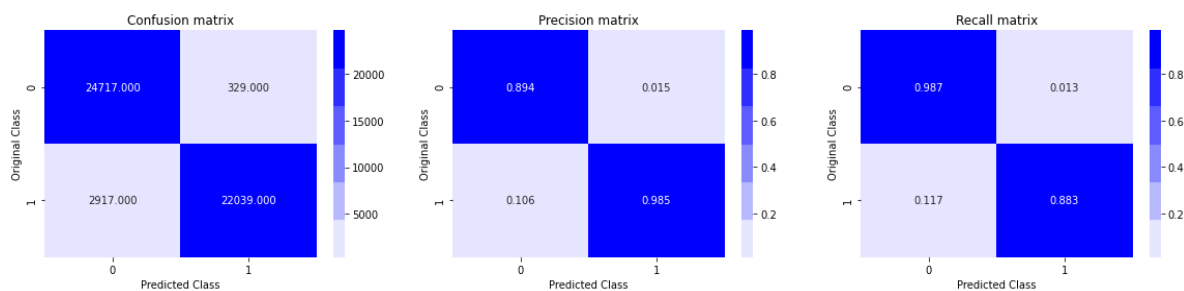
In []:

```
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

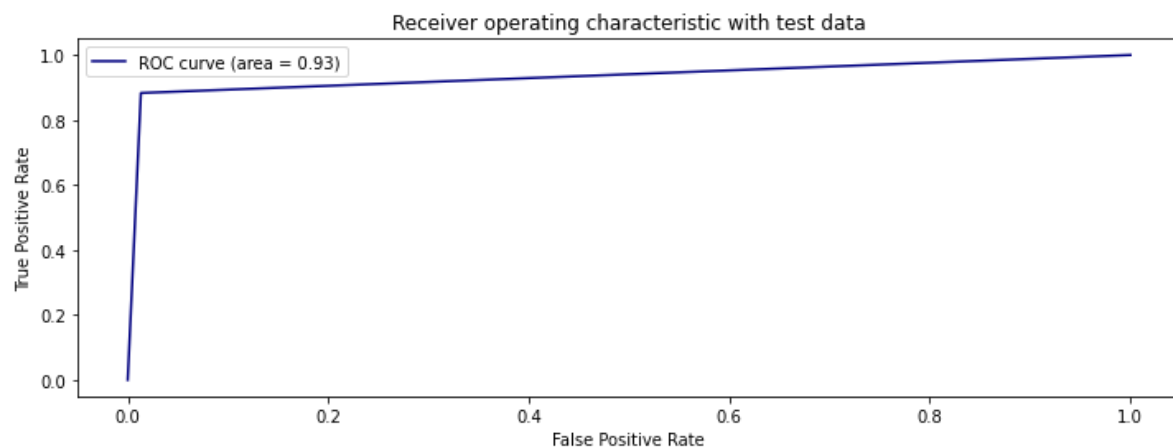


Test confusion_matrix



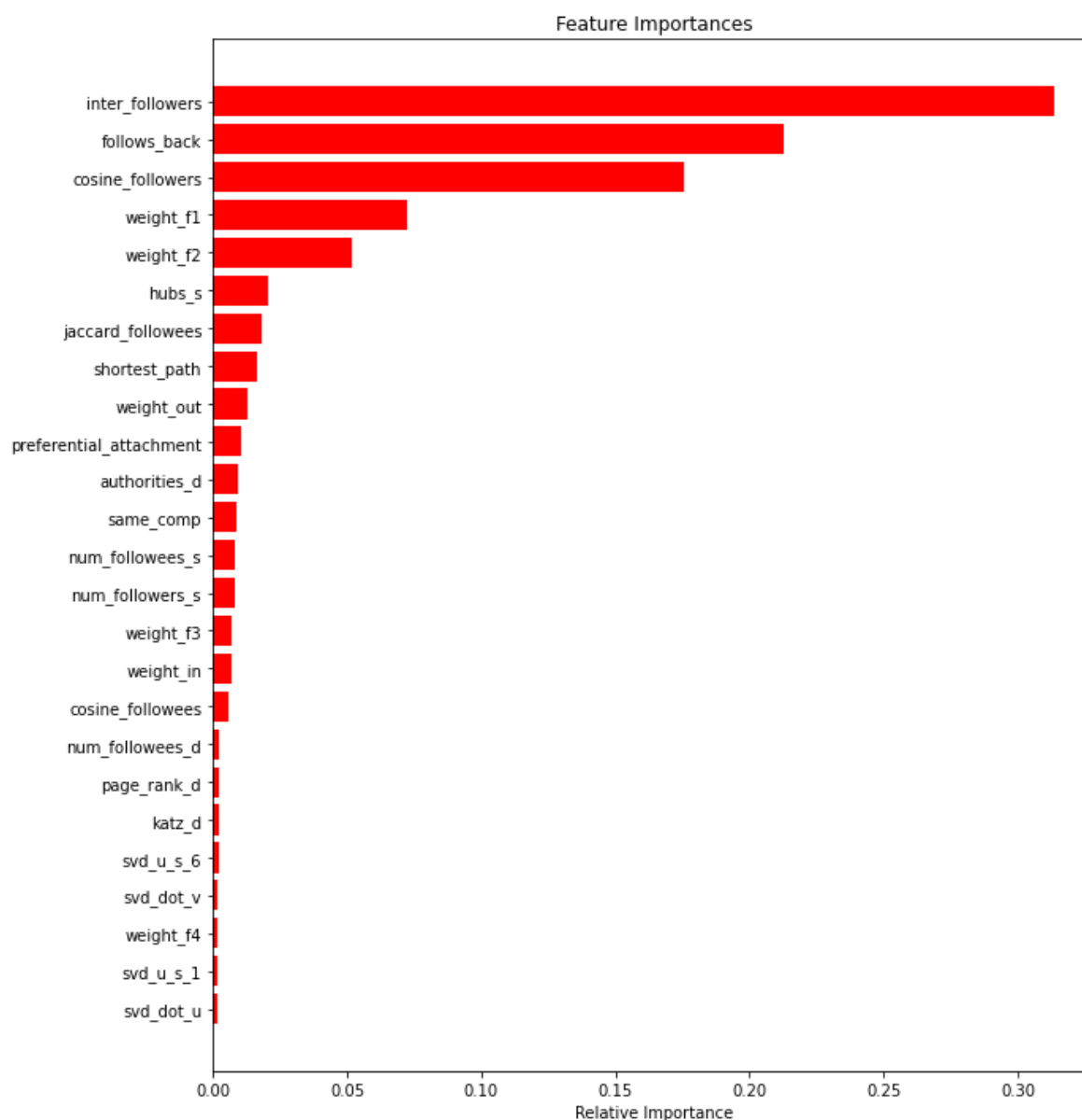
In []:

```
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()
```



In []:

```
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()
```



#####1. For Detailed Explanation of Steps upto addition of Preferential Attachment as a Feature, Please refer to the "**Facebook Case Study - Notes**" pdf attached with this assignment submission

#####2. Preferential Attachment:

#####2.1. One well-known concept in social networks is that users with many friends tend to create more connections in the future. This is due to the fact that in some social networks, like in finance, the rich get richer. We estimate how "rich" our two vertices are by calculating the multiplication between the number of friends ($|\Gamma(x)|$) or followers each vertex has. It may be noted that the similarity index does not require any node neighbor information; therefore, this similarity index has the lowest computational complexity.

#####3. SVD Dot Feature: 3.1. svd_dot is a Dot product between source node svd and destination node svd features.

#####4. Apply XGBoostClassifier on all these Features to predict whether link is present or not

#####5. Hyperparameter Tune the parameters of XGBoost Classifier to get the best possible metric value for this dataset and set of features

#####6. Results:

#####Train f1 score 0.980836113995383

#####Test f1 score 0.9314090102273688