

## Importing Libraries

In [94]:

```
1 import warnings
2 warnings.filterwarnings("ignore")
3
4 import csv
5 import pandas as pd
6 import datetime
7 import time
8 import numpy as np
9 import matplotlib
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 from matplotlib import rcParams
13 from sklearn.cluster import MiniBatchKMeans, KMean
14 import math
15 import pickle
16 import os
17 from xgboost import XGBClassifier
18 import networkx as nx
19 import pdb
20 import pickle
21 from pandas import HDFStore, DataFrame
22 from pandas import read_hdf
23 from scipy.sparse.linalg import svds, eigs
24 import gc
25 from tqdm import tqdm
26 from sklearn.ensemble import RandomForestClassifier
27 from sklearn.metrics import f1_score
```

## Reading Data

In [24]:

```
1
2 from pandas import read_hdf
3 df_final_train = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'train_df', mode='r')
4 df_final_test = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'test_df', mode='r')
```

## Adding Feature Preferential Attachment

**Preferential Attachment** 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.

$$\text{Score}(x, y) = |\Gamma(x)| \cdot |\Gamma(y)|$$

In [28]:

```

1 def preferential_attachment(df_final):
2     preferential_attachment_followers=[] # function for calculating preferential_atta
3     preferential_attachment_followees=[]
4     preferential_attachment_followers=df_final['num_followers_s']*df_final['num_follow
5     preferential_attachment_followees=df_final['num_followees_s']*df_final['num_follow
6     return preferential_attachment_followers,preferential_attachment_followees

```

In [109]:

```

1 #adding the feature into both x_test and x_train
2 (df_final_train[' preferential_attachment_followers'], df_final_train['preferential_at
3 (df_final_test[' preferential_attachment_followers'], df_final_test['preferential_atta
4

```

In [32]:

```

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

```

In [37]:

```

1 def svd(x, S):
2     try:
3         z = sadj_dict[x]
4         return S[z]
5     except:
6         return [0,0,0,0,0,0]

```

In [33]:

```

1 #for svd features to get feature vector creating a dict node val and inedx in svd vecto
2 sadj_col = sorted(train_graph.nodes())
3 sadj_dict = { val:idx for idx,val in enumerate(sadj_col)}

```

In [34]:

```

1 Adj = nx.adjacency_matrix(train_graph,nodelist=sorted(train_graph.nodes())).asfptype()

```

In [35]:

```

1 U, s, V = svds(Adj, k = 6)
2 print('Adjacency matrix Shape',Adj.shape)
3 print('U Shape',U.shape)
4 print('V Shape',V.shape)
5 print('s Shape',s.shape)

```

Adjacency matrix Shape (1780722, 1780722)

U Shape (1780722, 6)

V Shape (6, 1780722)

s Shape (6,)

## Adding Feature svd\_dot

**svd\_dot** is the Dot product between source node svd and destination node svd features

In [48]:

```

1 def np_dot(df_final):
2     # function product between source node and destination of svd features
3     np_dot_u=[]
4     np_dot_v=[]
5     for i,row in df_final.iterrows():
6         a=svd(row['source_node'],U)
7         b=svd(row['destination_node'],U)
8         np_dot_u.append(np.dot(a,b))
9         c=svd(row['source_node'],V.T)
10        d=svd(row['destination_node'],V.T)
11        np_dot_v.append(np.dot(c,d))
12
13    return np_dot_u,np_dot_v

```

In [51]:

```

1 (df_final_train['np_dot_u'],df_final_train['np_dot_v'])=np_dot(df_final_train)
2 (df_final_test['np_dot_u'],df_final_test['np_dot_v'])=np_dot(df_final_test)

```

In [110]:

```
1 df_final_train.columns
```

Out[110]:

```

Index(['jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_inde
x',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_ou
t',
      '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',
      'num_followers_d', ' preferential_attachment_followers',
      'preferential_attachment_followees', 'np_dot_u', 'np_dot_v'],
      dtype='object')

```

In [55]:

```

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

```

In [56]:

```

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 [95]:

```

1 estimators = [10,50,100,150,200]
2 train_scores = []
3 test_scores = []
4 for i in estimators:
5     clf = XGBClassifier(n_estimators=i, n_jobs=-1)
6     clf.fit(df_final_train,y_train)
7     train_sc = f1_score(y_train,clf.predict(df_final_train))
8     test_sc = f1_score(y_test,clf.predict(df_final_test))
9     test_scores.append(test_sc)
10    train_scores.append(train_sc)
11    print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
12 plt.plot(estimators,train_scores,label='Train Score')
13 plt.plot(estimators,test_scores,label='Test Score')
14 plt.xlabel('Estimators')
15 plt.ylabel('Score')
16 plt.legend()
17 plt.title('Estimators vs score at depth of 5')

```

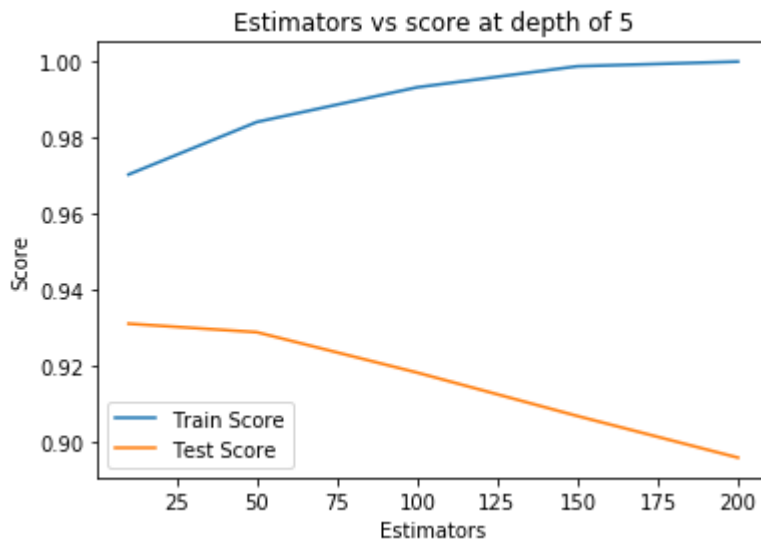
```

Estimators = 10 Train Score 0.970206342244571 test Score 0.9311625734966357
Estimators = 50 Train Score 0.9839552802538148 test Score 0.928928844033584
9
Estimators = 100 Train Score 0.9930318130320136 test Score 0.91835337412811
85
Estimators = 150 Train Score 0.9985403211293514 test Score 0.90697976128697
44
Estimators = 200 Train Score 0.9997601822614812 test Score 0.89609538784067
09

```

Out[95]:

Text(0.5,1,'Estimators vs score at depth of 5')



In [97]:

```

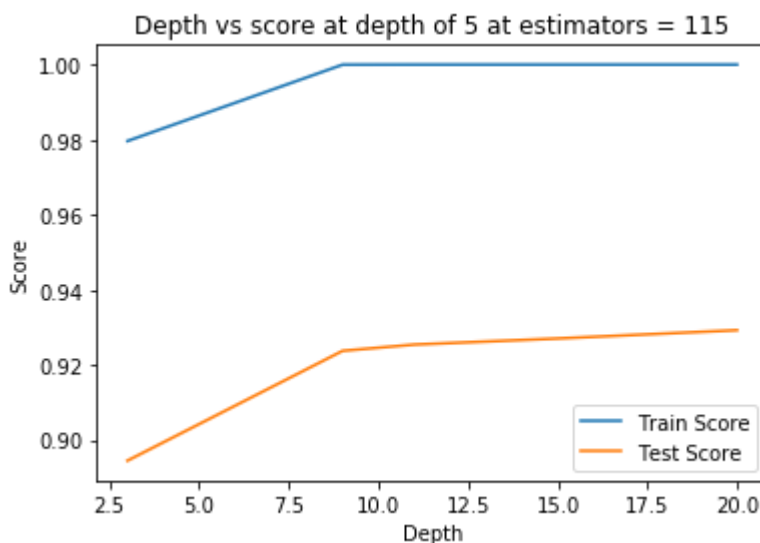
1 depths = [3,9,11,15,20]
2 train_scores = []
3 test_scores = []
4 for i in depths:
5     clf = XGBClassifier(max_depth=i, n_jobs=-1)
6     clf.fit(df_final_train,y_train)
7     train_sc = f1_score(y_train,clf.predict(df_final_train))
8     test_sc = f1_score(y_test,clf.predict(df_final_test))
9     test_scores.append(test_sc)
10    train_scores.append(train_sc)
11    print('depth = ',i,'Train Score',train_sc,'test Score',test_sc)
12 plt.plot(depths,train_scores,label='Train Score')
13 plt.plot(depths,test_scores,label='Test Score')
14 plt.xlabel('Depth')
15 plt.ylabel('Score')
16 plt.legend()
17 plt.title('Depth vs score at depth of 5 at estimators = 115')
18 plt.show()

```

```

depth = 3 Train Score 0.9796490520371117 test Score 0.8945009292664262
depth = 9 Train Score 0.9999800195808108 test Score 0.9237517300117108
depth = 11 Train Score 1.0 test Score 0.9254175343164337
depth = 15 Train Score 1.0 test Score 0.9270393550440225
depth = 20 Train Score 1.0 test Score 0.9292501745994793

```



In [98]:

```

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             'learning_rate': [0.1, 0.01, 0.05]
11             }
12
13 clf = XGBClassifier( n_jobs=-1)
14
15 rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
16                               n_iter=3,cv=3,scoring='f1',random_state=25,return_t
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 [0.98027002 0.97432382 0.97266916]

mean train scores [0.99934022 0.98917069 0.97726836]

In [99]:

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

```

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
              importance_type='gain', interaction_constraints='',
              learning_rate=0.1, max_delta_step=0, max_depth=12,
              min_child_weight=1, missing=nan, monotone_constraints='()',
              n_estimators=120, n_jobs=-1, num_parallel_tree=1,
              objective='binary:logistic', random_state=0, reg_alpha=0,
              reg_lambda=1, scale_pos_weight=1, subsample=1,
              tree_method='exact', validate_parameters=1, verbosity=None)

```

In [100]:

```
1 clf=rf_random.best_estimator_
```

In [101]:

```

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 [102]:

```

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 0.9984894713149604

Test f1 score 0.9269121331804476

In [103]:

```
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()
```

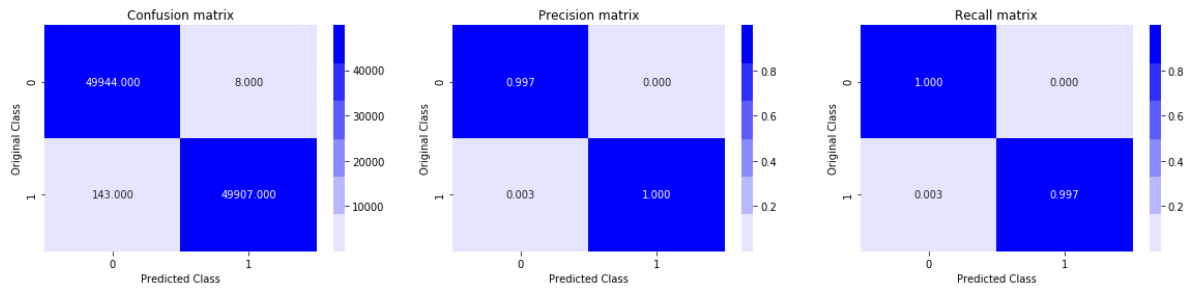
In [104]:

```

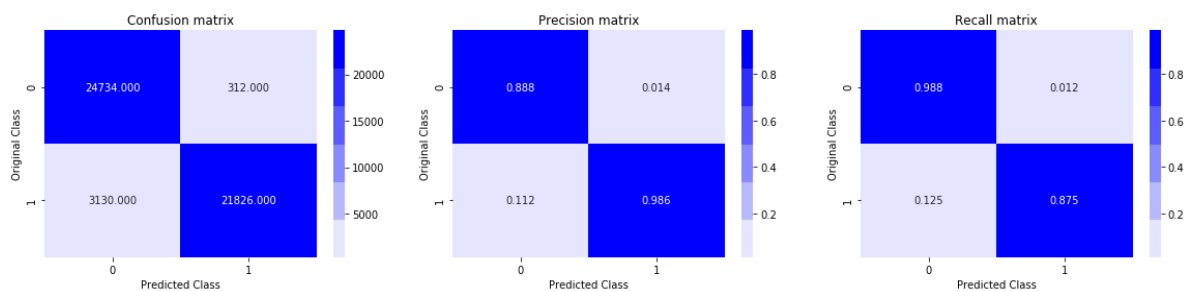
1 print('Train confusion_matrix')
2 plot_confusion_matrix(y_train,y_train_pred)
3 print('Test confusion_matrix')
4 plot_confusion_matrix(y_test,y_test_pred)

```

Train confusion\_matrix



Test confusion\_matrix

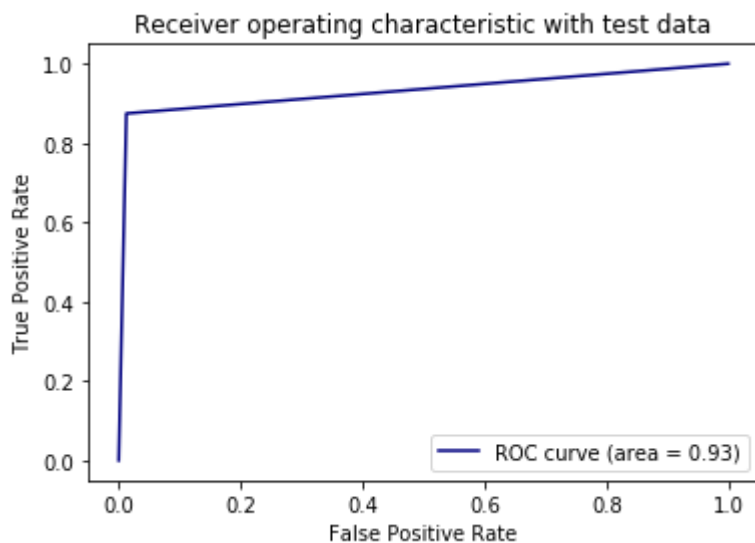


In [106]:

```

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()

```



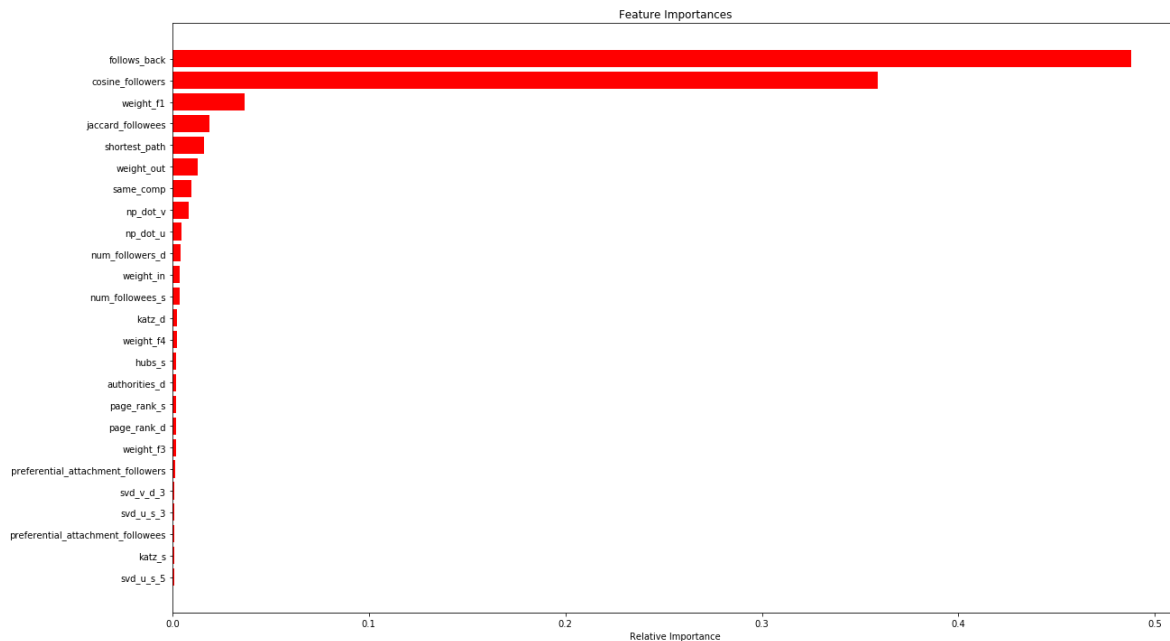


In [108]:

```

1 features = df_final_train.columns
2 importances = clf.feature_importances_
3 indices = (np.argsort(importances))[-25:]
4 plt.figure(figsize=(20,12))
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()

```



- for xgboost modelling ,model is overfitting also there decrease in precision and recall for a class
- still follow back is the most important feature
- new added feature svd\_dot have some on calssification

In [ ]:

1