

Social network Graph Link Prediction - Facebook Challenge

In [0]:

```
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
24 import warnings
25 import networkx as nx
26 import pdb
27 import pickle
28 from pandas import HDFStore, DataFrame
29 from pandas import read_hdf
30 from scipy.sparse.linalg import svds, eigs
31 import gc
32 from tqdm import tqdm
33 from sklearn.ensemble import RandomForestClassifier
34 from sklearn.metrics import f1_score
```

In []:

```
1 #reading
2 from pandas import read_hdf
3 df_final_train = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'train_
4 df_final_test = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'test_df
```

```
In [0]: 1 df_final_train.columns
```

```
Out[3]: 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')
```

```
In [0]: 1 y_train = df_final_train.indicator_link  
        2 y_test = df_final_test.indicator_link  
        3  
        4 df_final_train.drop(['source_node', 'destination_node', 'indicator_link'],axis=1)  
        5 df_final_test.drop(['source_node', 'destination_node', 'indicator_link'],axis=1)
```

```
In [0]: 1 df_final_train.drop(['source_node', 'destination_node', 'indicator_link'],axis=1)  
        2 df_final_test.drop(['source_node', 'destination_node', 'indicator_link'],axis=1)
```

```

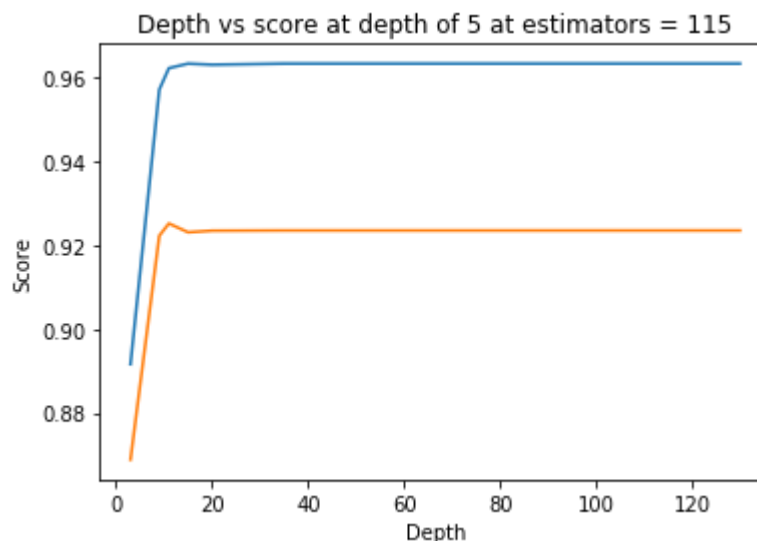
In [0]: 1 depths = [3,9,11,15,20,35,50,70,130]
2 train_scores = []
3 test_scores = []
4 for i in depths:
5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion
6         max_depth=i, max_features='auto', max_leaf_nodes=None,
7         min_impurity_decrease=0.0, min_impurity_split=None,
8         min_samples_leaf=52, min_samples_split=120,
9         min_weight_fraction_leaf=0.0, n_estimators=115, n_jobs=-1,random_
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

```



```

In [0]: 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,
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 [0.96225043 0.96215493 0.96057081 0.96194015 0.96330005]
mean train scores [0.96294922 0.96266735 0.96115674 0.96263457 0.96430539]

```

In [0]: 1 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)

```

```

In [0]: 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, min_impurity_split=None,
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)

```

```

In [0]: 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)

```

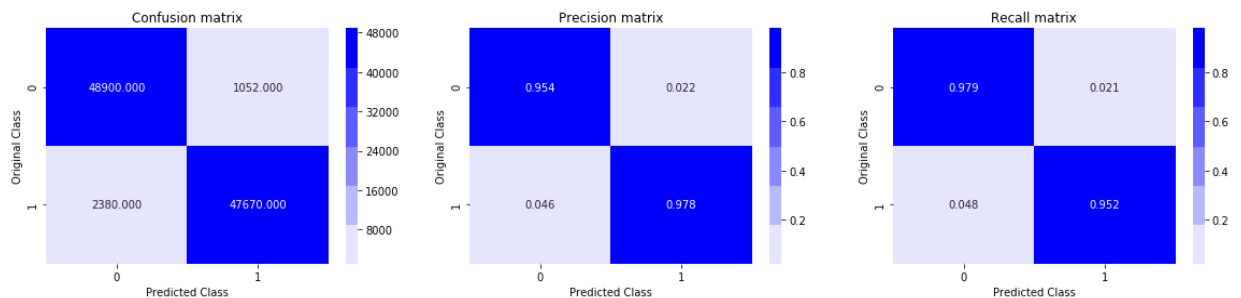
```
In [0]: 1 from sklearn.metrics import f1_score
2 print('Train f1 score',f1_score(y_train,y_train_pred))
3 print('Test f1 score',f1_score(y_test,y_test_pred))
4
5 print('Train confusion_matrix')
6 plot_confusion_matrix(y_train,y_train_pred)
7 print('Test confusion_matrix')
8 plot_confusion_matrix(y_test,y_test_pred)
```

Train f1 score 0.9652533106548414

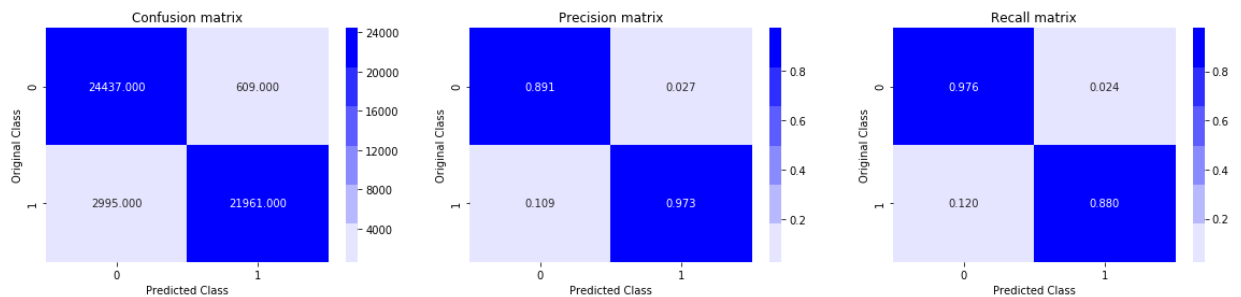
Test f1 score 0.9241678239279553

```
In [0]: 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)
```

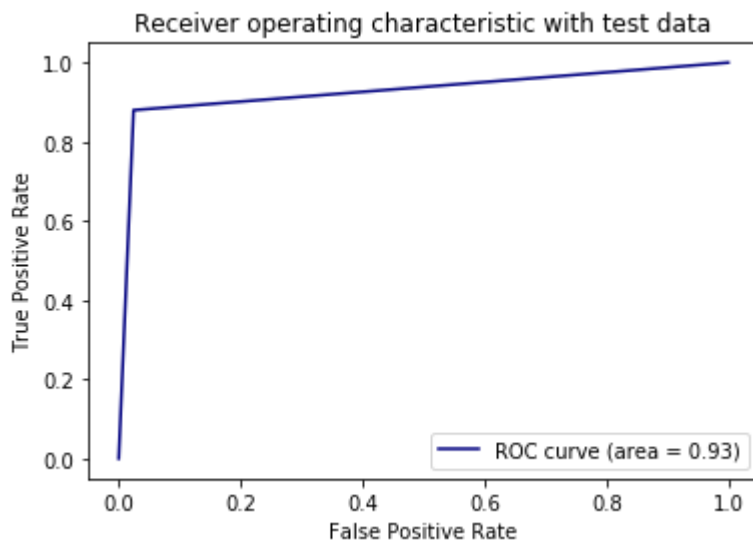
Train confusion_matrix



Test confusion_matrix



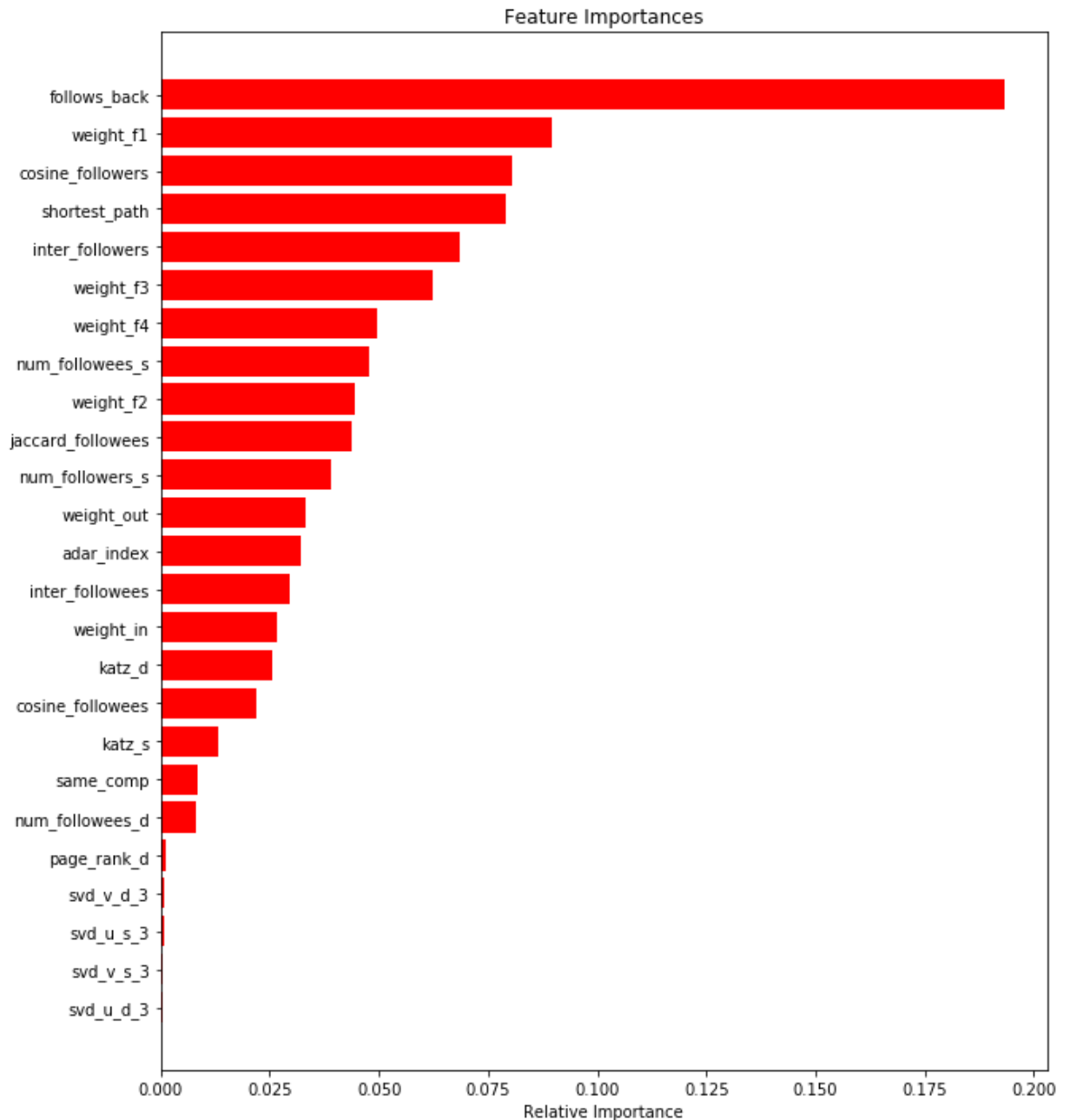
```
In [0]: 1 from sklearn.metrics import roc_curve, auc
2 fpr,tpr,ths = roc_curve(y_test,y_test_pred)
3 auc_sc = auc(fpr, tpr)
4 plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('Receiver operating characteristic with test data')
8 plt.legend()
9 plt.show()
```



```

In [0]: 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/> (<http://be.amazd.com/link-prediction/>)
2. Add feature called svd_dot. you can calculate svd_dot as Dot product between source 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 (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. Part 1: Feature engineering

In []: 1 In this part we will be doing the feature engineering


```
In [3]: 1 from google.colab import drive
        2 drive.mount('/content/drive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly (http s://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly)

Enter your authorization code:

.....

Mounted at /content/drive

```
In [0]: 1 #Loading the training data in graph
        2 import networkx as nx
        3 graph_data = nx.read_edgelist('drive/My Drive/facebook_casestdy/train_pos_aft
        4                                     nodetype = int)
```

```
In [0]: 1 #we get the final training and test data
        2 df_final_train = read_hdf('drive/My Drive/facebook_casestdy/storage_sample_st
        3 df_final_test = read_hdf('drive/My Drive/facebook_casestdy/storage_sample_sta
```

```
In [6]: 1 print('shape of final training data is:',df_final_train.shape)
        2 print('shape of final test data is:',df_final_test.shape)
```

shape of final training data is: (100002, 54)

shape of final test data is: (50002, 54)

1.1 feature engineering 1 - Preferential attachment

```

In [0]: 1 #we will write functions to compute the preferential attachment for followee
2
3 def followees_pa(user1,user2):
4     """This function computes the preferential attachment between followees of
5     destination
6     :input-user1,user2:the source and destination nodes
7     :output-returns the product of number of followees"""#where followees are
8     try:
9         user_1 = len(set(graph_data.successors(user_1)))#number of individuals fo
10        user_2 = len(set(graph_data.successors(user_2)))#number of individuals fo
11        return (user_1*user_2)
12    except:
13        return(0) #no pf if number of followees of any of them is 0
14
15
16
17 def followers_pa(user_1,user_2):
18     """This function computes the preferential attachment between followers of
19     destination
20     :input-user1,user2:the source and destination nodes
21     :output-returns the product of number of followers"""#where followers are
22     try:
23         user_1 = len(set(graph_data.predecessors(user_1)))#number of individuals
24         user_2 = len(set(graph_data.predecessors(user_2)))#number of individuals
25         return (user_1*user_2)
26     except:
27         return(0) #no pf if number of followers of any of them is 0
28

```

```

In [0]: 1 #added features for training data
2 df_final_train['pa_followees'] = df_final_train.apply(lambda x: followees_pa(x['source'],x['destination']))
3 df_final_train['pa_followers'] = df_final_train.apply(lambda x: followers_pa(x['source'],x['destination']))
4
5 #added features for test data
6 df_final_test['pa_followees'] = df_final_test.apply(lambda x: followees_pa(x['source'],x['destination']))
7 df_final_test['pa_followers'] = df_final_test.apply(lambda x: followers_pa(x['source'],x['destination']))
8
9

```

```

In [0]: 1 del graph_data #delete the graph data to release some memory

```

1.2 Feature engineering 2 - SVD_DOT

```
In [0]: 1 # for product of respective elements of right and left orthogonal vectors in
2
3 #For training data
4 svd_u_s_train = df_final_train[['svd_u_s_1','svd_u_s_2','svd_u_s_3','svd_u_s_4']
5 svd_v_s_train = df_final_train[['svd_v_s_1','svd_v_s_2','svd_v_s_3','svd_v_s_4']
6 svd_u_d_train = df_final_train[['svd_u_d_1','svd_u_d_2','svd_u_d_3','svd_u_d_4']
7 svd_v_d_train = df_final_train[['svd_v_d_1','svd_v_d_2','svd_v_d_3','svd_v_d_4']
8
9 #for test data
10 svd_u_s_test = df_final_test[['svd_u_s_1','svd_u_s_2','svd_u_s_3','svd_u_s_4']
11 svd_v_s_test = df_final_test[['svd_v_s_1','svd_v_s_2','svd_v_s_3','svd_v_s_4']
12 svd_u_d_test = df_final_test[['svd_u_d_1','svd_u_d_2','svd_u_d_3','svd_u_d_4']
13 svd_v_d_test = df_final_test[['svd_v_d_1','svd_v_d_2','svd_v_d_3','svd_v_d_4']
14
```

```
In [0]: 1 #function for computing the svd
2 def compute_svd_dot(df1,df2):
3     """Computes the respective right and left orthofonal vecotr values of sou
4     svd_dot = []
5     for i in range(0,df1.shape[0]):
6         svd_dot.append(np.dot(df1.ix[i].values,df2.ix[i].values))
7     #selectine the respective rows with .ix
8     return svd_dot
```

```
In [0]: 1 df_final_train['svd_dot_u'] = compute_svd_dot(svd_u_s_train,svd_u_d_train)
2 df_final_train['svd_dot_v'] = compute_svd_dot(svd_v_s_train,svd_v_d_train)
3 df_final_test['svd_dot_u'] = compute_svd_dot(svd_u_s_test,svd_u_d_test)
4 df_final_test['svd_dot_v'] = compute_svd_dot(svd_v_s_test,svd_v_d_test)
```

```
In [14]: 1 print('After all the featurization final shape of the training data is:',df_f
2 print('After all the featurization final shape of the test data is:',df_final
```

After all the featurization final shape of the training data is: (100002, 58)
 After all the featurization final shape of the test data is: (50002, 58)

```
In [15]: 1 df_final_test.columns
```

```
Out[15]: 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',
'pa_followees', 'pa_followers', 'svd_dot_u', 'svd_dot_v'],
dtype='object')
```

```
In [0]: 1 y_train = df_final_train.indicator_link #target variable for training data
2 y_test = df_final_test.indicator_link #target variable for test data
3
4 df_final_train.drop(['source_node', 'destination_node', 'indicator_link'],axis=
5 df_final_test.drop(['source_node', 'destination_node', 'indicator_link'],axis=
```

2. Machine learning models

Model 1 - Random Forest

2.1.1 Simply Tuning the number of estimators in RandomForest

```
In [ ]: 1 start = datetime.datetime.now()
2 estimators = [10,50,100,150,200,250,300,350,450,500]
3 train_scores = []
4 test_scores = []
5 for i in estimators:
6     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion
7         max_depth=5, max_features='auto', max_leaf_nodes=None,
8         min_impurity_decrease=0.0, min_impurity_split=None,
9         min_samples_leaf=52, min_samples_split=120,
10        min_weight_fraction_leaf=0.0, n_estimators=i, n_jobs=-1, random_st
11        #we are considering here
12        clf.fit(df_final_train,y_train)
13        train_sc = f1_score(y_train,clf.predict(df_final_train))
14        test_sc = f1_score(y_test,clf.predict(df_final_test))
15        test_scores.append(test_sc)
16        train_scores.append(train_sc)
17        print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
18
19 plt.figure(figsize = (10,5))
20 plt.plot(estimators,train_scores,label='Train Score')
21 plt.plot(estimators,test_scores,label='Test Score')
22 plt.grid()
23 plt.xlabel('Estimators',fontsize = 12)
24 plt.ylabel('Score',fontsize = 12)
25 plt.title('Estimators vs score at depth of 5',fontsize = 12)
26 print('total time taken:',datetime.datetime.now() - start)
```

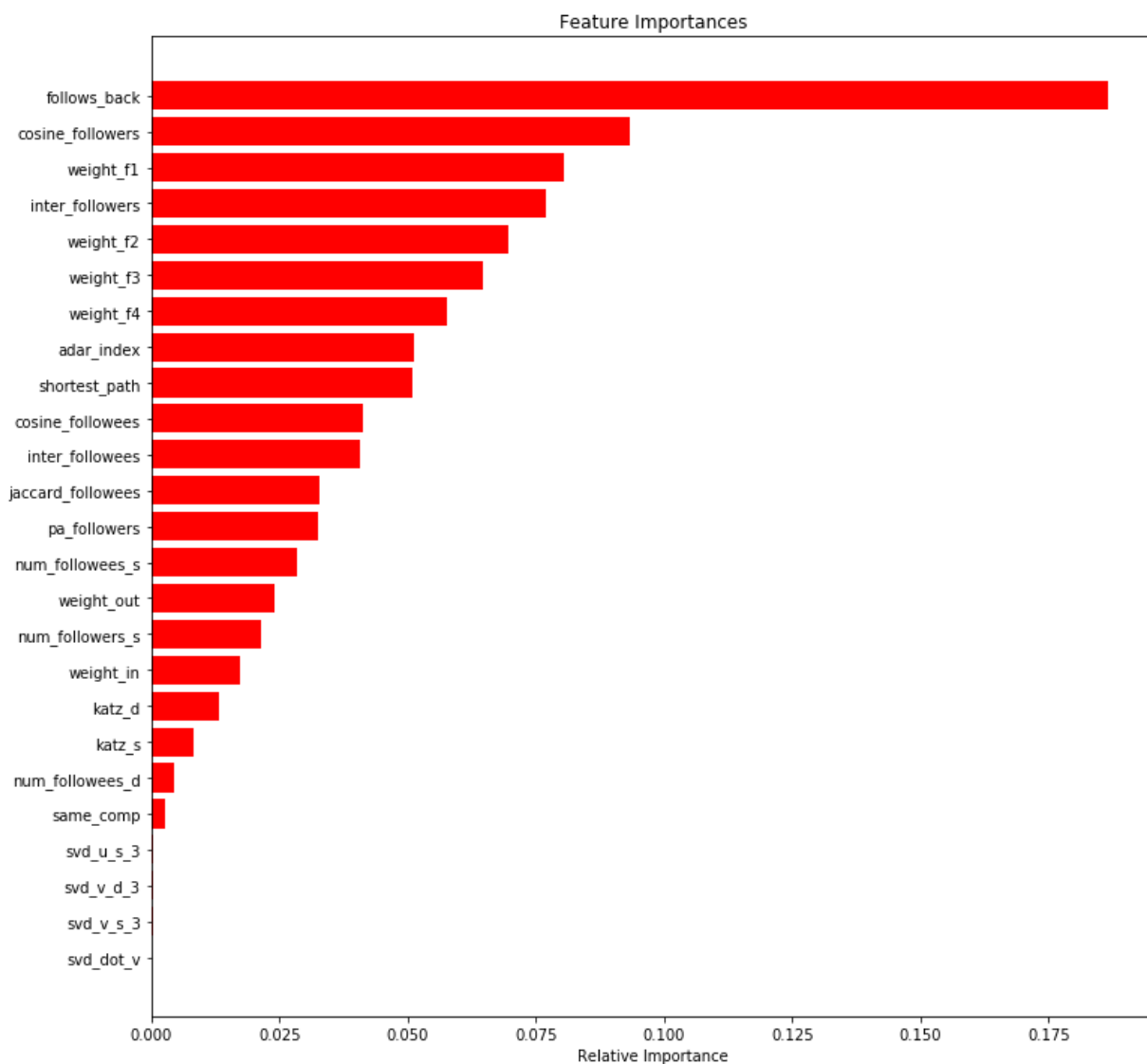
We see the best number of estimators we have is 200 with maximum depth 5

2.1.2 Best classifier and feature importances

```

In [28]: 1 #fitting the best etimator in the model for feature importances
2
3 clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
4                             max_depth=5, max_features='auto', max_leaf_nodes=None,
5                             min_impurity_decrease=0.0, min_impurity_split=None,
6                             min_samples_leaf=52, min_samples_split=120,
7                             min_weight_fraction_leaf=0.0, n_estimators=200, n_jobs=-1, random_state=42)
8 clf.fit(df_final_train,y_train)
9 features = df_final_train.columns
10 importances = clf.feature_importances_
11 indices = (np.argsort(importances))[-25:]#we are considering here top 25 features
12 plt.figure(figsize=(12,12))
13 plt.title('Feature Importances')
14 plt.barh(range(len(indices)), importances[indices], color='r', align='center')
15 plt.yticks(range(len(indices)), [features[i] for i in indices])
16 plt.xlabel('Relative Importance')
17 plt.show()

```



2.2.1 Tuning the depth in RandomForest

```

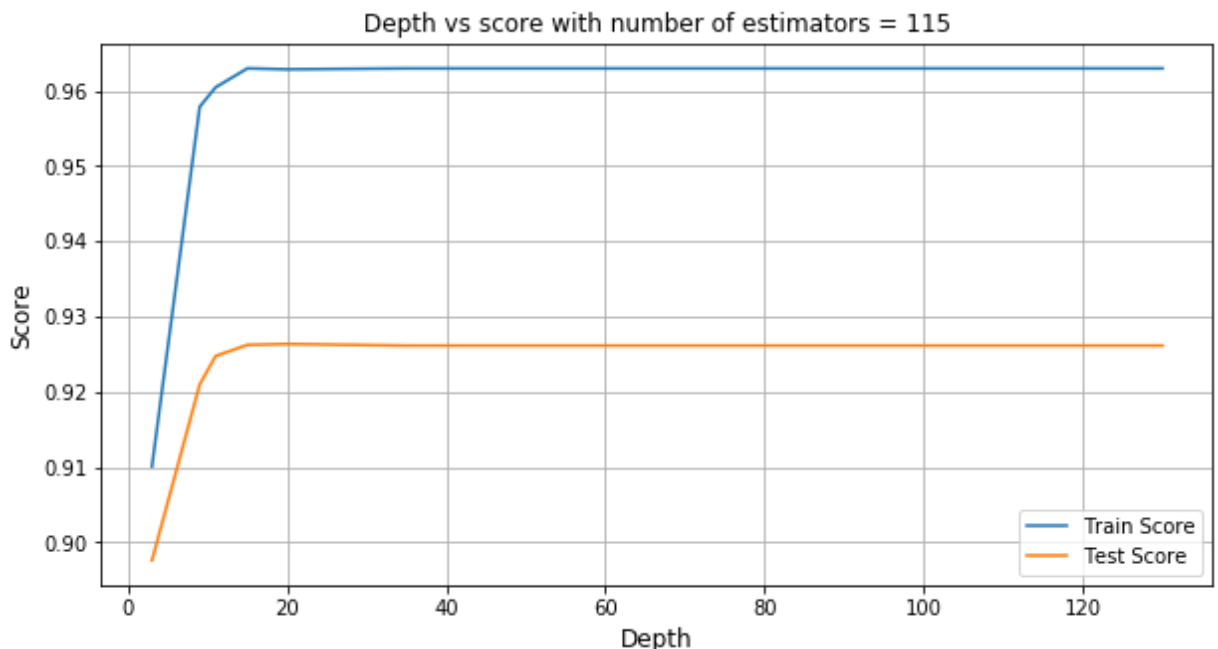
In [30]: 1 depths = [3,9,11,15,20,35,50,70,130]
2 train_scores = []
3 test_scores = []
4 for i in depths:
5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion
6         max_depth=i, max_features='auto', max_leaf_nodes=None,
7         min_impurity_decrease=0.0, min_impurity_split=None,
8         min_samples_leaf=52, min_samples_split=120,
9         min_weight_fraction_leaf=0.0, n_estimators=115, n_jobs=-1,random_
10     #we are considering here the
11     clf.fit(df_final_train,y_train)
12     train_sc = f1_score(y_train,clf.predict(df_final_train))
13     test_sc = f1_score(y_test,clf.predict(df_final_test))
14     test_scores.append(test_sc)
15     train_scores.append(train_sc)
16     print('depth = ',i,'Train Score',train_sc,'test Score',test_sc)
17
18 plt.figure(figsize = (10,5))
19 plt.plot(depths,train_scores,label='Train Score')
20 plt.plot(depths,test_scores,label='Test Score')
21 plt.grid()
22 plt.xlabel('Depth',fontsize = 12)
23 plt.ylabel('Score',fontsize = 12)
24 plt.title('Depth vs score with number of estimators = 115',fontsize = 12)
25 plt.legend(loc = 'best')
26 plt.show()

```

```

depth = 3 Train Score 0.9100877868473741 test Score 0.8975796072570266
depth = 9 Train Score 0.9579142403388036 test Score 0.9209777084168126
depth = 11 Train Score 0.9604904632152588 test Score 0.9247538086019695
depth = 15 Train Score 0.9630103985606685 test Score 0.9262395214625723
depth = 20 Train Score 0.962880942706217 test Score 0.9263556912709275
depth = 35 Train Score 0.9630058390454431 test Score 0.9261465283185655
depth = 50 Train Score 0.9630058390454431 test Score 0.9261465283185655
depth = 70 Train Score 0.9630058390454431 test Score 0.9261465283185655
depth = 130 Train Score 0.9630058390454431 test Score 0.9261465283185655

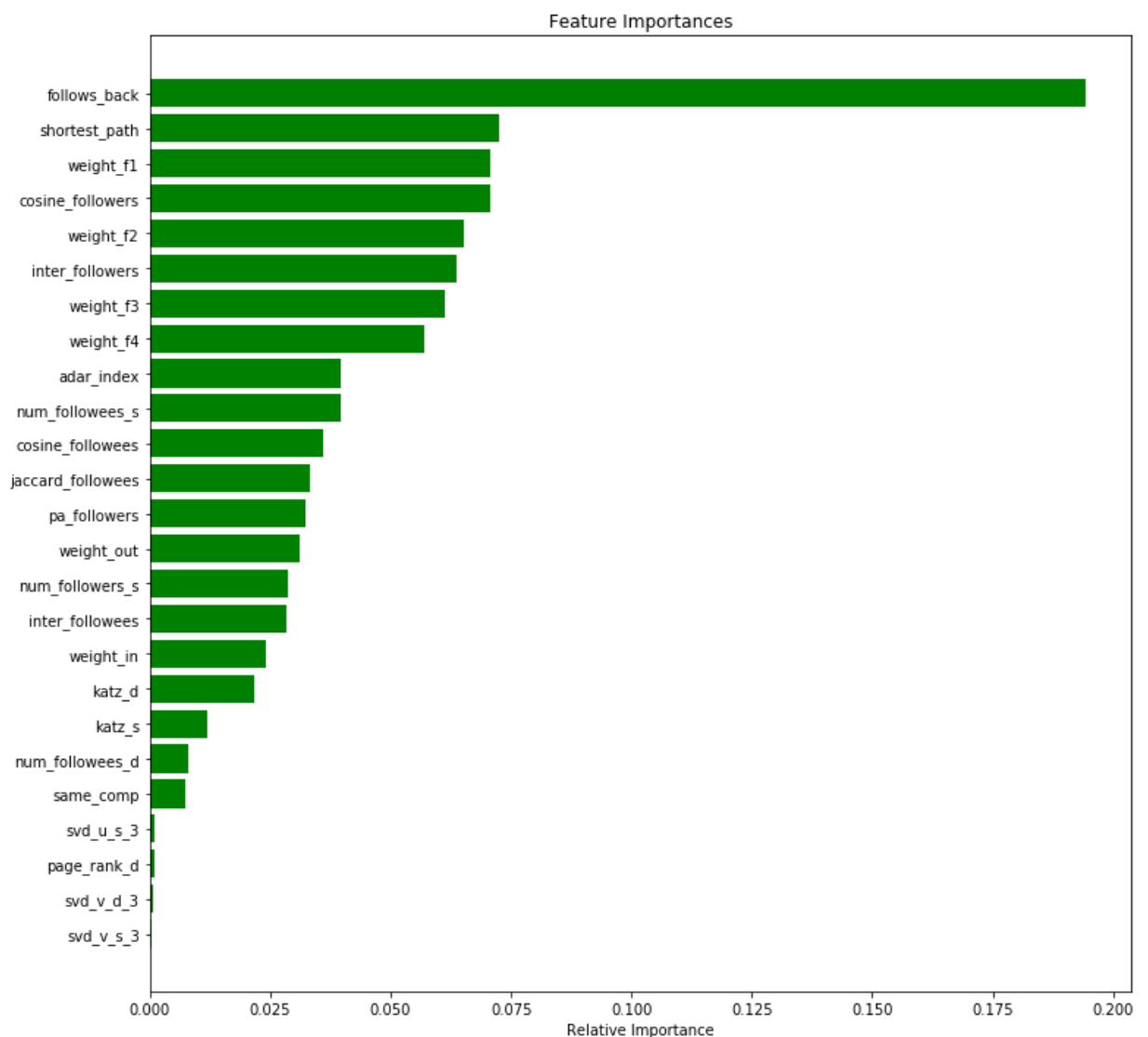
```



We see that the best score we get is at depth = 20

2.2.2 Best classifier and feature importances

```
In [32]: 1 clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
2           max_depth=20, max_features='auto', max_leaf_nodes=None,
3           min_impurity_decrease=0.0, min_impurity_split=None,
4           min_samples_leaf=52, min_samples_split=120,
5           min_weight_fraction_leaf=0.0, n_estimators=150, n_jobs=-1, random_state=None)
6 clf.fit(df_final_train, y_train)
7 features = df_final_train.columns
8 importances = clf.feature_importances_
9 indices = (np.argsort(importances))[-25:]
10 plt.figure(figsize=(12,12))
11 plt.title('Feature Importances')
12 plt.barh(range(len(indices)), importances[indices], color='g', align='center')
13 plt.yticks(range(len(indices)), [features[i] for i in indices])
14 plt.xlabel('Relative Importance')
15 plt.show()
```



2.3.1 Tuning the parameters for RandomForest using Randomized Search Cross validation

```
In [17]: 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 #for uniform discrete random va
6 from scipy.stats import uniform
7
8 param_dist = {"n_estimators":sp_randint(105,125),#number of estimator tress
9              "max_depth": sp_randint(10,15),#the depth of the the trees
10             "min_samples_split": sp_randint(110,190), #number of splits
11             "min_samples_leaf": sp_randint(25,65)} #number of leafs
12 start = datetime.datetime.now()
13 clf = RandomForestClassifier(random_state=25,n_jobs=-1)
14
15 rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
16                               n_iter=5,cv=10,scoring='f1',random_state=2
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'])
21 print('total time taken for computation is:',datetime.datetime.now() - start)
```

Fitting 10 folds for each of 5 candidates, totalling 50 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent worker s.

[Parallel(n_jobs=1)]: Done 50 out of 50 | elapsed: 6.8min finished

mean test scores [0.9615244 0.96117253 0.95967696 0.9607424 0.96274533]
 mean train scores [0.96229498 0.96182509 0.96001391 0.96160653 0.96378369]
 total time taken for computation is: 0:06:58.170276

```
In [18]: 1 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)
```

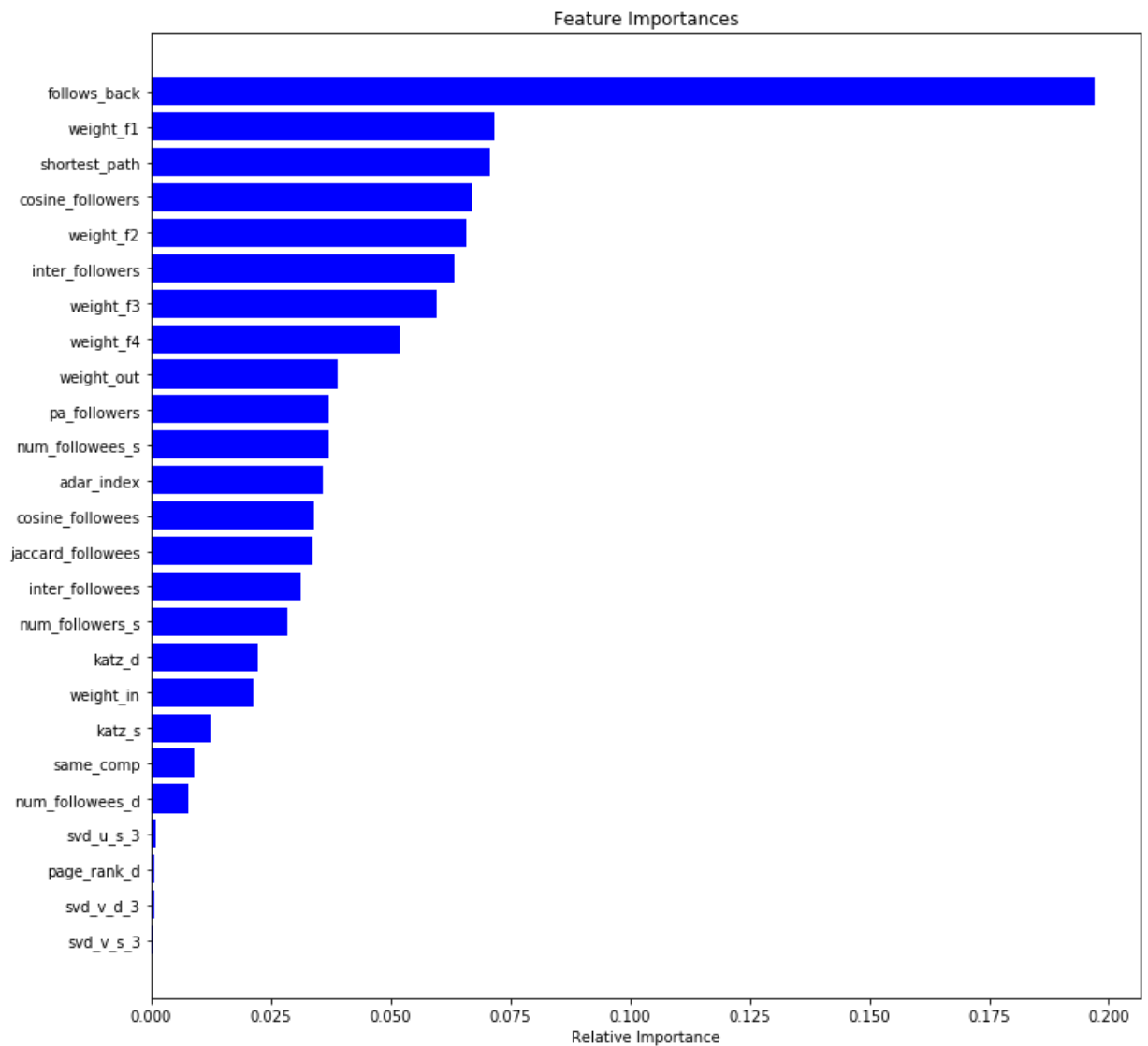
2.3.2 Best Classifier and feature importances

In [28]:

```

1  #training with the best model for feature importances
2
3  clf = rf_random.best_estimator_ #we get the best model
4  clf.fit(df_final_train,y_train)
5  features = df_final_train.columns
6  importances = clf.feature_importances_
7  indices = (np.argsort(importances))[-25:]
8  plt.figure(figsize=(12,12))
9  plt.title('Feature Importances')
10 plt.barh(range(len(indices)), importances[indices], color='b', align='center')
11 plt.yticks(range(len(indices)), [features[i] for i in indices])
12 plt.xlabel('Relative Importance')
13 plt.show()

```



```

In [0]: 1 #function for calculating the confusion,precison and recall matrix
2 from sklearn.metrics import confusion_matrix
3 def plot_confusion_matrix(test_y, predict_y):
4     C = confusion_matrix(test_y, predict_y)#confusion matrix
5
6     A = (((C.T)/(C.sum(axis=1))).T) #recallmatrix
7
8     B = (C/C.sum(axis=0)) #precision matrix
9     plt.figure(figsize=(20,4))
10
11     labels = [0,1]
12     # representing A in heatmap format
13     cmap=sns.light_palette("blue")
14     plt.subplot(1, 3, 1)
15     sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytic
16     plt.xlabel('Predicted Class')
17     plt.ylabel('Original Class')
18     plt.title("Confusion matrix")
19
20     plt.subplot(1, 3, 2)
21     sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytic
22     plt.xlabel('Predicted Class')
23     plt.ylabel('Original Class')
24     plt.title("Precision matrix")
25
26     plt.subplot(1, 3, 3)
27     # representing B in heatmap format
28     sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytic
29     plt.xlabel('Predicted Class')
30     plt.ylabel('Original Class')
31     plt.title("Recall matrix")
32
33     plt.show()

```

2.3.3 Confusion matrix and ROC curve

```

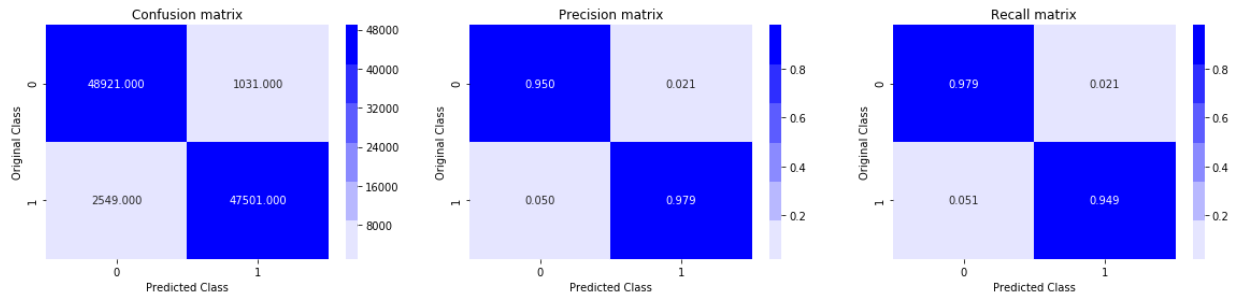
In [24]: 1 from sklearn.metrics import f1_score
2 print('Train f1 score',f1_score(y_train,y_train_pred))#best score with classi
3 print('Test f1 score',f1_score(y_test,y_test_pred))#best test score with clas
4
5 print('Train confusion_matrix')
6 plot_confusion_matrix(y_train,y_train_pred)
7 print('Test confusion_matrix')
8 plot_confusion_matrix(y_test,y_test_pred)

```

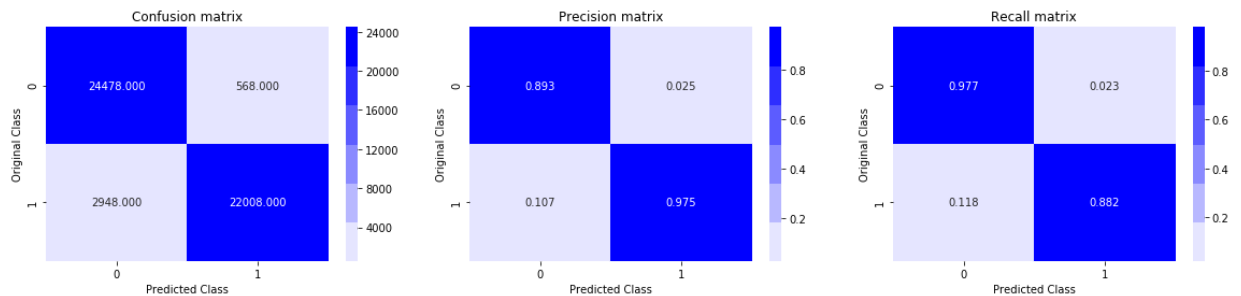
Train f1 score 0.9636850540666653

Test f1 score 0.9260287806109567

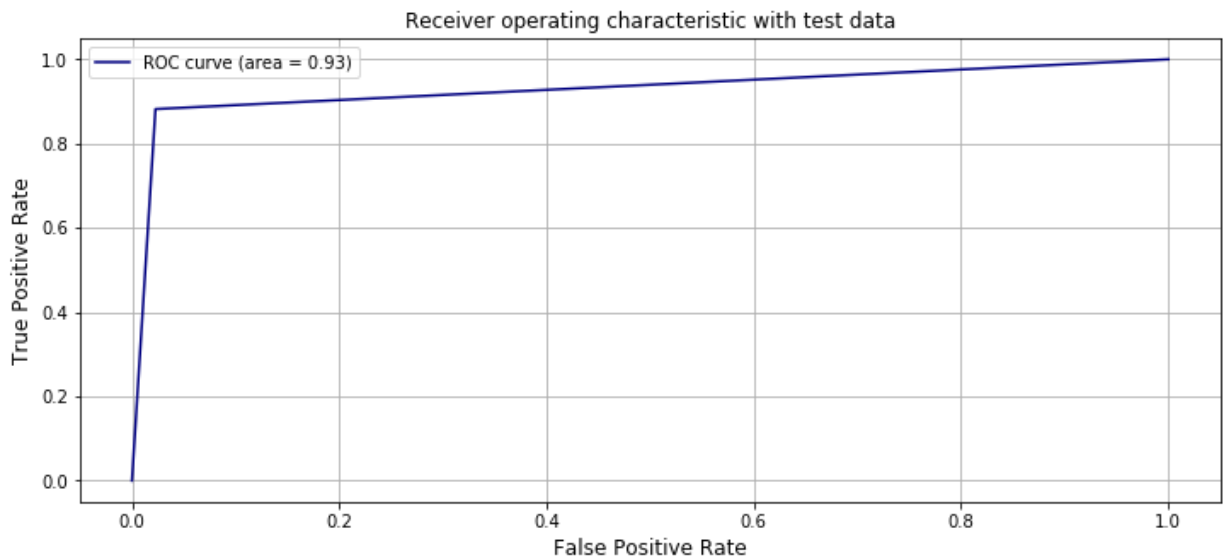
Train confusion_matrix



Test confusion_matrix



```
In [25]: 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.figure(figsize = (12,5))
5 plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
6 plt.grid()
7 plt.xlabel('False Positive Rate',fontsize = 12)
8 plt.ylabel('True Positive Rate',fontsize = 12)
9 plt.title('Receiver operating characteristic with test data',fontsize = 12)
10 plt.legend(loc = 'best')
11 plt.show()
```



Model 2 - XGBoost

3.1.1 Tuning the number of estimators in XGBoost

```

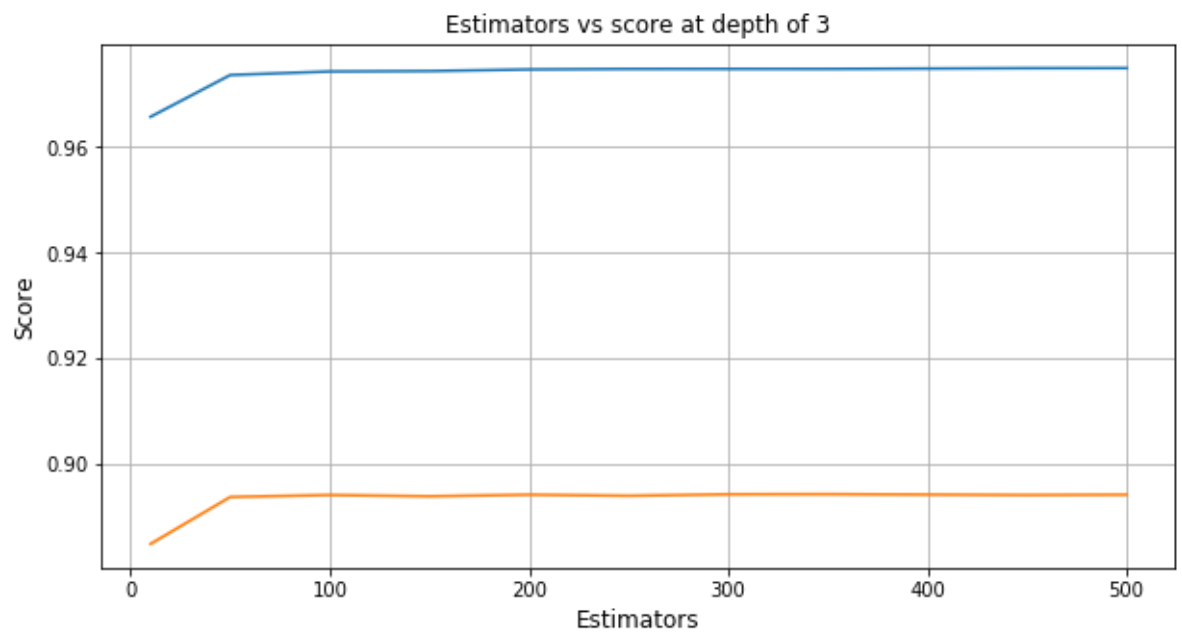
In [31]: 1 start = datetime.datetime.now()
2 estimators = [10,50,100,150,200,250,300,350,450,500]
3 train_scores = []
4 test_scores = []
5 for i in estimators:
6     clf = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1
7         colsample_bynode=1, colsample_bytree=0.3, eta=0.02, gamma=0,
8         learning_rate=0.7, max_delta_step=0, max_depth=3,
9         min_child_weight=1, missing=None, n_estimators=i, n_jobs=1,
10        nthread=None, objective='binary:logistic', random_state=0,
11        reg_alpha=93, reg_lambda=1, reg_lambda=77, scale_pos_weight=1,
12        seed=None, silent=None, subsample=0.6, verbosity=1)
13    clf.fit(df_final_train,y_train)
14    train_sc = f1_score(y_train,clf.predict(df_final_train))
15    test_sc = f1_score(y_test,clf.predict(df_final_test))
16    test_scores.append(test_sc)
17    train_scores.append(train_sc)
18    print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
19
20 plt.figure(figsize = (10,5))
21 plt.plot(estimators,train_scores,label='Train Score')
22 plt.plot(estimators,test_scores,label='Test Score')
23 plt.grid()
24 plt.xlabel('Estimators',fontsize = 12)
25 plt.ylabel('Score',fontsize = 12)
26 plt.title('Estimators vs score at depth of 3',fontsize = 12)
27 print('total time taken:',datetime.datetime.now() - start)

```

```

Estimators = 10 Train Score 0.9657533556878027 test Score 0.8849496002780675
Estimators = 50 Train Score 0.9736547367996206 test Score 0.8938228565203181
Estimators = 100 Train Score 0.9743248153082233 test Score 0.8941925912997949
Estimators = 150 Train Score 0.9744029715162404 test Score 0.8939512961508248
Estimators = 200 Train Score 0.9747113909744086 test Score 0.8942511181411584
Estimators = 250 Train Score 0.9747983972709198 test Score 0.894048840102132
Estimators = 300 Train Score 0.9748036583149947 test Score 0.8943245543214994
Estimators = 350 Train Score 0.974798905967724 test Score 0.8943511450381679
Estimators = 450 Train Score 0.9749500413798671 test Score 0.8941925912997949
Estimators = 500 Train Score 0.9749669555741659 test Score 0.8942536255588267
total time taken: 0:02:43.025765

```



3.2 Tuning the depth in the Classifier

```

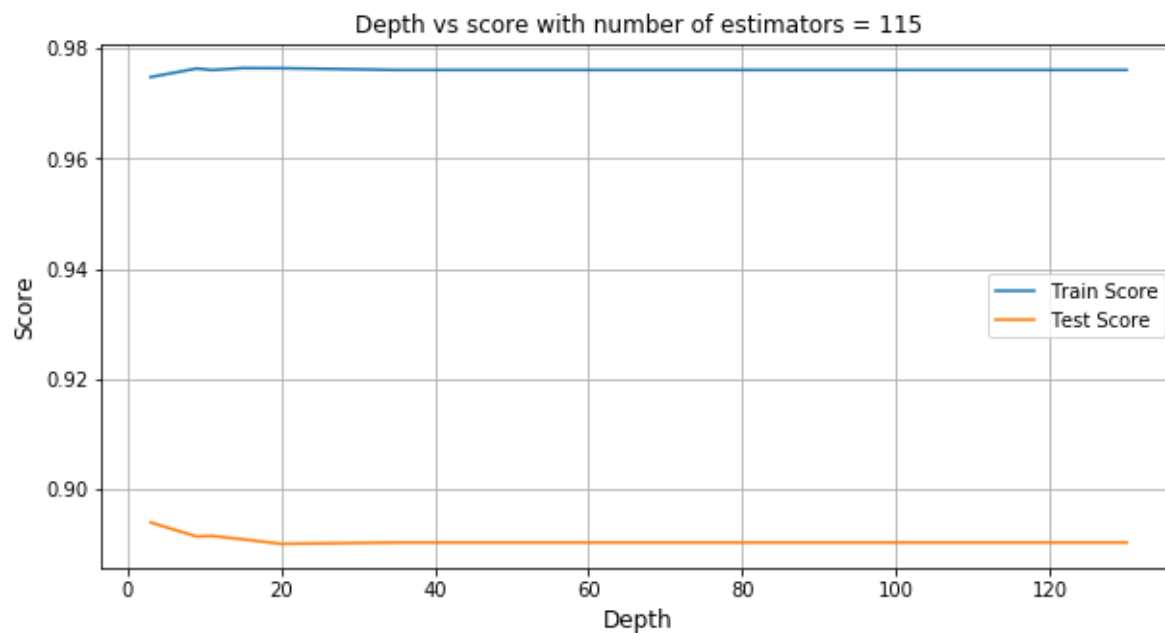
In [20]: 1 from xgboost import XGBClassifier
2 depths = [3,9,11,15,20,35,50,70,130]
3 train_scores = []
4 test_scores = []
5 for i in depths:
6     clf = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1
7         colsample_bynode=1, colsample_bytree=0.3, eta=0.02, gamma=0,
8         learning_rate=0.7, max_delta_step=0, max_depth=i,
9         min_child_weight=1, missing=None, n_estimators=250, n_jobs=1,
10        nthread=None, objective='binary:logistic', random_state=0,
11        reg_alpha=93, reg_lambda=1, reg_lambda=77, scale_pos_weight=1,
12        seed=None, silent=None, subsample=0.6, verbosity=1)
13    clf.fit(df_final_train,y_train)
14    train_sc = f1_score(y_train,clf.predict(df_final_train))
15    test_sc = f1_score(y_test,clf.predict(df_final_test))
16    test_scores.append(test_sc)
17    train_scores.append(train_sc)
18    print('depth = ',i,'Train Score',train_sc,'test Score',test_sc)
19
20 plt.figure(figsize = (10,5))
21 plt.plot(depths,train_scores,label='Train Score')
22 plt.plot(depths,test_scores,label='Test Score')
23 plt.grid()
24 plt.xlabel('Depth',fontsize = 12)
25 plt.ylabel('Score',fontsize = 12)
26 plt.title('Depth vs score with number of estimators = 115',fontsize = 12)
27 plt.legend(loc = 'best')
28 plt.show()

```

```

depth = 3 Train Score 0.9747983972709198 test Score 0.894048840102132
depth = 9 Train Score 0.9763506358603024 test Score 0.8915062386644231
depth = 11 Train Score 0.9760751059963659 test Score 0.8916099971574137
depth = 15 Train Score 0.9764268991684831 test Score 0.890990281061203
depth = 20 Train Score 0.9764059111312857 test Score 0.890168853069359
depth = 35 Train Score 0.976087658406651 test Score 0.8904166575438442
depth = 50 Train Score 0.976087658406651 test Score 0.8904166575438442
depth = 70 Train Score 0.976087658406651 test Score 0.8904166575438442
depth = 130 Train Score 0.976087658406651 test Score 0.8904166575438442

```

3.3 Tuning the XGBClassifier using Randomized Search Cross validation

```

In [21]: 1 #importing all the libraries
2 from scipy.stats import randint as sp_randint
3 from xgboost import XGBClassifier
4 from sklearn.model_selection import RandomizedSearchCV, StratifiedKFold
5
6 #hypertuning the model
7 n_estimators = [500,600,750,900,1100,1250]#tuning the number of estimators
8 learning_rate = [0.0001,0.003,0.05,0.7]#tuning the learning rate
9 colsample_bytree = [0.3,0.5,0.7]#sampling rate of the columns
10 subsample = [0.3,0.5,0.6,0.9]#subsampling rate
11 reg_alpha = sp_randint(0,600) #L1 regularization on weights
12 reg_lambda = sp_randint(0,600) #L2 regularization on weights
13
14 start = datetime.datetime.now()
15 param_grid = dict(learning_rate=learning_rate,
16                   n_estimators=n_estimators,
17                   colsample_bytree = colsample_bytree,
18                   subsample = subsample,
19                   reg_alpha = reg_alpha, reg_lambda = reg_lambda)
20
21 model = XGBClassifier(learning_rate = 0.1, objective='binary:logistic', eta = 0.1)
22 random_search = RandomizedSearchCV(model, param_distributions = param_grid, s
23
24
25 #training the model
26
27 #start = dt.datetime.now()
28 random_search.fit(df_final_train,y_train)
29
30 #print("\nTimeTaken: ",dt.datetime.now() - start)
31
32 #printing the best hyperparameters
33 print('Best hyperparameters are:', random_search.best_params_)
34 print('mean test scores', random_search.cv_results_['mean_test_score'])
35 print('mean train scores', random_search.cv_results_['mean_train_score'])
36 print('total time taken for computation is:', datetime.datetime.now() - start)

```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent worker s.

[Parallel(n_jobs=1)]: Done 50 out of 50 | elapsed: 60.4min finished

Best hyperparameters are: {'colsample_bytree': 0.3, 'learning_rate': 0.7, 'n_estimators': 750, 'reg_alpha': 93, 'reg_lambda': 77, 'subsample': 0.6}

mean test scores [0.92928731 0.93000827 0.96310668 0.86977193 0.96815021 0.92927881

0.96207715 0.97308008 0.91276454 0.91798703]

mean train scores [0.9294994 0.93008222 0.9636366 0.86998695 0.96844007 0.92936659

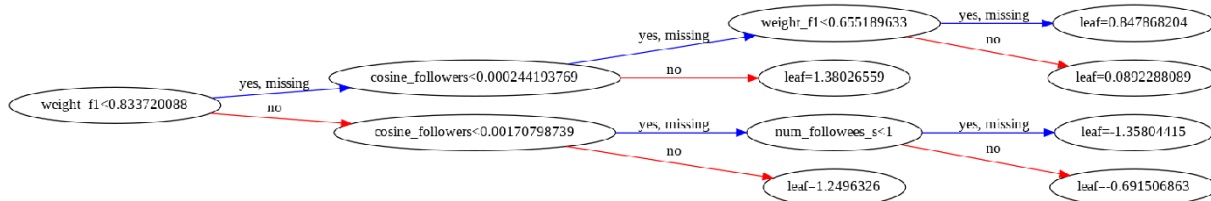
0.9623515 0.97410101 0.91275506 0.91800895]

total time taken for computation is: 1:01:05.735260

```
In [22]: 1 print('the best etimator is :',random_search.best_estimator_)
```

```
the best etimator is : XGBClassifier(base_score=0.5, booster='gbtree', colsampl
e_bylevel=1,
                                   colsample_bynode=1, colsample_bytree=0.3, eta=0.02, gamma=0,
                                   learning_rate=0.7, max_delta_step=0, max_depth=3,
                                   min_child_weight=1, missing=None, n_estimators=750, n_jobs=1,
                                   nthread=None, objective='binary:logistic', random_state=0,
                                   reg_alpha=93, reg_lambda=1, reg_lamda=77, scale_pos_weight=1,
                                   seed=None, silent=None, subsample=0.6, verbosity=1)
```

```
In [25]: 1 xgb = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
2                                   colsample_bynode=1, colsample_bytree=0.3, eta=0.02, gamma=0,
3                                   learning_rate=0.7, max_delta_step=0, max_depth=3,
4                                   min_child_weight=1, missing=None, n_estimators=750, n_jobs=1,
5                                   nthread=None, objective='binary:logistic', random_state=0,
6                                   reg_alpha=93, reg_lambda=1, reg_lamda=77, scale_pos_weight=1,
7                                   seed=None, silent=None, subsample=0.6, verbosity=1) #best class
8
9 xgb.fit(df_final_train,y_train)
10 from xgboost import plot_tree
11
12 plot_tree(xgb,rankdir='LR')#visualizes the bossted decision trees
13 fig = plt.gcf()
14 fig.set_size_inches(25,20)
15 fig.savefig('tree.png')
```



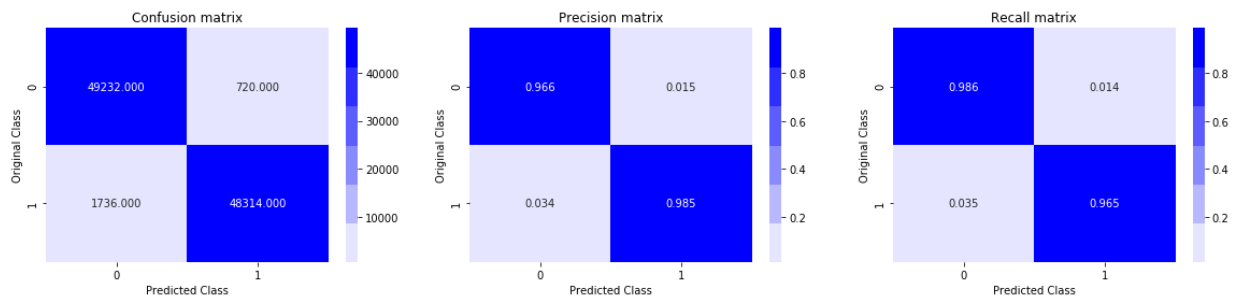
3.3.1 Confusion matrix and ROC curve

```
In [28]: 1 y_train_pred = xgb.predict(df_final_train)
2 y_test_pred = xgb.predict(df_final_test)
3 print('Train f1 score',f1_score(y_train,y_train_pred))
4 print('Test f1 score',f1_score(y_test,y_test_pred))
5
6 print('Train confusion_matrix')
7 plot_confusion_matrix(y_train,y_train_pred)
8 print('Test confusion_matrix')
9 plot_confusion_matrix(y_test,y_test_pred)
```

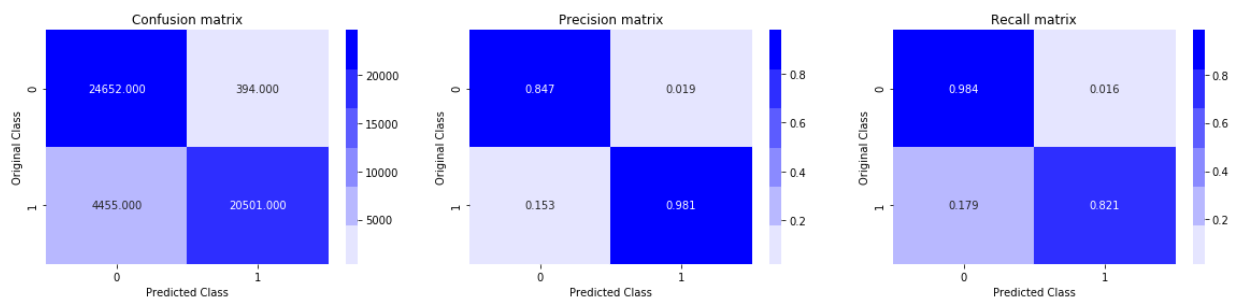
Train f1 score 0.9752129506277502

Test f1 score 0.8942444003402324

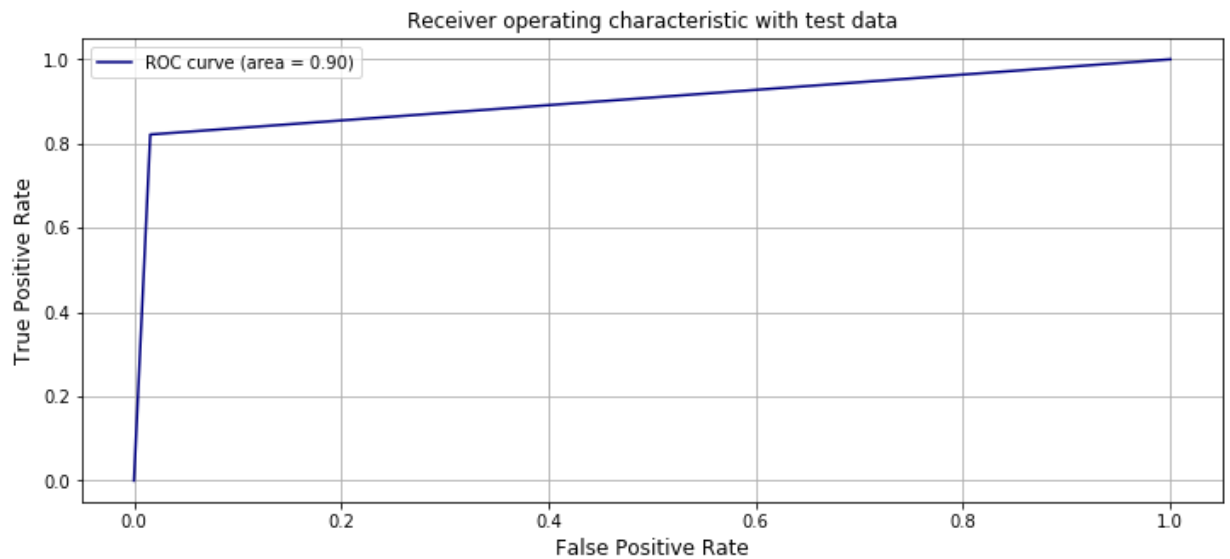
Train confusion_matrix



Test confusion_matrix



```
In [29]: 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.figure(figsize = (12,5))
5 plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
6 plt.grid()
7 plt.xlabel('False Positive Rate',fontsize = 12)
8 plt.ylabel('True Positive Rate',fontsize = 12)
9 plt.title('Receiver operating characteristic with test data',fontsize = 12)
10 plt.legend(loc = 'best')
11 plt.show()
```



Conclusion

```
In [11]: 1 from prettytable import PrettyTable
2
3 table_1 =PrettyTable()
4 table_1.field_names = ["Model", "Train F1", "Test F1"]
5 table_1.add_row(["Random Forest (tuning n_estimators)", 0.965, 0.924])
6 table_1.add_row(["Random Forest (tuning maximum depth)",0.962, 0.926])
7 table_1.add_row(["XGBoost (tuning n_estimators)", 0.974,0.894])
8 table_1.add_row(["XGBoost (tuning maximum depth)",0.974,0.894])
9 print('\t\tModels without hyperparameter tuning')
10 print(table_1)
11 print('\n\n')
12
13 table_2 =PrettyTable()
14 table_2.field_names = ["Model", "Train F1", "Test F1"]
15 table_2.add_row(["Random Forest", 0.963,0.926])
16 table_2.add_row(["XGboost", 0.975, 0.894])
17
18 print('Models with hyperparameter tuning using RandomizedCV')
19 print(table_2)
```

Models without hyperparameter tuning

Model	Train F1	Test F1
Random Forest (tuning n_estimators)	0.965	0.924
Random Forest (tuning maximum depth)	0.964	0.926
XGBoost (tuning n_estimators)	0.974	0.894
XGBoost (tuning maximum depth)	0.974	0.894

Models with hyperparameter tuning using RandomizedCV

Model	Train F1	Test F1
Random Forest	0.963	0.926
XGboost	0.975	0.894

In this case we are trying to find missing links which we will use for recommending friends or followers or in general recommending connections. Facebook created this to solve their business problem. The dataset is fairly reasonably sized to get a feel and understanding of how Network Analysis works with Machine Learning in real world scenarios. Facebook has given this problem based on directed graphs which represent users on social network. The dataset only has source node information and destination node information. Here each node represents a user. Basically a link between source to destination means that a user follows another user.

In the given data, only those source and destination nodes are given for which an edge exists. There is no information about the nodes which does not have an edge between them. So, in order to map this problem to a binary classification problem of whether or not an edge exists in the graph, we need to create training and testing sample which has a class label of 0 (0 means that there are no edges present between source to destination).

NOTE: In the given dataset, we have roughly 9.43 million edges and 1.93 million nodes (vertices or users). For the given data, all the links are present and hence the class label will be 1. However, for classification we also need 0 class labels. How do we generate the 0 class labels?

Create the same number of 0 labeled pairs of vertices that we have for class 1 labeled pairs of vertices. Randomly sample a pair of vertices. Check if the path length is greater than 2. Check if no edge connection exists between the pairs of vertices. If both the above conditions are satisfied then we will have a new pair of edges which will have a class label 0. Coming to business constraints, there are no low latency requirements. We need to use probability estimates to predict edges between two nodes. This will give us more interpretability in terms of which edge connections are more important. The metric we have chosen is F1 score and binary confusion matrix.

We curated features like number of followers and followees of each node, whether or not a node is followed back by any other nodes, page rank of individual nodes, katz score, jaccard index, preferential attachment, svd features, svd dot features, adar index and so on. There were a total of 59 features on which we train and test our model.

There is not timestamp provided for this data. Ideally, if you think about it, we have the dataset for a given time step t . However, the graph is evolving and changing over time. After 30 days the edge connections might change, because people might have new followers and they may even start to follow new people. In the real world, we would split the data according to time. But, since we do not have any information about time stamp, we will split the data randomly in 80:20 ratio. 80% for training data and 20% for cross validation data.

Hyperparameter tuning using Randomized Searchc cross validaion was done for both the models and our f1 score imporved by a small margin for Randomr forest model ,where we got f1 score for training data 0.963 and for test data the score was 0.926

Another important observation was that the while calculating the feature importances we find that the feature 'follows_back' is most important in almost all of the models