

# Social network Graph Link Prediction - Facebook Challenge

```
In [1]: # Importing Libraries
# please do go through this python notebook:
import warnings
warnings.filterwarnings("ignore")

import csv
import pandas as pd#pandas to create small dataframes
import datetime #Convert to unix time
import time #Convert to unix time
# if numpy is not installed already : pip3 install numpy
import numpy as np#Do arithmetic operations on arrays
# matplotlib: used to plot graphs
import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots
from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
import math
import pickle
import os
# to install xgboost: pip3 install xgboost
import xgboost as xgb

import warnings
import networkx as nx
import pdb
import pickle
from pandas import HDFStore, DataFrame
from pandas import read_hdf
from scipy.sparse.linalg import svds, eigs
import gc
from tqdm import tqdm
from sklearn.ensemble import RandomForestClassifier
```

```
In [2]: #reading
from pandas import read_hdf
df_final_train = read_hdf('C:/Users/PareshBhatia/Downloads/Learning/Data Science/analytics vidya/facebook/d
df_final_test = read_hdf('C:/Users/PareshBhatia/Downloads/Learning/Data Science/analytics vidya/facebook/d
```

```
In [3]: Out[3]: Index(['source_node', 'destination_node', 'indicator_link',
      'jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followers_d',
      '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', 'svd_u_dot', 'svd_v_dot'],
      dtype='object')
```

```
In [4]: y_train = df_final_train.indicator_link
```

```
In [5]: df_final_train.drop(['source_node', 'destination_node', 'indicator_link'], axis=1, inplace=True)
```

## Random Forest Model

```

In [6]: estimators = [10,50,100,250,450]
train_scores = []
test_scores = []
for i in estimators:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=5, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=52, min_samples_split=120,
                                min_weight_fraction_leaf=0.0, n_estimators=i, n_jobs=-1, random_state=25, verbose=0, warm_start=False)
    clf.fit(df_final_train, y_train)
    train_sc = f1_score(y_train, clf.predict(df_final_train))
    test_sc = f1_score(y_test, clf.predict(df_final_test))
    test_scores.append(test_sc)
    train_scores.append(train_sc)
    print('Estimators = ', i, 'Train Score', train_sc, 'test Score', test_sc)
plt.plot(estimators, train_scores, label='Train Score')
plt.plot(estimators, test_scores, label='Test Score')
plt.xlabel('Estimators')
plt.ylabel('Score')

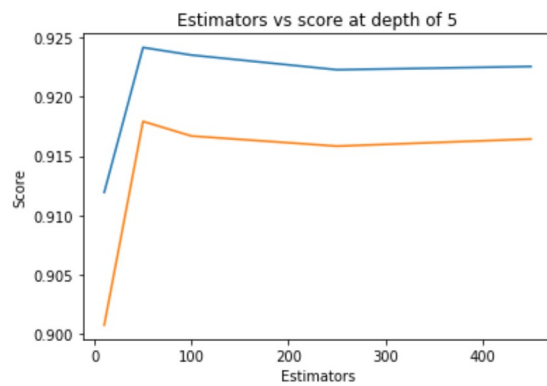
```

```

Estimators = 10 Train Score 0.9119491648210627 test Score 0.9007777220270948
Estimators = 50 Train Score 0.9241336247268187 test Score 0.9179152957941281
Estimators = 100 Train Score 0.9234945962532636 test Score 0.9166841289132895
Estimators = 250 Train Score 0.9222600493487836 test Score 0.9158345221112696
Estimators = 450 Train Score 0.9225338914225025 test Score 0.91643062461971

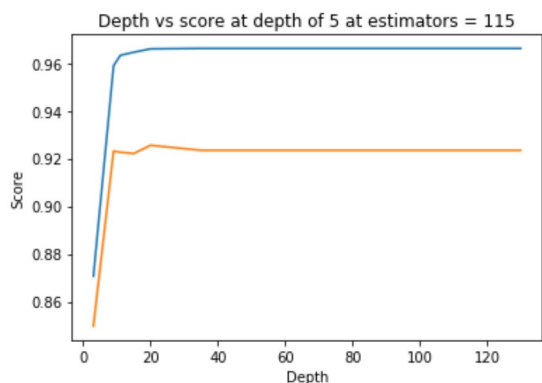
```

Out[6]: Text(0.5,1,'Estimators vs score at depth of 5')



```
In [7]: depths = [3,9,11,15,20,35,50,70,130]
train_scores = []
test_scores = []
for i in depths:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=i, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=52, min_samples_split=120,
                                min_weight_fraction_leaf=0.0, n_estimators=115, n_jobs=-1, random_state=25, verbose=0, warm_start=True)
    clf.fit(df_final_train, y_train)
    train_sc = f1_score(y_train, clf.predict(df_final_train))
    test_sc = f1_score(y_test, clf.predict(df_final_test))
    test_scores.append(test_sc)
    train_scores.append(train_sc)
    print('depth = ', i, 'Train Score', train_sc, 'test Score', test_sc)
plt.plot(depths, train_scores, label='Train Score')
plt.plot(depths, test_scores, label='Test Score')
plt.xlabel('Depth')
plt.ylabel('Score')
plt.title('Depth vs score at depth of 5 at estimators = 115')
```

```
depth = 3 Train Score 0.8706410024764314 test Score 0.8497735006610747
depth = 9 Train Score 0.95889991737139 test Score 0.9231158854441467
depth = 11 Train Score 0.963235741676304 test Score 0.9226263222217541
depth = 15 Train Score 0.9645529178630646 test Score 0.9221057947228886
depth = 20 Train Score 0.9659667054596974 test Score 0.9255337031433856
depth = 35 Train Score 0.9661550291906549 test Score 0.9234331908801919
depth = 50 Train Score 0.9661550291906549 test Score 0.9234331908801919
depth = 70 Train Score 0.9661550291906549 test Score 0.9234331908801919
depth = 130 Train Score 0.9661550291906549 test Score 0.9234331908801919
```



```
In [33]: from sklearn.metrics import f1_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint
from scipy.stats import uniform

param_dist = {"n_estimators": sp_randint(105, 125),
              "max_depth": sp_randint(10, 15),
              "min_samples_split": sp_randint(110, 190),
              "min_samples_leaf": sp_randint(25, 65)}

clf = RandomForestClassifier(random_state=25, n_jobs=-1)

rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
                               n_iter=5, cv=10, scoring='f1', random_state=25)

rf_random.fit(df_final_train, y_train)

mean test scores [0.96433296 0.96384632 0.9615574 0.96388255 0.96490786]
```

```
In [34]: 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 [11]:
In [12]: clf.fit(df_final_train,y_train)
          y_train_pred = clf.predict(df_final_train)
In [13]: from sklearn.metrics import f1_score
          print('Train f1 score',f1_score(y_train,y_train_pred))
```

Train f1 score 0.9659810206706569  
Test f1 score 0.9223890669816095

```
In [14]: from sklearn.metrics import confusion_matrix
          def plot_confusion_matrix(test_y, predict_y):
              C = confusion_matrix(test_y, predict_y)

              A = ((C.T)/(C.sum(axis=1))).T

              B = (C/C.sum(axis=0))
              plt.figure(figsize=(20,4))

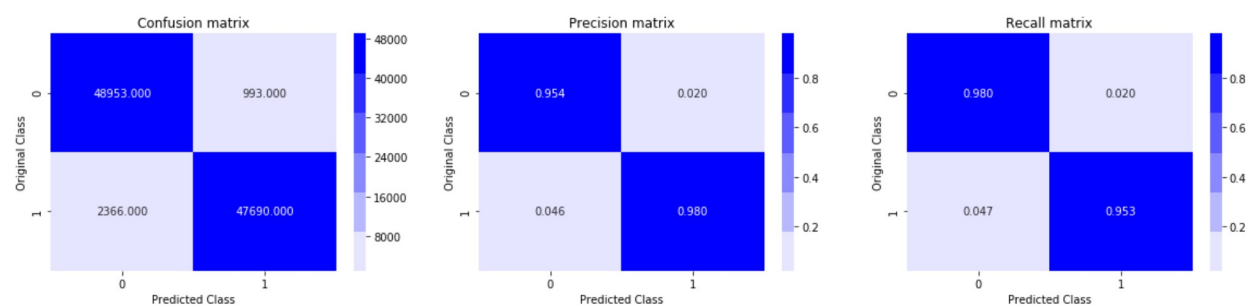
              labels = [0,1]
              # representing A in heatmap format
              cmap=sns.light_palette("blue")
              plt.subplot(1, 3, 1)
              sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
              plt.xlabel('Predicted Class')
              plt.ylabel('Original Class')
              plt.title("Confusion matrix")

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

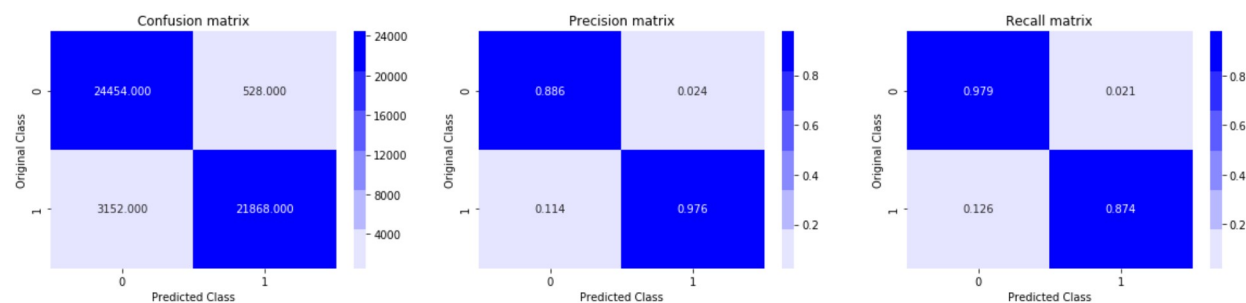
              plt.subplot(1, 3, 3)
              # representing B in heatmap format
              sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
              plt.xlabel('Predicted Class')
              plt.ylabel('Original Class')
              plt.title("Recall matrix")
```

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

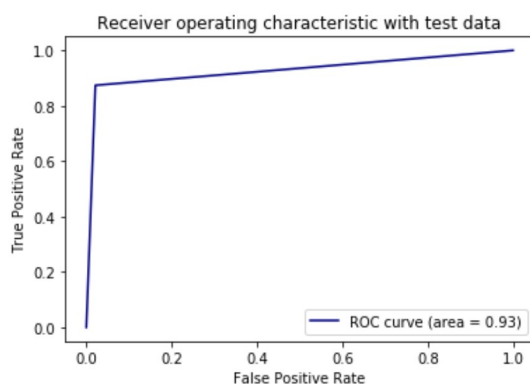
Train confusion\_matrix



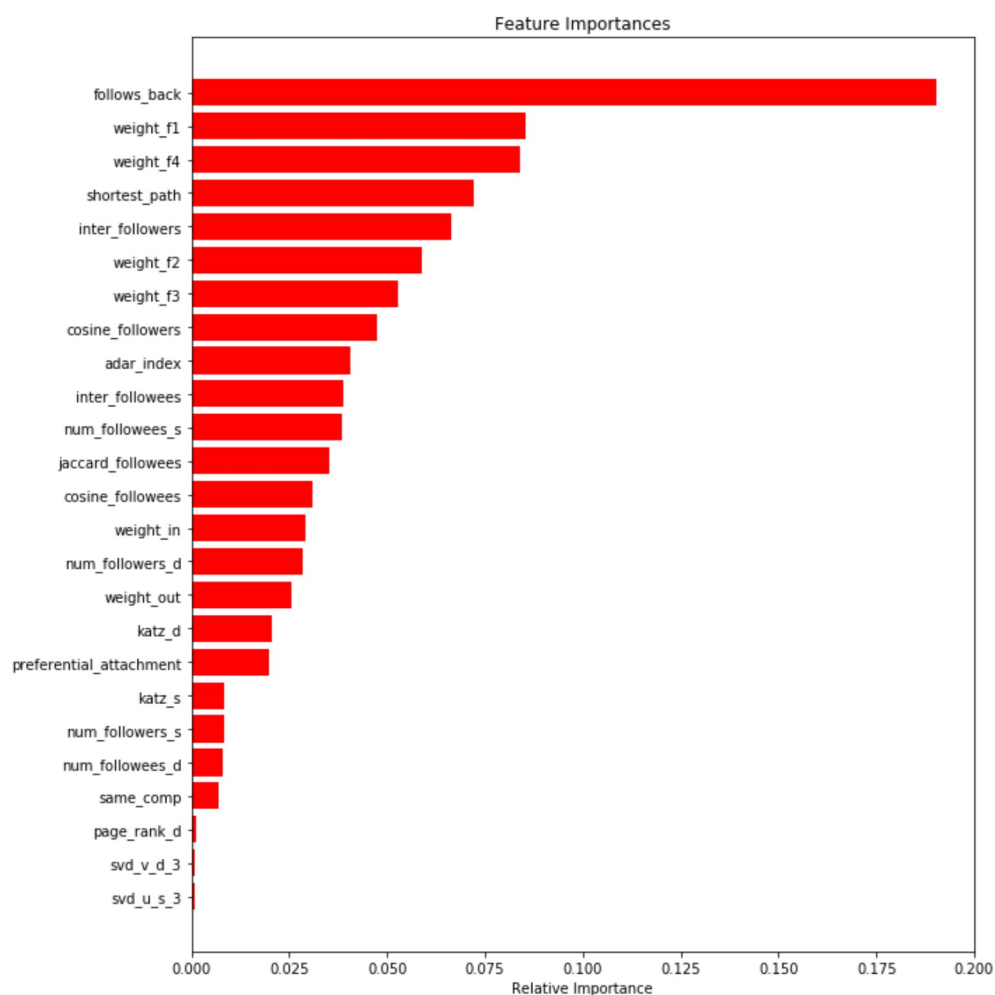
Test confusion\_matrix



```
In [16]: from sklearn.metrics import roc_curve, auc
fpr,tpr,ths = roc_curve(y_test,y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic with test data')
plt.legend()
```



```
In [17]: features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-25:]
plt.figure(figsize=(10,12))
plt.title('Feature Importances')
plt.barh(range(len(indices)), importances[indices], color='r', align='center')
plt.yticks(range(len(indices)), [features[i] for i in indices])
plt.xlabel('Relative Importance')
```



## XGboost model

```
In [22]: %timeit
from xgboost import XGBClassifier, train, DMatrix
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import StratifiedKFold
# params = {}
# params['objective'] = 'binary:logistic'
# params['eval_metric'] = 'logloss'
# params['eta'] = 0.02
# params['max_depth'] = 4

# d_train = xgb.DMatrix(X_train, label=y_train)
# d_test = xgb.DMatrix(X_test, label=y_test)

# watchlist = [(d_train, 'train'), (d_test, 'valid')]

# bst = xgb.train(params, d_train, 400, watchlist, early_stopping_rounds=20, verbose_eval=10)

# xgdmatrix = xgb.DMatrix(X_train, y_train)
# predict_y = bst.predict(d_test)
# print("The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
#####

# A parameter grid for XGBoost

params = {
    'min_child_weight': [1, 5, 10],
    'gamma': [0.5, 1, 1.5, 2, 5],
    'subsample': [0.6, 0.8, 1.0],
    'colsample_bytree': [0.6, 0.8, 1.0],
    'max_depth': [4, 5, 6, 8]
}

xgb = XGBClassifier(learning_rate=0.02, n_estimators=1000, objective='binary:logistic',
                    silent=True, n_thread=2)

folds = 5 # number of folds

skf = StratifiedKFold(n_splits=folds, shuffle = True, random_state = 1001)

random_search = RandomizedSearchCV(xgb, param_distributions=params, scoring='f1', n_jobs=4, cv=skf.split(d:

# Here we go
# start_time = timer(None) # timing starts from this point for "start_time" variable
random_search.fit(df_final_train, y_train)

print('\n best model :')
print(random_search.best_estimator_)

print('\n Best hyperparameters:')
print(random_search.best_params_)

print('\n Best f1 score:')
print(random_search.best_score_)
```

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

```
[Parallel(n_jobs=4)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=4)]: Done 24 tasks | elapsed: 64.9min
[Parallel(n_jobs=4)]: Done 50 out of 50 | elapsed: 133.5min finished
```

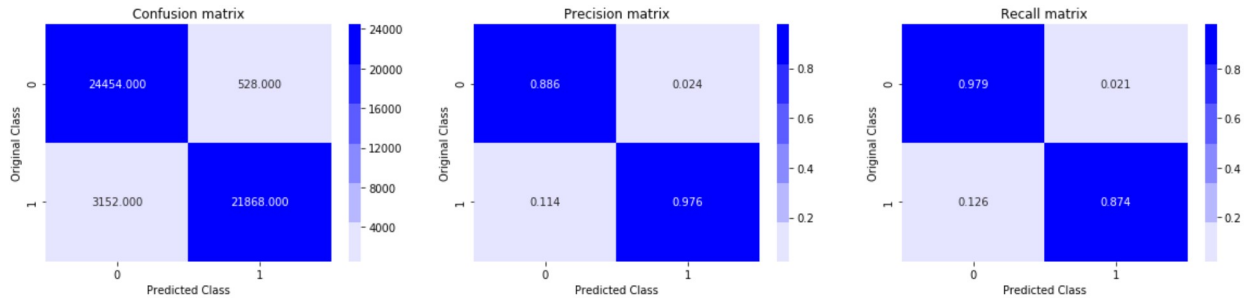
```
best model :
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bytree=0.8, gamma=2, learning_rate=0.02,
              max_delta_step=0, max_depth=8, min_child_weight=5, missing=None,
              n_estimators=1000, n_jobs=1, n_thread=2, nthread=None,
              objective='binary:logistic', random_state=0, reg_alpha=0,
              reg_lambda=1, scale_pos_weight=1, seed=None, silent=True,
              subsample=1.0)
```

```
Best hyperparameters:
{'subsample': 1.0, 'min_child_weight': 5, 'max_depth': 8, 'gamma': 2, 'colsample_bytree': 0.8}
```

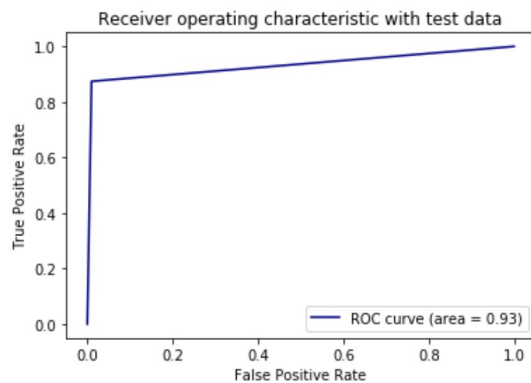
```
Best f1 score:
0.9823191083629064
```

```
In [28]: # train model using best parameters given by random search
xgb = random_search.best_estimator_
xgb.fit(df_final_train,y_train)
y_pred_train = xgb.predict(df_final_train)
print("The train f1 score is:",f1_score(y_train, y_pred_train))
y_pred_test = xgb.predict(df_final_test)
print("The test f1 score is:",f1_score(y_test, y_pred_test))
```

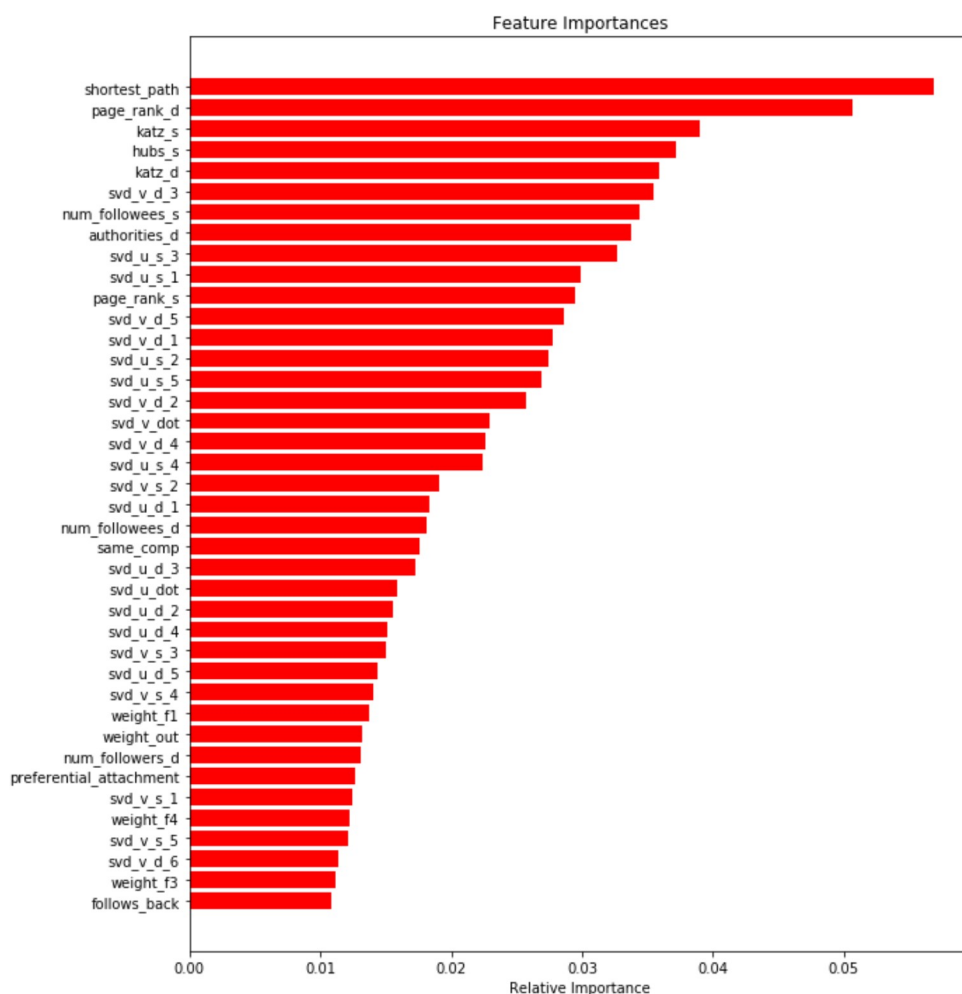
The train f1 score is: 0.992600761981151  
The test f1 score is: 0.9281486667090943



```
In [30]: fpr,tpr,ths = roc_curve(y_test,y_pred_test)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic with test data')
plt.legend()
```



```
In [32]: ▶ features = df_final_train.columns
importances = xgb.feature_importances_
indices = (np.argsort(importances))[-40:]
plt.figure(figsize=(10,12))
plt.title('Feature Importances')
plt.barh(range(len(indices)), importances[indices], color='r', align='center')
plt.yticks(range(len(indices)), [features[i] for i in indices])
plt.xlabel('Relative Importance')
```



## Conclusion

1. Preferential attachment proves to be important feature in random forest model.
2. Shortest path is most important feature for XGBoost.
3. follows\_back is most important feature for RandomForest model.
4. XgBoost performs marginally better than random forest with best F1 score of 0.9281.

In [ ]: ▶