

# Social network Graph Link Prediction - Facebook Challenge

```
In [8]: #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
from sklearn.metrics import f1_score
```

## 1. Reading Data

```
In [11]: startTime = datetime.datetime.now()
print("Current Time = ",startTime)

if os.path.isfile('data/after_eda/train_pos_after_eda.csv'):
    train_graph=nx.read_edgelist('data/after_eda/train_pos_after_eda.csv',delimit
    print(nx.info(train_graph))
else:
    print("please run the FB_EDA.ipynb or download the files from drive")

print("\nTime taken for creation of dataframe is {}".format(datetime.datetime.now
```

Current Time = 2019-05-24 21:09:06.868196

Name:

Type: DiGraph

Number of nodes: 1780722

Number of edges: 7550015

Average in degree: 4.2399

Average out degree: 4.2399

Time taken for creation of dataframe is 0:04:00.803636

## 2. Similarity measures

### 2.1 Jaccard Distance:

<http://www.statisticshowto.com/jaccard-index/> (<http://www.statisticshowto.com/jaccard-index/>)

$$j = \frac{|X \cap Y|}{|X \cup Y|}$$

```
In [0]: #for followees
def jaccard_for_followees(a,b):
    try:
        if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.succes
            return 0
        sim = (len(set(train_graph.successors(a)).intersection(set(train_graph.su
            (len(set(train_graph.successors(a)).union(set
    except:
        return 0
    return sim
```

```
In [0]: #one test case
print(jaccard_for_followees(273084,1505602))

0.0
```

```
In [0]: #node 1635354 not in graph
print(jaccard_for_followees(273084,1505602))

0.0
```

```
In [0]: #for followers
def jaccard_for_followers(a,b):
    try:
        if len(set(train_graph.predecessors(a))) == 0 | len(set(g.predecessors(b))) == 0:
            return 0
        sim = (len(set(train_graph.predecessors(a)).intersection(set(train_graph.predecessors(b)))) /
                (len(set(train_graph.predecessors(a)).union(set(train_graph.predecessors(b))))))
        return sim
    except:
        return 0
```

```
In [0]: print(jaccard_for_followers(273084,470294))

0
```

```
In [0]: #node 1635354 not in graph
print(jaccard_for_followers(669354,1635354))

0
```

## 2.2 Cosine distance

$$\text{CosineDistance} = \frac{|X \cap Y|}{|X| \cdot |Y|}$$

```
In [0]: #for followees
def cosine_for_followees(a,b):
    try:
        if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.successors(b))) == 0:
            return 0
        sim = (len(set(train_graph.successors(a)).intersection(set(train_graph.successors(b)))) /
                (math.sqrt(len(set(train_graph.successors(a))) * len(set(train_graph.successors(b)))))
        return sim
    except:
        return 0
```

```
In [0]: print(cosine_for_followees(273084,1505602))

0.0
```

```
In [0]: print(cosine_for_followees(273084,1635354))

0
```

```
In [0]: def cosine_for_followers(a,b):
        try:

            if len(set(train_graph.predecessors(a))) == 0 | len(set(train_graph.pred
                return 0
            sim = (len(set(train_graph.predecessors(a)).intersection(set(train_graph.
                (math.sqrt(len(set(train_graph.predecessors(
            return sim
        except:
            return 0
```

```
In [0]: print(cosine_for_followers(2,470294))
```

```
0.02886751345948129
```


```
In [0]: print(cosine_for_followers(669354,1635354))
```

```
0
```

### 3. Ranking Measures

[https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx.algorithms.link\\_analysis.pagerank\\_alg.pagerank.html](https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx.algorithms.link_analysis.pagerank_alg.pagerank.html)  
[https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx.algorithms.link\\_analysis.pagerank\\_alg.pagerank.html](https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx.algorithms.link_analysis.pagerank_alg.pagerank.html)

PageRank computes a ranking of the nodes in the graph  $G$  based on the structure of the incoming links.

 Mathematical PageRanks for a simple network, expressed as percentages. (Google uses a logarithmic scale.) Page C has a higher PageRank than Page E, even though there are fewer links to C; the one link to C comes from an important page and hence is of high value. If web surfers who start on a random page have an 85% likelihood of choosing a random link from the page they are currently visiting, and a 15% likelihood of jumping to a page chosen at random from the entire web, they will reach Page E 8.1% of the time. **(The 15% likelihood of jumping to an arbitrary page corresponds to a damping factor of 85%.) Without damping, all web surfers would eventually end up on Pages A, B, or C, and all other pages would have PageRank zero. In the presence of damping, Page A effectively links to all pages in the web, even though it has no outgoing links of its own.**

#### 3.1 Page Ranking

<https://en.wikipedia.org/wiki/PageRank> (<https://en.wikipedia.org/wiki/PageRank>)

```
In [0]: if not os.path.isfile('data/fea_sample/page_rank.p'):
        pr = nx.pagerank(train_graph, alpha=0.85)
        pickle.dump(pr, open('data/fea_sample/page_rank.p', 'wb'))
    else:
        pr = pickle.load(open('data/fea_sample/page_rank.p', 'rb'))
```

```
In [0]: print('min', pr[min(pr, key=pr.get)])
        print('max', pr[max(pr, key=pr.get)])
        print('mean', float(sum(pr.values())) / len(pr))
```

```
min 1.6556497245737814e-07
max 2.7098251341935827e-05
mean 5.615699699389075e-07
```

```
In [0]: #for imputing to nodes which are not there in Train data
        mean_pr = float(sum(pr.values())) / len(pr)
        print(mean_pr)
```

```
5.615699699389075e-07
```

## 4. Other Graph Features

### 4.1 Shortest path:

Getting Shortest path between two nodes, if nodes have direct path i.e directly connected then we are removing that edge and calculating path.

```
In [0]: #if has direct edge then deleting that edge and calculating shortest path
        def compute_shortest_path_length(a,b):
            p=-1
            try:
                if train_graph.has_edge(a,b):
                    train_graph.remove_edge(a,b)
                    p= nx.shortest_path_length(train_graph,source=a,target=b)
                    train_graph.add_edge(a,b)
            else:
                p= nx.shortest_path_length(train_graph,source=a,target=b)
            return p
        except:
            return -1
```

```
In [0]: #testing
        compute_shortest_path_length(77697, 826021)
```

```
Out[21]: 10
```

```
In [0]: #testing
        compute_shortest_path_length(669354,1635354)
```

```
Out[22]: -1
```

## 4.2 Checking for same community

```
In [0]: #getting weekly connected edges from graph
wcc=list(nx.weakly_connected_components(train_graph))
def belongs_to_same_wcc(a,b):
    index = []
    if train_graph.has_edge(b,a):
        return 1
    if train_graph.has_edge(a,b):
        for i in wcc:
            if a in i:
                index=i
                break
        if (b in index):
            train_graph.remove_edge(a,b)
            if compute_shortest_path_length(a,b)==-1:
                train_graph.add_edge(a,b)
                return 0
            else:
                train_graph.add_edge(a,b)
                return 1
        else:
            return 0
    else:
        for i in wcc:
            if a in i:
                index=i
                break
        if(b in index):
            return 1
        else:
            return 0
```

```
In [0]: belongs_to_same_wcc(861, 1659750)
```

```
Out[24]: 0
```

```
In [0]: belongs_to_same_wcc(669354,1635354)
```

```
Out[25]: 0
```

## 4.3 Adamic/Adar Index:

Adamic/Adar measures is defined as inverted sum of degrees of common neighbours for given two vertices.

$$A(x, y) = \sum_{u \in N(x) \cap N(y)} \frac{1}{\log(\ln(u))}$$

```
In [0]: #adar index
def calc_adar_in(a,b):
    sum=0
    try:
        n=list(set(train_graph.successors(a)).intersection(set(train_graph.succes
        if len(n)!=0:
            for i in n:
                sum=sum+(1/np.log10(len(list(train_graph.predecessors(i)))))
            return sum
        else:
            return 0
    except:
        return 0
```

```
In [0]: calc_adar_in(1,189226)
```

```
Out[27]: 0
```

```
In [0]: calc_adar_in(669354,1635354)
```

```
Out[28]: 0
```

## 4.4 Is person was following back:

```
In [0]: def follows_back(a,b):
        if train_graph.has_edge(b,a):
            return 1
        else:
            return 0
```

```
In [0]: follows_back(1,189226)
```

```
Out[30]: 1
```

```
In [0]: follows_back(669354,1635354)
```

```
Out[31]: 0
```

## 4.5 Katz Centrality:

[https://en.wikipedia.org/wiki/Katz\\_centrality\\_\(https://en.wikipedia.org/wiki/Katz\\_centrality\)](https://en.wikipedia.org/wiki/Katz_centrality_(https://en.wikipedia.org/wiki/Katz_centrality))

<https://www.geeksforgeeks.org/katz-centrality-centrality-measure/>

(<https://www.geeksforgeeks.org/katz-centrality-centrality-measure/>) Katz centrality computes the centrality for a node based on the centrality of its neighbors. It is a generalization of the eigenvector centrality. The Katz centrality for node  $i$  is

$$x_i = \alpha \sum_j A_{ij} x_j + \beta,$$

where  $A$  is the adjacency matrix of the graph  $G$  with eigenvalues

$$\lambda$$

The parameter

$$\beta$$

controls the initial centrality and

$$\alpha < \frac{1}{\lambda_{\max}}.$$

```
In [0]: if not os.path.isfile('data/fea_sample/katz.p'):
        katz = nx.katz_katz centrality(train_graph,alpha=0.005,beta=1)
        pickle.dump(katz,open('data/fea_sample/katz.p','wb'))
    else:
        katz = pickle.load(open('data/fea_sample/katz.p','rb'))
```

```
In [0]: print('min',katz[min(katz, key=katz.get)])
        print('max',katz[max(katz, key=katz.get)])
        print('mean',float(sum(katz.values())) / len(katz))
```

```
min 0.0007313532484065916
max 0.003394554981699122
mean 0.0007483800935562018
```

```
In [0]: mean_katz = float(sum(katz.values())) / len(katz)
        print(mean_katz)
```

```
0.0007483800935562018
```

## 4.6 Hits Score

The HITS algorithm computes two numbers for a node. Authorities estimates the node value based on the incoming links. Hubs estimates the node value based on outgoing links.

[https://en.wikipedia.org/wiki/HITS\\_algorithm](https://en.wikipedia.org/wiki/HITS_algorithm) ([https://en.wikipedia.org/wiki/HITS\\_algorithm](https://en.wikipedia.org/wiki/HITS_algorithm))

```
In [0]: if not os.path.isfile('data/fea_sample/hits.p'):
        hits = nx.hits(train_graph, max_iter=100, tol=1e-08, nstart=None, normalized=True)
        pickle.dump(hits,open('data/fea_sample/hits.p','wb'))
    else:
        hits = pickle.load(open('data/fea_sample/hits.p','rb'))
```

```
In [0]: print('min',hits[0][min(hits[0], key=hits[0].get)])
        print('max',hits[0][max(hits[0], key=hits[0].get)])
        print('mean',float(sum(hits[0].values())) / len(hits[0]))
```

```
min 0.0
max 0.004868653378780953
mean 5.615699699344123e-07
```

## 5. Featurization



## 5. 1 Reading a sample of Data from both train and test

```
In [0]: import random
if os.path.isfile('data/after_eda/train_after_eda.csv'):
    filename = "data/after_eda/train_after_eda.csv"
    # you uncomment this line, if you dont know the length of the file name
    # here we have hardcoded the number of lines as 15100030
    # n_train = sum(1 for line in open(filename)) #number of records in file (excl
    n_train = 15100028
    s = 100000 #desired sample size
    skip_train = sorted(random.sample(range(1,n_train+1),n_train-s))
    #https://stackoverflow.com/a/22259008/4084039
```

```
In [0]: if os.path.isfile('data/after_eda/train_after_eda.csv'):
    filename = "data/after_eda/test_after_eda.csv"
    # you uncomment this line, if you dont know the length of the file name
    # here we have hardcoded the number of lines as 3775008
    # n_test = sum(1 for line in open(filename)) #number of records in file (excl
    n_test = 3775006
    s = 50000 #desired sample size
    skip_test = sorted(random.sample(range(1,n_test+1),n_test-s))
    #https://stackoverflow.com/a/22259008/4084039
```

```
In [0]: print("Number of rows in the train data file:", n_train)
print("Number of rows we are going to eliminate in train data are",len(skip_train))
print("Number of rows in the test data file:", n_test)
print("Number of rows we are going to eliminate in test data are",len(skip_test))
```

Number of rows in the train data file: 15100028  
 Number of rows we are going to eliminate in train data are 15000028  
 Number of rows in the test data file: 3775006  
 Number of rows we are going to eliminate in test data are 3725006

```
In [0]: df_final_train = pd.read_csv('data/after_eda/train_after_eda.csv', skiprows=skip_train)
df_final_train['indicator_link'] = pd.read_csv('data/train_y.csv', skiprows=skip_train)
print("Our train matrix size ",df_final_train.shape)
df_final_train.head(2)
```

Our train matrix size (100002, 3)

```
Out[49]:
```

	source_node	destination_node	indicator_link
0	273084	1505602	1
1	832016	1543415	1

```
In [0]: df_final_test = pd.read_csv('data/after_eda/test_after_eda.csv', skiprows=skip_te
df_final_test['indicator_link'] = pd.read_csv('data/test_y.csv', skiprows=skip_te
print("Our test matrix size ",df_final_test.shape)
df_final_test.head(2)
```

Our test matrix size (50002, 3)

```
Out[50]:
```

	source_node	destination_node	indicator_link
0	848424	784690	1
1	483294	1255532	1

## 5.2 Adding a set of features

we will create these each of these features for both train and test data points

1. jaccard\_followers
2. jaccard\_followees
3. cosine\_followers
4. cosine\_followees
5. num\_followers\_s
6. num\_followees\_s
7. num\_followers\_d
8. num\_followees\_d
9. inter\_followers
10. inter\_followees

```
In [0]: if not os.path.isfile('data/fea_sample/storage_sample_stage1.h5'):
        #mapping jaccrd followers to train and test data
        df_final_train['jaccard_followers'] = df_final_train.apply(lambda row:
                                                                    jaccard_for_followers(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)
        df_final_test['jaccard_followers'] = df_final_test.apply(lambda row:
                                                                    jaccard_for_followers(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)

        #mapping jaccrd followees to train and test data
        df_final_train['jaccard_followees'] = df_final_train.apply(lambda row:
                                                                    jaccard_for_followees(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)
        df_final_test['jaccard_followees'] = df_final_test.apply(lambda row:
                                                                    jaccard_for_followees(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)

        #mapping cosine followers to train and test data
        df_final_train['cosine_followers'] = df_final_train.apply(lambda row:
                                                                    cosine_for_followers(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)
        df_final_test['cosine_followers'] = df_final_test.apply(lambda row:
                                                                    cosine_for_followers(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)

        #mapping cosine followees to train and test data
        df_final_train['cosine_followees'] = df_final_train.apply(lambda row:
                                                                    cosine_for_followees(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)
        df_final_test['cosine_followees'] = df_final_test.apply(lambda row:
                                                                    cosine_for_followees(row['source_node'],
                                                                    row['target_node']),
                                                                    axis=1)
```

```
In [0]: def compute_features_stage1(df_final):
    #calculating no of followers followees for source and destination
    #calculating intersection of followers and followees for source and destination
    num_followers_s=[]
    num_followees_s=[]
    num_followers_d=[]
    num_followees_d=[]
    inter_followers=[]
    inter_followees=[]
    for i,row in df_final.iterrows():
        try:
            s1=set(train_graph.predecessors(row['source_node']))
            s2=set(train_graph.successors(row['source_node']))
        except:
            s1 = set()
            s2 = set()
        try:
            d1=set(train_graph.predecessors(row['destination_node']))
            d2=set(train_graph.successors(row['destination_node']))
        except:
            d1 = set()
            d2 = set()
        num_followers_s.append(len(s1))
        num_followees_s.append(len(s2))

        num_followers_d.append(len(d1))
        num_followees_d.append(len(d2))

        inter_followers.append(len(s1.intersection(d1)))
        inter_followees.append(len(s2.intersection(d2)))

    return num_followers_s, num_followers_d, num_followees_s, num_followees_d, inter_followers, inter_followees
```

```
In [0]: if not os.path.isfile('data/fea_sample/storage_sample_stage1.h5'):
    df_final_train['num_followers_s'], df_final_train['num_followers_d'], \
    df_final_train['num_followees_s'], df_final_train['num_followees_d'], \
    df_final_train['inter_followers'], df_final_train['inter_followees']= compute_features_stage1(df_final_train)

    df_final_test['num_followers_s'], df_final_test['num_followers_d'], \
    df_final_test['num_followees_s'], df_final_test['num_followees_d'], \
    df_final_test['inter_followers'], df_final_test['inter_followees']= compute_features_stage1(df_final_test)

    hdf = HDFStore('data/fea_sample/storage_sample_stage1.h5')
    hdf.put('train_df',df_final_train, format='table', data_columns=True)
    hdf.put('test_df',df_final_test, format='table', data_columns=True)
    hdf.close()
else:
    df_final_train = read_hdf('data/fea_sample/storage_sample_stage1.h5', 'train_df')
    df_final_test = read_hdf('data/fea_sample/storage_sample_stage1.h5', 'test_df')
```

## 5.3 Adding new set of features

we will create these each of these features for both train and test data points

1. adar index
2. is following back
3. belongs to same weakly connect components
4. shortest path between source and destination

```
In [0]: if not os.path.isfile('data/fea_sample/storage_sample_stage2.h5'):
#mapping adar index on train
df_final_train['adar_index'] = df_final_train.apply(lambda row: calc_adar_in(
#mapping adar index on test
df_final_test['adar_index'] = df_final_test.apply(lambda row: calc_adar_in(row)

#-----
#mapping followback or not on train
df_final_train['follows_back'] = df_final_train.apply(lambda row: follows_back(
#mapping followback or not on test
df_final_test['follows_back'] = df_final_test.apply(lambda row: follows_back(row)

#-----
#mapping same component of wcc or not on train
df_final_train['same_comp'] = df_final_train.apply(lambda row: belongs_to_same_
##mapping same component of wcc or not on test
df_final_test['same_comp'] = df_final_test.apply(lambda row: belongs_to_same_

#-----
#mapping shortest path on train
df_final_train['shortest_path'] = df_final_train.apply(lambda row: compute_sh
#mapping shortest path on test
df_final_test['shortest_path'] = df_final_test.apply(lambda row: compute_shor

hdf = HDFStore('data/fea_sample/storage_sample_stage2.h5')
hdf.put('train_df',df_final_train, format='table', data_columns=True)
hdf.put('test_df',df_final_test, format='table', data_columns=True)
hdf.close()
else:
df_final_train = read_hdf('data/fea_sample/storage_sample_stage2.h5', 'train_
df_final_test = read_hdf('data/fea_sample/storage_sample_stage2.h5', 'test_df
```

## 5.4 Adding new set of features

we will create these each of these features for both train and test data points

1. Weight Features
  - weight of incoming edges
  - weight of outgoing edges
  - weight of incoming edges + weight of outgoing edges
  - weight of incoming edges *weight of outgoing edges*
  - 2weight of incoming edges + weight of outgoing edges
  - weight of incoming edges + 2\*weight of outgoing edges
2. Page Ranking of source



```
In [0]: if not os.path.isfile('data/fea_sample/storage_sample_stage3.h5'):
        #mapping to pandas train
        df_final_train['weight_in'] = df_final_train.destination_node.apply(lambda x: W
        df_final_train['weight_out'] = df_final_train.source_node.apply(lambda x: Wei

        #mapping to pandas test
        df_final_test['weight_in'] = df_final_test.destination_node.apply(lambda x: W
        df_final_test['weight_out'] = df_final_test.source_node.apply(lambda x: Weigh

        #some features engineerings on the in and out weights
        df_final_train['weight_f1'] = df_final_train.weight_in + df_final_train.weigh
        df_final_train['weight_f2'] = df_final_train.weight_in * df_final_train.weigh
        df_final_train['weight_f3'] = (2*df_final_train.weight_in + 1*df_final_train.
        df_final_train['weight_f4'] = (1*df_final_train.weight_in + 2*df_final_train.

        #some features engineerings on the in and out weights
        df_final_test['weight_f1'] = df_final_test.weight_in + df_final_test.weight_o
        df_final_test['weight_f2'] = df_final_test.weight_in * df_final_test.weight_o
        df_final_test['weight_f3'] = (2*df_final_test.weight_in + 1*df_final_test.wei
        df_final_test['weight_f4'] = (1*df_final_test.weight_in + 2*df_final_test.wei
```

```

In [4]: if not os.path.isfile('data/fea_sample/storage_sample_stage3.h5'):

    #page rank for source and destination in Train and Test
    #if anything not there in train graph then adding mean page rank
    df_final_train['page_rank_s'] = df_final_train.source_node.apply(lambda x: pr.get(x, 0))
    df_final_train['page_rank_d'] = df_final_train.destination_node.apply(lambda x: pr.get(x, 0))

    df_final_test['page_rank_s'] = df_final_test.source_node.apply(lambda x: pr.get(x, 0))
    df_final_test['page_rank_d'] = df_final_test.destination_node.apply(lambda x: pr.get(x, 0))
    #=====

    #Katz centrality score for source and destination in Train and test
    #if anything not there in train graph then adding mean katz score
    df_final_train['katz_s'] = df_final_train.source_node.apply(lambda x: katz.get(x, 0))
    df_final_train['katz_d'] = df_final_train.destination_node.apply(lambda x: katz.get(x, 0))

    df_final_test['katz_s'] = df_final_test.source_node.apply(lambda x: katz.get(x, 0))
    df_final_test['katz_d'] = df_final_test.destination_node.apply(lambda x: katz.get(x, 0))
    #=====

    #Hits algorithm score for source and destination in Train and test
    #if anything not there in train graph then adding 0
    df_final_train['hubs_s'] = df_final_train.source_node.apply(lambda x: hits[0].get(x, 0))
    df_final_train['hubs_d'] = df_final_train.destination_node.apply(lambda x: hits[0].get(x, 0))

    df_final_test['hubs_s'] = df_final_test.source_node.apply(lambda x: hits[0].get(x, 0))
    df_final_test['hubs_d'] = df_final_test.destination_node.apply(lambda x: hits[0].get(x, 0))
    #=====

    #Hits algorithm score for source and destination in Train and Test
    #if anything not there in train graph then adding 0
    df_final_train['authorities_s'] = df_final_train.source_node.apply(lambda x: hits[1].get(x, 0))
    df_final_train['authorities_d'] = df_final_train.destination_node.apply(lambda x: hits[1].get(x, 0))

    df_final_test['authorities_s'] = df_final_test.source_node.apply(lambda x: hits[1].get(x, 0))
    df_final_test['authorities_d'] = df_final_test.destination_node.apply(lambda x: hits[1].get(x, 0))
    #=====

    hdf = HDFStore('data/fea_sample/storage_sample_stage3.h5')
    hdf.put('train_df', df_final_train, format='table', data_columns=True)
    hdf.put('test_df', df_final_test, format='table', data_columns=True)
    hdf.close()
else:
    df_final_train = read_hdf('data/fea_sample/storage_sample_stage3.h5', 'train_df')
    df_final_test = read_hdf('data/fea_sample/storage_sample_stage3.h5', 'test_df')

```

## 5.5 Adding new set of features

we will create these each of these features for both train and test data points

1. SVD features for both source and destination



```
In [12]: def svd(x, S):  
        try:  
            z = sadj_dict[x]  
            return S[z]  
        except:  
            return [0,0,0,0,0,0]
```

```
In [13]: #for svd features to get feature vector creating a dict node val and index in svd  
sadj_col = sorted(train_graph.nodes())  
sadj_dict = { val:idx for idx,val in enumerate(sadj_col)}
```

```
In [14]: Adj = nx.adjacency_matrix(train_graph,nodelist=sorted(train_graph.nodes())).asfptype
```

```
In [15]: U, s, V = svds(Adj, k = 6)  
print('Adjacency matrix Shape',Adj.shape)  
print('U Shape',U.shape)  
print('V Shape',V.shape)  
print('s Shape',s.shape)
```

```
Adjacency matrix Shape (1780722, 1780722)  
U Shape (1780722, 6)  
V Shape (6, 1780722)  
s Shape (6,)
```

```

In [16]: if not os.path.isfile('data/fea_sample/storage_sample_stage4.h5'):
#=====

df_final_train[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
df_final_train.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)

df_final_train[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
df_final_train.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)
#=====

df_final_train[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
df_final_train.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)

df_final_train[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
df_final_train.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
#=====

df_final_test[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
df_final_test.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)

df_final_test[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
df_final_test.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)

#=====

df_final_test[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
df_final_test.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)

df_final_test[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
df_final_test.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
#=====

hdf = HDFStore('data/fea_sample/storage_sample_stage4.h5')
hdf.put('train_df', df_final_train, format='table', data_columns=True)
hdf.put('test_df', df_final_test, format='table', data_columns=True)
hdf.close()

```

```

In [0]: # prepared and stored the data from machine learning models
# please check the FB_Models.ipynb

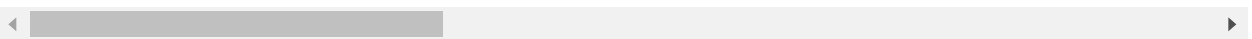
```

```
In [17]: #df_final_train
df_final_train.ix[:, 'weight_f1:'][:10]
```

Out[17]:

	weight_f1	weight_f2	weight_f3	weight_f4	page_rank_s	page_rank_d	katz_s	katz_d	h
0	0.627964	0.094491	1.005929	0.877964	2.045290e-06	3.459963e-07	0.000773	0.000756	1.94
1	0.229598	0.013030	0.332196	0.356598	2.353458e-07	6.427660e-07	0.000845	0.001317	3.90
2	0.339999	0.028653	0.525694	0.494302	6.211019e-07	5.179801e-07	0.000885	0.000855	7.73
3	0.696923	0.117851	0.985599	1.105172	2.998153e-07	1.704245e-06	0.000739	0.000773	5.44
4	1.301511	0.301511	2.301511	1.603023	4.349180e-07	2.089590e-07	0.000751	0.000735	3.88
5	0.617739	0.095346	0.933967	0.919250	5.942343e-07	1.143388e-06	0.000767	0.000766	9.05
6	1.447214	0.447214	1.894427	2.447214	2.848986e-07	1.128758e-06	0.000735	0.000750	1.60
7	0.853553	0.176777	1.353553	1.207107	6.694862e-07	5.254600e-07	0.000763	0.000743	5.60
8	0.583489	0.084515	0.850750	0.899717	1.466870e-06	1.373409e-06	0.000757	0.000781	5.73
9	0.930904	0.204124	1.508254	1.284457	6.630224e-07	2.618341e-07	0.000758	0.000739	2.49

10 rows × 42 columns



## Preferential\_attachment

```
In [18]: def followee_preferential_attachment(user1,user2):
    try:
        user_1 = len(set(train_graph.successors(user1)))
        user_2 = len(set(train_graph.successors(user2)))
        return(user_1*user_2)
    except:
        return(0)

def follower_preferential_attachment(user1,user2):
    try:
        user_1 = len(set(train_graph.predecessors(user1)))
        user_2 = len(set(train_graph.predecessors(user2)))
        return(user_1*user_2)
    except:
        return(0)
```

```

In [19]: startTime = datetime.datetime.now()
print("Current Time = ",startTime)

if not os.path.isfile('data/fea_sample/storage_sample_stage5.h5'):
    #=====

    df_final_train[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
                    'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
                    'svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
                    'svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
                    'followee_preferential_attachment', 'follower_preferential_attachment']] = df_final_train[['source_node', 'destination_node']].apply(lambda x: svd(x, U)).apply(pd.Series)

    df_final_train[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
                    'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
                    'svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
                    'svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
                    'followee_preferential_attachment', 'follower_preferential_attachment']] = df_final_train[['source_node', 'destination_node']].apply(lambda x: svd(x, V.T)).apply(pd.Series)

    df_final_test[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
                   'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
                   'svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
                   'svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
                   'followee_preferential_attachment', 'follower_preferential_attachment']] = df_final_test[['source_node', 'destination_node']].apply(lambda x: svd(x, U)).apply(pd.Series)

    df_final_test[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
                   'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
                   'svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
                   'svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
                   'followee_preferential_attachment', 'follower_preferential_attachment']] = df_final_test[['source_node', 'destination_node']].apply(lambda x: svd(x, V.T)).apply(pd.Series)

    hdf = HDFStore('data/fea_sample/storage_sample_stage5.h5')
    hdf.put('train_df',df_final_train, format='table', data_columns=True)
    hdf.put('test_df',df_final_test, format='table', data_columns=True)
    hdf.close()
else:
    df_final_train = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'train_df')
    df_final_test = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'test_df')

print("Time taken for creation of dataframe is {}".format(datetime.datetime.now() - startTime))

```

Current Time = 2019-05-24 21:18:39.178276

Time taken for creation of dataframe is 0:00:02.927990

```
In [30]: # for Train data

x1 = list(df_final_train['svd_u_s_1'])
x2 = list(df_final_train['svd_u_s_2'])
x3 = list(df_final_train['svd_u_s_3'])
x4 = list(df_final_train['svd_u_s_4'])
x5 = list(df_final_train['svd_u_s_5'])
x6 = list(df_final_train['svd_u_s_6'])

x7 = list(df_final_train['svd_u_d_1'])
x8 = list(df_final_train['svd_u_d_2'])
x9 = list(df_final_train['svd_u_d_3'])
x10 = list(df_final_train['svd_u_d_4'])
x11 = list(df_final_train['svd_u_d_5'])
x12 = list(df_final_train['svd_u_d_6'])

y1 = list(df_final_train['svd_v_s_1'])
y2 = list(df_final_train['svd_v_s_2'])
y3 = list(df_final_train['svd_v_s_3'])
y4 = list(df_final_train['svd_v_s_4'])
y5 = list(df_final_train['svd_v_s_5'])
y6 = list(df_final_train['svd_v_s_6'])

y7 = list(df_final_train['svd_v_d_1'])
y8 = list(df_final_train['svd_v_d_2'])
y9 = list(df_final_train['svd_v_d_3'])
y10 = list(df_final_train['svd_v_d_4'])
y11 = list(df_final_train['svd_v_d_5'])
y12 = list(df_final_train['svd_v_d_6'])

print(np.shape(x1))
print(np.shape(x2))
print(np.shape(x3))
print(np.shape(x4))
print(np.shape(x5))
print(np.shape(x6))
print(np.shape(x7))
print(np.shape(x8))
print(np.shape(x9))
print(np.shape(x10))
print(np.shape(x11))
print(np.shape(x12))

print(np.shape(y1))
print(np.shape(y2))
print(np.shape(y3))
print(np.shape(y4))
print(np.shape(y5))
print(np.shape(y6))
print(np.shape(y7))
print(np.shape(y8))
print(np.shape(y9))
print(np.shape(y10))
print(np.shape(y11))
print(np.shape(y12))
```

```

train_u_source = []
train_u_destination = []
train_v_source = []
train_v_destination = []
train_u_s_dot = []
train_u_d_dot = []

for loop1 in range(0,len(x1)):
    train_u_source.append(x1[loop1])
    train_u_source.append(x2[loop1])
    train_u_source.append(x3[loop1])
    train_u_source.append(x4[loop1])
    train_u_source.append(x5[loop1])
    train_u_source.append(x6[loop1])

    train_u_destination.append(x7[loop1])
    train_u_destination.append(x8[loop1])
    train_u_destination.append(x9[loop1])
    train_u_destination.append(x10[loop1])
    train_u_destination.append(x11[loop1])
    train_u_destination.append(x12[loop1])

    dot_product = np.dot(train_u_source[loop1],train_u_destination[loop1])
    train_u_s_dot.append(dot_product)

for loop2 in range(0,len(y1)):
    train_v_source.append(y1[loop2])
    train_v_source.append(y2[loop2])
    train_v_source.append(y3[loop2])
    train_v_source.append(y4[loop2])
    train_v_source.append(y5[loop2])
    train_v_source.append(y6[loop2])

    train_v_destination.append(y7[loop2])
    train_v_destination.append(y8[loop2])
    train_v_destination.append(y9[loop2])
    train_v_destination.append(y10[loop2])
    train_v_destination.append(y11[loop2])
    train_v_destination.append(y12[loop2])

    dot_product = np.dot(train_v_source[loop2],train_v_destination[loop2])
    train_u_d_dot.append(dot_product)

print("*****")
print(np.shape(train_u_s_dot))
print(np.shape(train_u_d_dot))

```

```

(100002,)
(100002,)
(100002,)
(100002,)
(100002,)
(100002,)
(100002,)
(100002,)
(100002,)

```

```
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
(100002,)  
*****  
(100002,)  
(100002,)
```

In [31]: *# for Test data*

```
x1 = list(df_final_test['svd_u_s_1'])
x2 = list(df_final_test['svd_u_s_2'])
x3 = list(df_final_test['svd_u_s_3'])
x4 = list(df_final_test['svd_u_s_4'])
x5 = list(df_final_test['svd_u_s_5'])
x6 = list(df_final_test['svd_u_s_6'])

x7 = list(df_final_test['svd_u_d_1'])
x8 = list(df_final_test['svd_u_d_2'])
x9 = list(df_final_test['svd_u_d_3'])
x10 = list(df_final_test['svd_u_d_4'])
x11 = list(df_final_test['svd_u_d_5'])
x12 = list(df_final_test['svd_u_d_6'])

y1 = list(df_final_test['svd_v_s_1'])
y2 = list(df_final_test['svd_v_s_2'])
y3 = list(df_final_test['svd_v_s_3'])
y4 = list(df_final_test['svd_v_s_4'])
y5 = list(df_final_test['svd_v_s_5'])
y6 = list(df_final_test['svd_v_s_6'])

y7 = list(df_final_test['svd_v_d_1'])
y8 = list(df_final_test['svd_v_d_2'])
y9 = list(df_final_test['svd_v_d_3'])
y10 = list(df_final_test['svd_v_d_4'])
y11 = list(df_final_test['svd_v_d_5'])
y12 = list(df_final_test['svd_v_d_6'])

print(np.shape(x1))
print(np.shape(x2))
print(np.shape(x3))
print(np.shape(x4))
print(np.shape(x5))
print(np.shape(x6))
print(np.shape(x7))
print(np.shape(x8))
print(np.shape(x9))
print(np.shape(x10))
print(np.shape(x11))
print(np.shape(x12))

print(np.shape(y1))
print(np.shape(y2))
print(np.shape(y3))
print(np.shape(y4))
print(np.shape(y5))
print(np.shape(y6))
print(np.shape(y7))
print(np.shape(y8))
print(np.shape(y9))
print(np.shape(y10))
print(np.shape(y11))
print(np.shape(y12))
```



```

test_u_source = []
test_u_destination = []
test_v_source = []
test_v_destination = []
test_v_s_dot = []
test_v_d_dot = []

for loop3 in range(0,len(x1)):
    test_u_source.append(x1[loop3])
    test_u_source.append(x2[loop3])
    test_u_source.append(x3[loop3])
    test_u_source.append(x4[loop3])
    test_u_source.append(x5[loop3])
    test_u_source.append(x6[loop3])

    test_u_destination.append(x7[loop3])
    test_u_destination.append(x8[loop3])
    test_u_destination.append(x9[loop3])
    test_u_destination.append(x10[loop3])
    test_u_destination.append(x11[loop3])
    test_u_destination.append(x12[loop3])

    dot_product = np.dot(test_u_source[loop3],test_u_destination[loop3])
    test_v_s_dot.append(dot_product)

for loop4 in range(0,len(y1)):
    test_v_source.append(y1[loop4])
    test_v_source.append(y2[loop4])
    test_v_source.append(y3[loop4])
    test_v_source.append(y4[loop4])
    test_v_source.append(y5[loop4])
    test_v_source.append(y6[loop4])

    test_v_destination.append(y7[loop4])
    test_v_destination.append(y8[loop4])
    test_v_destination.append(y9[loop4])
    test_v_destination.append(y10[loop4])
    test_v_destination.append(y11[loop4])
    test_v_destination.append(y12[loop4])

    dot_product = np.dot(test_v_source[loop4],test_v_destination[loop4])
    test_v_d_dot.append(dot_product)

print("*****")
print(np.shape(test_v_s_dot))
print(np.shape(test_v_d_dot))

```

```

(50002,)
(50002,)
(50002,)
(50002,)
(50002,)
(50002,)
(50002,)
(50002,)
(50002,)

```



In [33]: `df_final_test.ix[:, 'adar_index':][:10]`

Out[33]:

	adar_index	follows_back	same_comp	shortest_path	weight_in	weight_out	weight_f1	weight_
0	0.000000	1	1	2	0.258199	0.377964	0.636163	0.09756
1	0.000000	1	1	7	0.235702	0.707107	0.942809	0.16666
2	0.000000	0	1	5	0.301511	0.242536	0.544047	0.07311
3	0.000000	0	1	3	0.162221	0.301511	0.463733	0.04897
4	6.136433	0	1	2	0.188982	0.250000	0.438982	0.04724
5	0.000000	0	0	-1	0.588969	0.301511	0.890481	0.17758
6	0.000000	1	1	-1	1.000000	0.353553	1.353553	0.35355
7	3.095903	1	1	2	0.250000	0.288675	0.538675	0.07216
8	0.000000	0	0	-1	0.588969	0.301511	0.890481	0.17758
9	0.000000	1	1	-1	0.377964	1.000000	1.377964	0.37796

10 rows × 48 columns

In [24]: `df_final_test.ix[:, 'weight_f4':][:10]`

3	0.765244	7.000791e-07	1.657820e-06	0.000778	0.000884	9.876114e-14	1.039593e-13	1.5116
4	0.688982	7.103008e-07	1.181606e-06	0.000779	0.000840	1.557332e-15	1.096037e-14	5.1808
5	1.191992	3.803655e-07	5.615700e-07	0.000739	0.000748	2.494904e-16	0.000000e+00	3.4428
6	1.707107	1.251767e-06	1.655650e-07	0.000762	0.000731	2.188066e-14	5.745713e-16	8.0050
7	0.827350	1.319717e-06	1.790599e-06	0.000777	0.000789	1.672452e-17	9.345255e-14	6.7223
8	1.191992	3.329817e-07	5.615700e-07	0.000747	0.000748	6.356182e-17	0.000000e+00	6.8316
9	0.377964	2.848986e-07	1.000000e-06	0.000735	0.000751	0.000000e+00	0.000000e+00	1.1010

In [25]: `df_final_test.ix[:, 'svd_u_s_1'][:10]`

3	-4.491249e-11	9.917580e-13	7.891237e-06	9.458626e-12	2.715849e-11	1.822079e-12	-1.070734e-10	2.8880
4	-9.778569e-13	5.724678e-13	4.951878e-06	1.396835e-12	2.153250e-11	2.872774e-14	-2.275782e-11	1.0076
5	-8.977767e-15	4.541045e-15	9.100025e-14	2.419543e-13	4.586857e-14	4.601913e-15	0.000000e+00	0.0000
6	-2.143381e-11	1.069559e-09	3.714264e-11	2.544021e-11	2.935397e-11	4.036789e-13	-5.328096e-14	4.1412
7	-2.047625e-14	3.044747e-14	8.187949e-13	4.402437e-15	1.083419e-14	3.136596e-16	-1.012196e-11	2.1269
8	-9.936400e-14	1.079878e-12	5.993042e-12	1.324362e-13	2.225664e-08	1.171559e-15	0.000000e+00	0.0000
9	-4.079425e-20	-1.681818e-20	-1.769057e-19	9.391413e-20	1.097935e-19	-1.897890e-19	-1.277149e-17	4.5740

10 rows × 30 columns

In [26]: `df_final_test.ix[:, 'svd_u_d_4'][:10]`

Out[26]:

	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
0	1.166038e-13	2.253356e-11	3.222301e-15	-2.148853e-13	1.883259e-13	5.904813e-11	2.701538e-12	4.300000e+00
1	1.907402e-12	3.797448e-11	4.993662e-14	-4.054500e-13	2.895773e-13	2.545371e-10	2.248600e-14	3.600000e+00
2	1.693277e-13	2.712767e-13	7.596286e-15	-4.148335e-13	4.618813e-12	1.122147e-05	1.778931e-12	2.700000e+00
3	7.229387e-12	3.563661e-11	1.917977e-12	-8.942576e-12	5.535692e-12	5.223671e-06	7.917210e-13	4.000000e+00
4	1.394103e-10	2.961998e-11	2.022055e-13	-3.804983e-12	1.593137e-13	1.035013e-06	1.361585e-13	1.100000e+00
5	0.000000e+00	0.000000e+00	0.000000e+00	-2.065320e-20	1.431037e-20	1.688691e-19	4.502556e-19	5.300000e+00
6	3.072229e-13	2.571312e-13	1.059874e-14	-5.259869e-12	4.156213e-09	1.005038e-09	3.214437e-11	3.200000e+00

```
In [34]: df_final_test.ix[:, 'svd_v_d_1:'][:10]
```

```
Out[34]:
```

svd_v_d_4	svd_v_d_5	svd_v_d_6	followee_preferential_attachment	follower_preferential_atta
2.486660e-09	2.771146e-09	1.727695e-12		54
6.665782e-12	1.495980e-10	9.836744e-14		19
1.630555e-13	3.954762e-13	3.875775e-14		144
4.384839e-12	1.239415e-11	6.483494e-13		340
3.637116e-12	3.948477e-12	2.415871e-13		405
0.000000e+00	0.000000e+00	0.000000e+00		0
0.000000e+00	0.000000e+00	0.000000e+00		7
5.675210e-13	4.877886e-13	3.437278e-14		275
0.000000e+00	0.000000e+00	0.000000e+00		0
6.131609e-14	1.297121e-13	6.960836e-16		0

```
In [20]: y_train = df_final_train.indicator_link
y_test = df_final_test.indicator_link
```

```
In [35]: df_final_train.columns
```

```
Out[35]: Index(['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',
               'followee_preferential_attachment', 'follower_preferential_attachment',
               'train_u_s_dot', 'train_u_d_dot', 's_dot', 'd_dot'],
              dtype='object')
```

```
In [ ]: df_final_train.drop(['source_node', 'destination_node', 'indicator_link', 'train_u_s_dot',
                             'train_u_d_dot', 's_dot', 'd_dot'], axis=1)
df_final_test.drop(['source_node', 'destination_node', 'indicator_link', 'test_v_s_dot',
                    'test_v_d_dot', 's_dot', 'd_dot'], axis=1)
```

```
In [39]: df_final_train.drop(['train_u_s_dot', 'train_u_d_dot'], axis=1, inplace=True)
df_final_test.drop(['test_v_s_dot', 'test_v_d_dot'], axis=1, inplace=True)
```

```
In [40]: df_final_train.columns
```

```
Out[40]: Index(['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',
'followee_preferential_attachment', 'follower_preferential_attachment',
's_dot', 'd_dot'],
dtype='object')
```

```
In [41]: df_final_test.columns
```

```
Out[41]: Index(['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',
'followee_preferential_attachment', 'follower_preferential_attachment',
's_dot', 'd_dot'],
dtype='object')
```

## Random Forest

```

In [162]: startTime = datetime.datetime.now()
print("Current Time = ",startTime)

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, random_state=25, verbose=0)
    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')
plt.title('Estimators vs score at depth of 5')

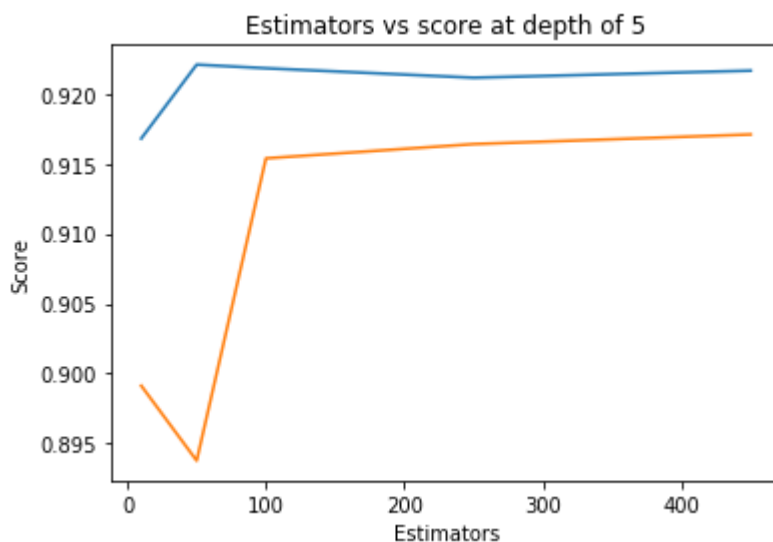
print("Time taken for creation of dataframe is {}".format(datetime.datetime.now()-startTime))

```

```

Current Time = 2019-05-24 17:15:40.888059
Estimators = 10 Train Score 0.9168054647804179 test Score 0.8991258141926637
Estimators = 50 Train Score 0.9220895584588777 test Score 0.8937688248990665
Estimators = 100 Train Score 0.9218287697728152 test Score 0.9153860707603692
Estimators = 250 Train Score 0.9211508895233912 test Score 0.9164059383071205
Estimators = 450 Train Score 0.9216563839799227 test Score 0.9170952910797137
Time taken for creation of dataframe is 0:06:27.128323

```



```

In [163]: startTime = datetime.datetime.now()
print("Current Time = ",startTime)

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, random_state=25, verbose=0)
    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')
plt.show()

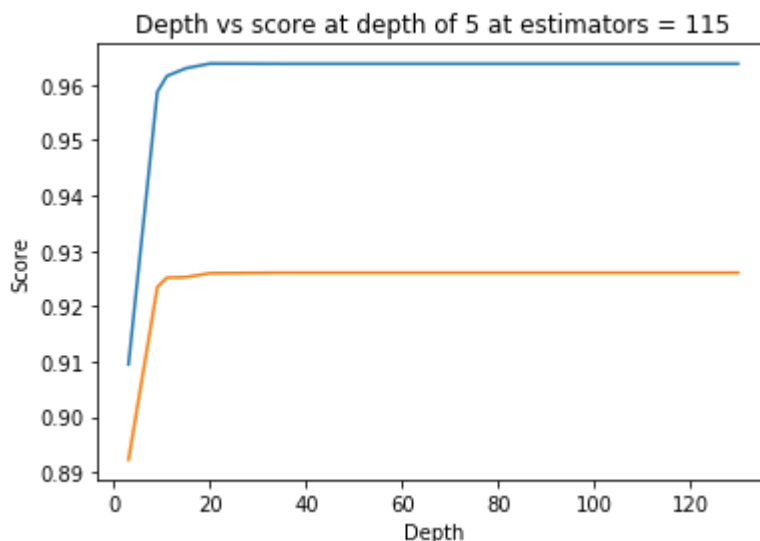
print("Time taken for creation of dataframe is {}".format(datetime.datetime.now()-startTime))

```

```

Current Time = 2019-05-24 17:22:55.575829
depth = 3 Train Score 0.9094969294708766 test Score 0.8923076923076922
depth = 9 Train Score 0.9586727982787979 test Score 0.9234631450149783
depth = 11 Train Score 0.9615533269594579 test Score 0.92516748830742
depth = 15 Train Score 0.9629599545122248 test Score 0.9252643552260185
depth = 20 Train Score 0.963790564217155 test Score 0.9259657330720111
depth = 35 Train Score 0.963760461525961 test Score 0.9260593800703796
depth = 50 Train Score 0.963760461525961 test Score 0.9260593800703796
depth = 70 Train Score 0.963760461525961 test Score 0.9260593800703796
depth = 130 Train Score 0.963760461525961 test Score 0.9260593800703796

```



Time taken for creation of dataframe is 0:11:28.420432



```

In [164]: 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

startTime = datetime.datetime.now()
print("Current Time = ",startTime)

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)

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)
print('mean test scores',rf_random.cv_results_['mean_test_score'])
print('mean train scores',rf_random.cv_results_['mean_train_score'])

print("*****")
print(rf_random.best_estimator_)

clf = 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)

clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)

from sklearn.metrics import f1_score
print('\nTrain f1 score',f1_score(y_train,y_train_pred))
print('Test f1 score',f1_score(y_test,y_test_pred))

print("Time taken for creation of dataframe is {}".format(datetime.datetime.now())

```

Current Time = 2019-05-24 17:34:24.043261

mean test scores [0.96230736 0.96209145 0.96068497 0.96182116 0.96318534]

mean train scores [0.9633391 0.96287272 0.96126794 0.96251656 0.96423855]

\*\*\*\*\*

```

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=None,
                        oob_score=False, random_state=25, verbose=0, warm_start=False)

```

Train f1 score 0.9643266955735856

Test f1 score 0.9263264402706634

Time taken for creation of dataframe is 0:59:16.060665

```
In [46]: 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")

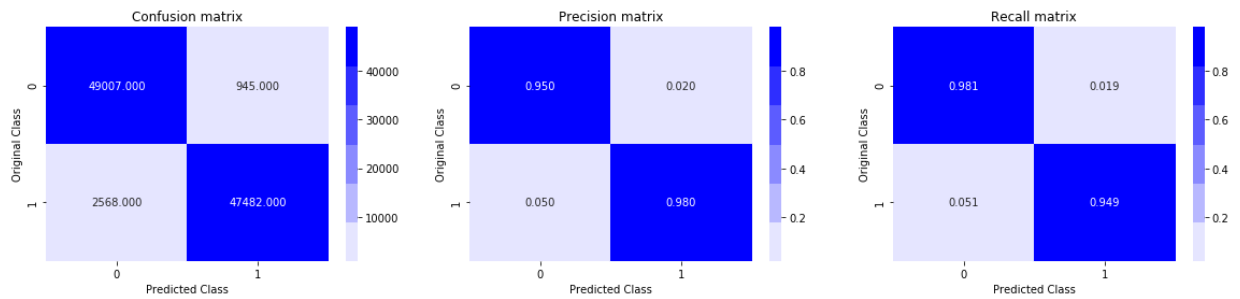
    plt.subplot(1, 3, 2)
    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Precision matrix")

    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")

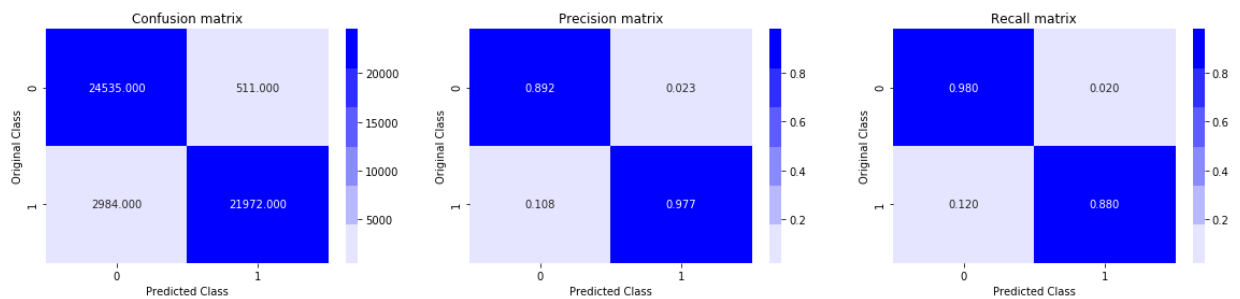
    plt.show()
```

```
In [166]: print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

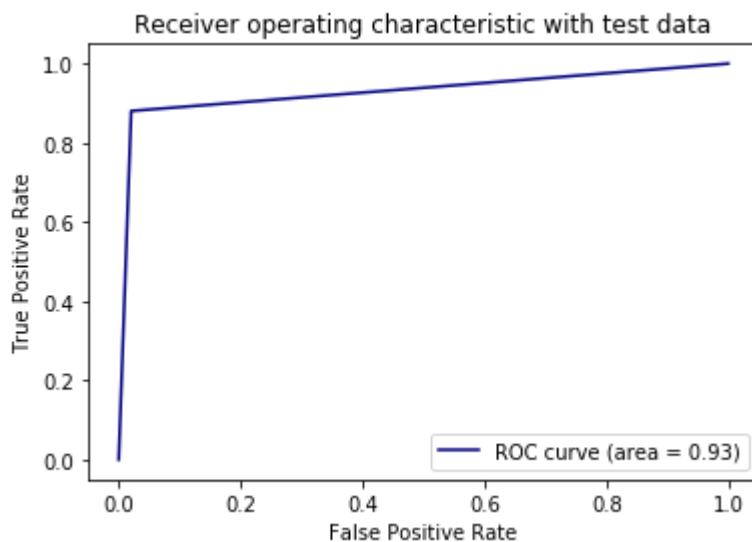
Train confusion\_matrix



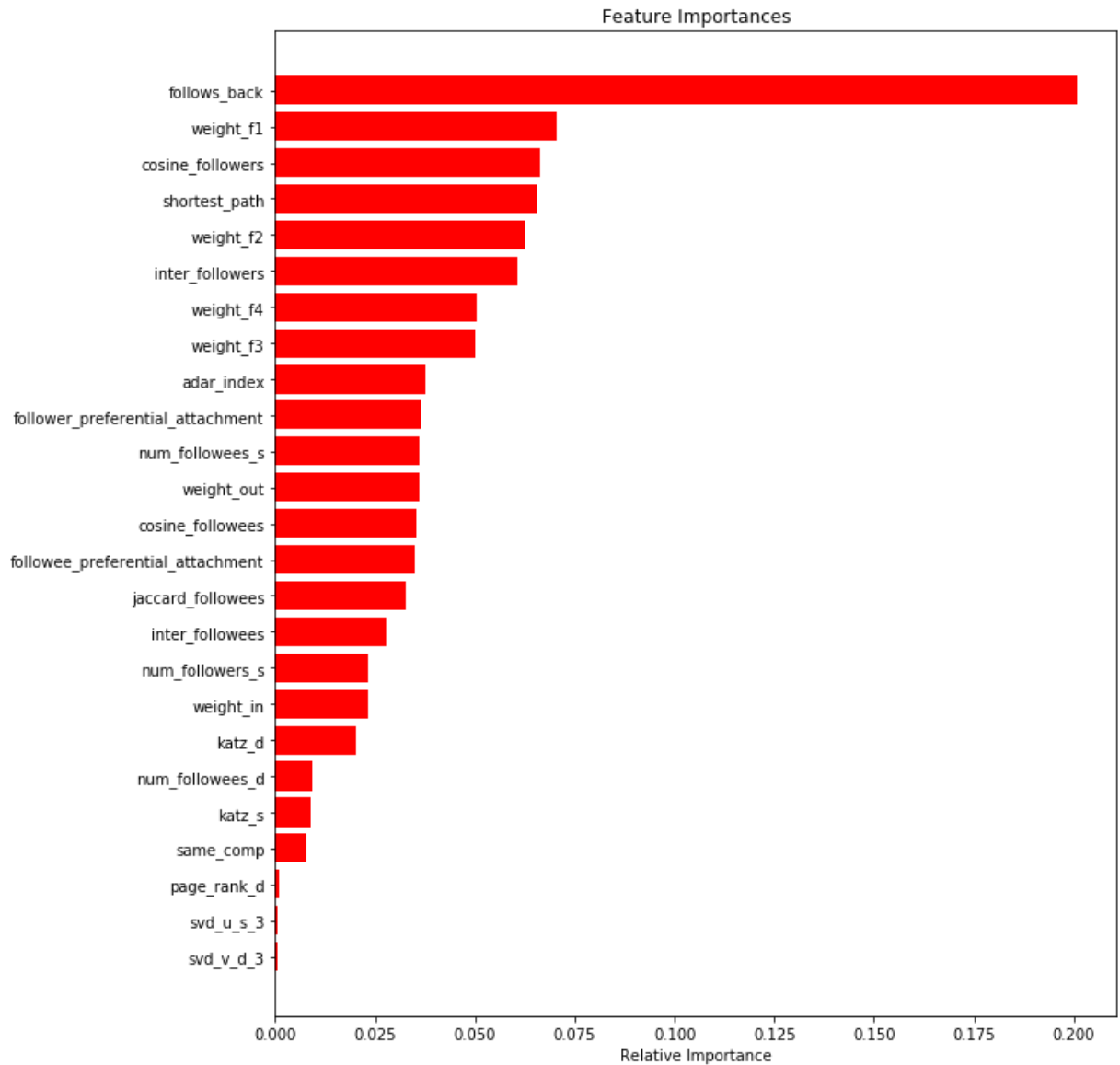
Test confusion\_matrix



```
In [167]: 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()
plt.show()
```



```
In [168]: 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')
plt.show()
```



## Hyperparameter tuning XGBoost

```

In [43]: startTime = datetime.datetime.now()
print("Current Time = ",startTime)

import xgboost as xgb
from sklearn.model_selection import RandomizedSearchCV
from sklearn.metrics import f1_score,make_scorer

min_child_weight = [2,4,6]
max_depth = [2,4,6]
n_estimators =[100,200,300]
learning_rate = [0.1,0.2,0.3]

scorer = make_scorer(f1_score)
tuned_parameters = {
    'min_child_weight':min_child_weight,
    'max_depth':max_depth,
    'n_estimators': n_estimators,
    'learning_rate':learning_rate}

clf = xgb.XGBClassifier()
model_gbt = RandomizedSearchCV(clf,tuned_parameters,scoring =scorer,cv=3,pre_dispatch=5)
model_gbt.fit(df_final_train,y_train)
print(model_gbt.best_estimator_)

best_min_child_weight_xgb = model_gbt.best_estimator_.min_child_weight
best_max_depth_xgb = model_gbt.best_params_["max_depth"]
best_n_estimators_xgb = model_gbt.best_estimator_.n_estimators
best_learning_rate_xgb = model_gbt.best_estimator_.learning_rate

print("\nbest_min_child_weight_xgb = ", best_min_child_weight_xgb)
print("best_max_depth_xgb = ",best_max_depth_xgb)
print("best_n_estimators_xgb = ", best_n_estimators_xgb)
print("best_learning_rate_xgb = ",best_learning_rate_xgb)

print("Time taken for creation of dataframe is {}".format(datetime.datetime.now()-startTime))

Current Time = 2019-05-24 21:25:09.249187
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bytree=1, gamma=0, learning_rate=0.3, max_delta_step=0,
              max_depth=4, min_child_weight=4, missing=None, n_estimators=300,
              n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=True, subsample=1)

best_min_child_weight_xgb = 4
best_max_depth_xgb = 4
best_n_estimators_xgb = 300
best_learning_rate_xgb = 0.3
Time taken for creation of dataframe is 1:55:39.179742

```

```
In [44]: startTime = datetime.datetime.now()
print("Current Time = ",startTime)

xgb_best = xgb.XGBClassifier(objective='binary:logistic',learning_rate = best_lea
                             min_child_weight = best_min_child_weight_xgb,n_estim
                             max_depth = best_max_depth_xgb)
xgb_best.fit(df_final_train,y_train)
pred_train = xgb_best.predict(df_final_train)
pred_test = xgb_best.predict(df_final_test)

train_score = f1_score(y_train,pred_train)
test_score = f1_score(y_test,pred_test)
print('\nTrain Score: ',train_score)
print('Test Score: ',test_score)

print("Time taken for creation of dataframe is {}".format(datetime.datetime.now()

```

Current Time = 2019-05-24 23:20:48.456948

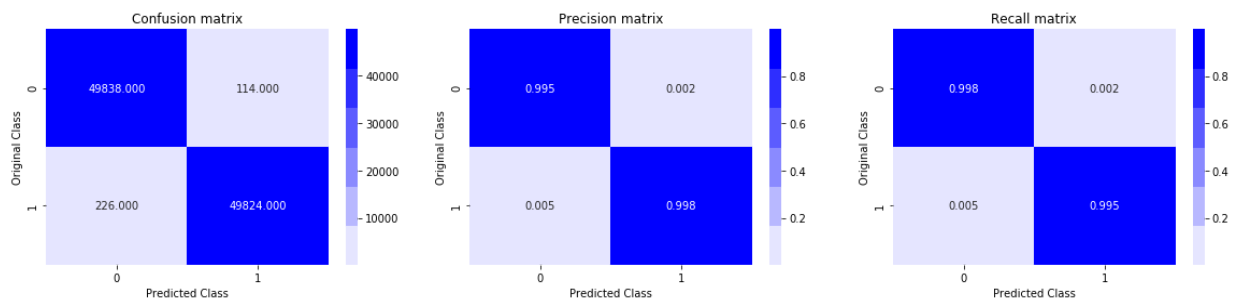
Train Score: 0.9965995919510341

Test Score: 0.875133547008547

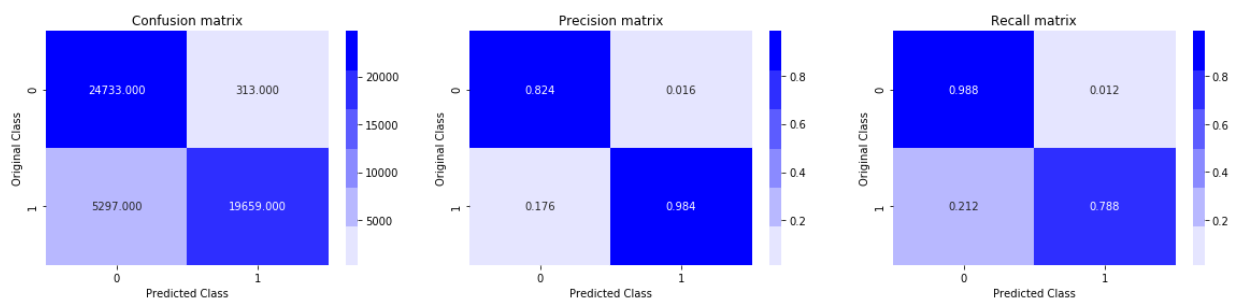
Time taken for creation of dataframe is 0:06:15.302619

```
In [47]: print('Train confusion_matrix')
plot_confusion_matrix(y_train,pred_train)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,pred_test)
```

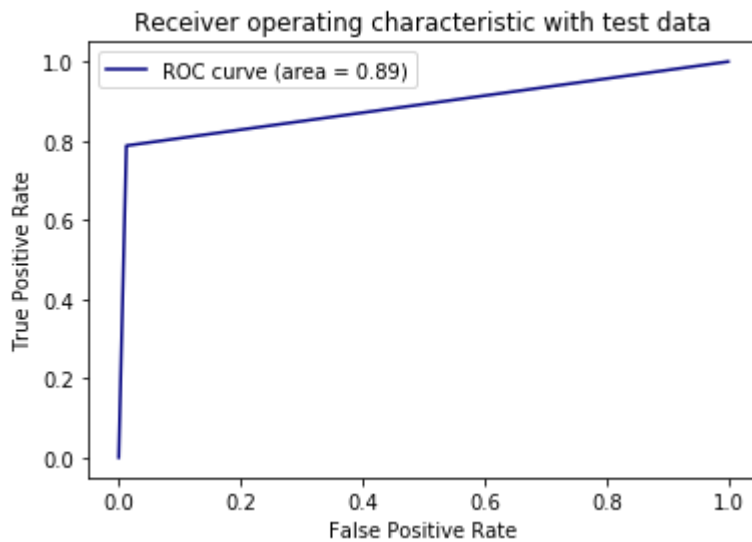
Train confusion\_matrix



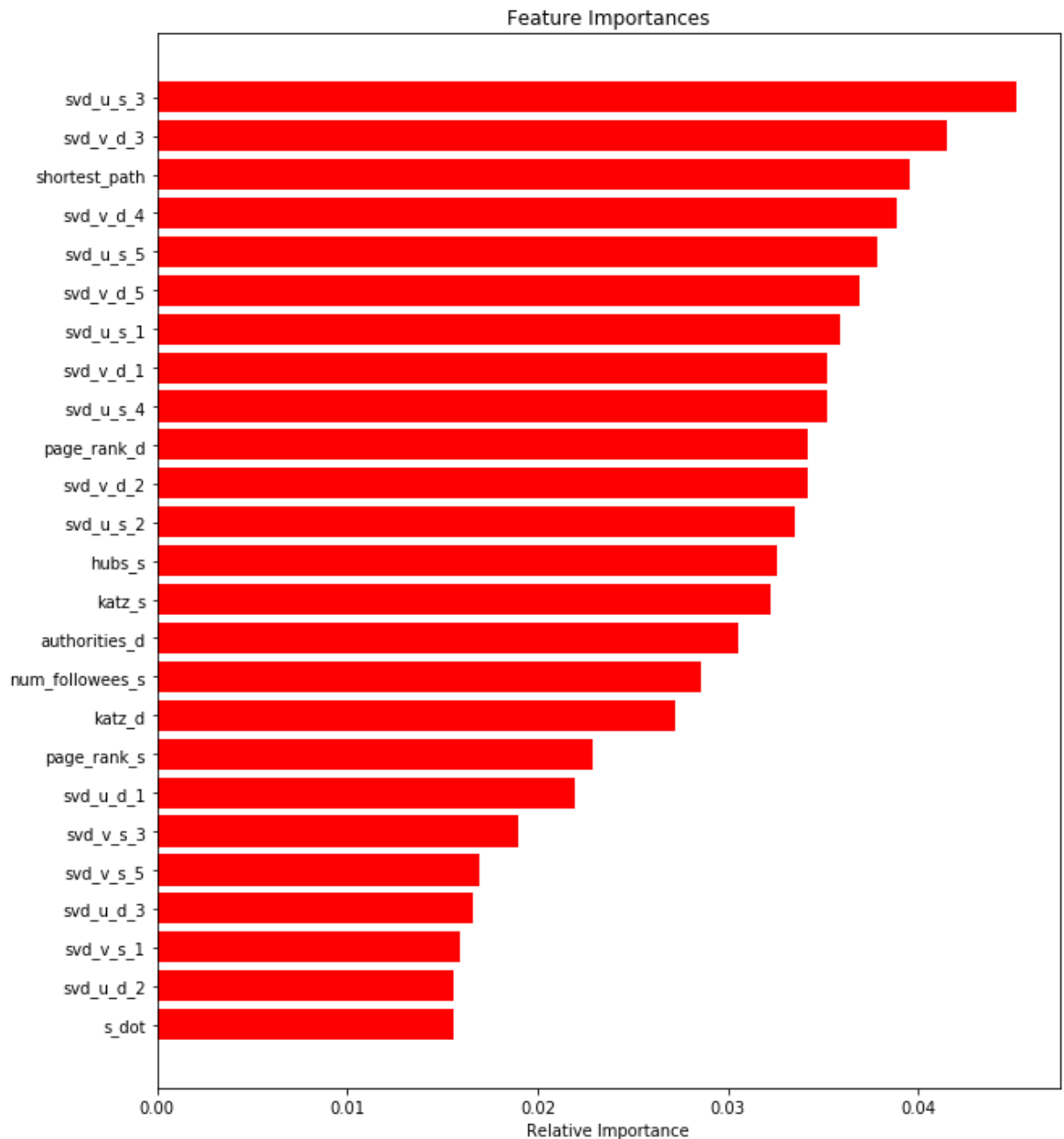
Test confusion\_matrix



```
In [48]: from sklearn.metrics import roc_curve, auc
fpr, tpr, ths = roc_curve(y_test, 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()
plt.show()
```



```
In [49]: features = df_final_train.columns
importances = xgb_best.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')
plt.show()
```





```
In [50]: from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Model \ Parameters", "Train f1_score", "Test f1_score"]
x.add_row(["RandomForest: ", 0.9643266955735856, 0.9263264402706634])
x.add_row(["XGBClassifier: ", 0.9965995919510341, 0.875133547008547])
print(x)
```

Model \ Parameters	Train f1_score	Test f1_score
RandomForest:	0.9643266955735856	0.9263264402706634
XGBClassifier:	0.9965995919510341	0.875133547008547

In [ ]: