

Social network Graph Link Prediction - Facebook Challenge

Problem statement:

Given a directed social graph, have to predict missing links to recommend users (Link Prediction in graph)

Data Overview

Taken data from facebook's recruiting challenge on kaggle <https://www.kaggle.com/c/FacebookRecruiting> (<https://www.kaggle.com/c/FacebookRecruiting>)
data contains two columns source and destination eac edge in graph

- Data columns (total 2 columns):
- source_node int64
- destination_node int64

Mapping the problem into supervised learning problem:

- Generated training samples of good and bad links from given directed graph and for each link got some features like no of followers, is he followed back, page rank, katz score, adar index, some svd fetures of adj matrix, some weight features etc. and trained ml model based on these features to predict link.
- Some reference papers and videos :
 - <https://www.cs.cornell.edu/home/kleinber/link-pred.pdf> (<https://www.cs.cornell.edu/home/kleinber/link-pred.pdf>)
 - <https://www3.nd.edu/~dial/publications/lichtenwalter2010new.pdf> (<https://www3.nd.edu/~dial/publications/lichtenwalter2010new.pdf>)
 - <https://www.youtube.com/watch?v=2M77Hgy17cg> (<https://www.youtube.com/watch?v=2M77Hgy17cg>)

Business objectives and constraints:

- No low-latency requirement.
- Probability of prediction is useful to recommend ighest probability links

Performance metric for supervised learning:

- Both precision and recall is important so F1 score is good choice
- Confusion matrix

```

In [1]: ▶ 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 tqdm
21 import os
22 # to install xgboost: pip3 install xgboost
23 import xgboost as xgb
24
25 import warnings
26 import networkx as nx
27 import pdb
28 from sklearn.ensemble import RandomForestClassifier
29 from sklearn.metrics import f1_score

```

```

In [2]: ▶ 1 ## Reading the graph
2 if not os.path.isfile('train_woheader.csv'):
3     traincsv = pd.read_csv('train.csv')
4     print(traincsv[traincsv.isna().any(1)])
5     print(traincsv.info())
6     print('Number of duplicates entry',sum(traincsv.duplicated()))
7     traincsv.to_csv('train_woheader.csv',header=False,index=False)
8     print('Saved the graph into file')
9     g=nx.read_edgelist('train_woheader.csv',delimiter=',',
10                      create_using=nx.DiGraph(),nodetype=int)
11 else:
12     g=nx.read_edgelist('train_woheader.csv',delimiter=',',
13                      create_using=nx.DiGraph(),nodetype=int)
14     print(nx.info(g))
15

```

Name:

Type: DiGraph

Number of nodes: 1862220

Number of edges: 9437519

Average in degree: 5.0679

Average out degree: 5.0679

```

In [3]: ► 1  ## Lets create a subgraph and visualize it
          2  if not os.path.isfile('train_woheader_sample.csv'):
          3      pd.read_csv('train.csv',nrows=50).to_csv('train_woheader_sample.csv',header=False,index=False)
          4
          5  subgraph = nx.read_edgelist('train_woheader_sample.csv',delimiter=',',
          6                          create_using=nx.DiGraph(),nodetype=int)
          7  pos = nx.spring_layout(subgraph)
          8  nx.draw(subgraph,pos,node_color='#A0CBE2',edge_color='#BB0000',width=2,
          9                          edge_cmap=plt.cm.Blues,with_labels=True)
         10  plt.savefig('graph_sample.pdf')
         11  print(nx.info(subgraph))

```

Name:
 Type: DiGraph
 Number of nodes: 66
 Number of edges: 50
 Average in degree: 0.7576
 Average out degree: 0.7576



1. Exploratory Data Analysis

```

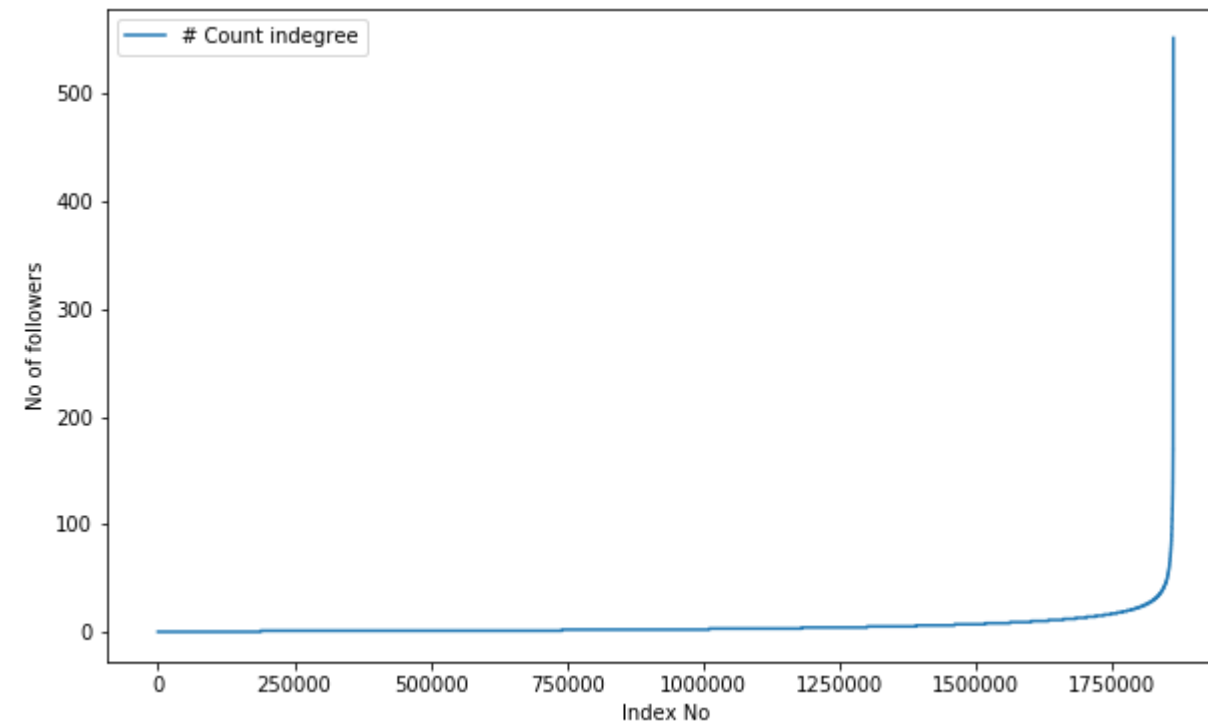
In [4]: ► 1  ## Number of unique persons
          2  print('The Number of unique persons in the node ',len(g.nodes()))

```

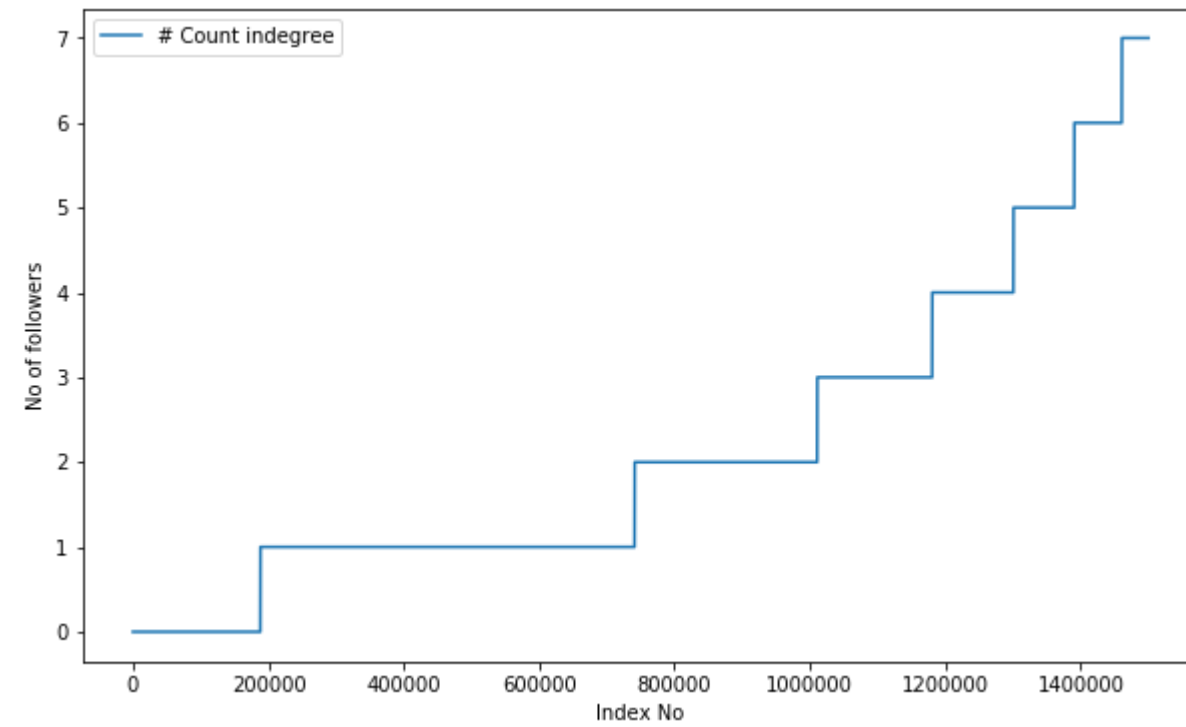
The Number of unique persons in the node 1862220

1.1 Number of Followers for each person

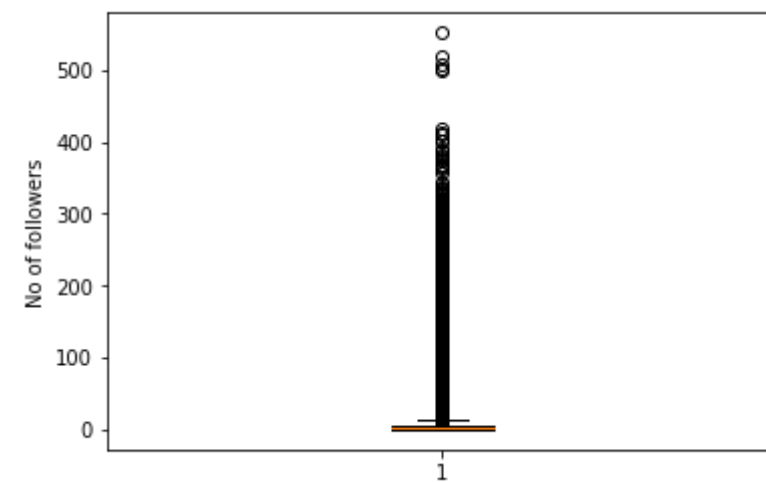
```
In [5]: 1 ## Number of followers of each person
2 indegree_dist = list(dict(g.in_degree()).values())
3 indegree_dist.sort()
4 plt.figure(figsize=(10,6))
5 plt.plot(indegree_dist,label='# Count indegree')
6 plt.xlabel('Index No')
7 plt.ylabel('No of followers')
8 plt.legend()
9 plt.show()
```



```
In [6]: 1  ## Lets zoom the portion where there is lot of followers
2  indegree_dist = list(dict(g.in_degree()).values())
3  indegree_dist.sort()
4  plt.figure(figsize=(10,6))
5  plt.plot(indegree_dist[0:1500000],label='# Count indegree')
6  plt.xlabel('Index No')
7  plt.ylabel('No of followers')
8  plt.legend()
9  plt.show()
```



```
In [7]: 1  plt.boxplot(indegree_dist)
2  plt.ylabel('No of followers')
3  plt.show()
```



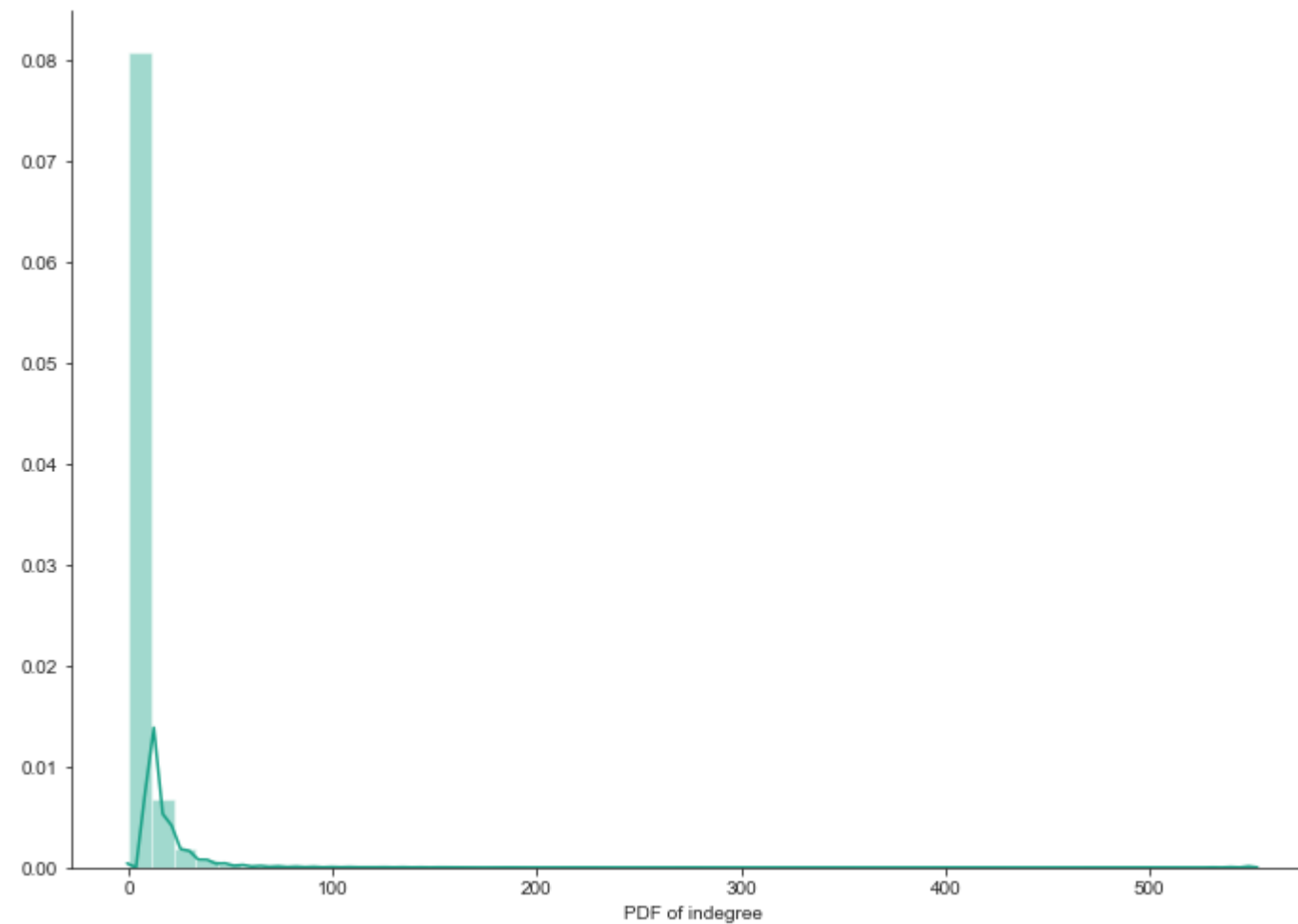
```
In [8]: 1  ### 90th to 100th percentile
        2  for i in range(0,11):
        3      print(90+i,'percentile value is',np.percentile(indegree_dist,90+i))
        4  print('*'*50)
        5  for i in range(10,110,10):
        6      print(99+(i/100),'percentile value is',np.percentile(indegree_dist,99+(i/100)))
```

```
90 percentile value is 12.0
91 percentile value is 13.0
92 percentile value is 14.0
93 percentile value is 15.0
94 percentile value is 17.0
95 percentile value is 19.0
96 percentile value is 21.0
97 percentile value is 24.0
98 percentile value is 29.0
99 percentile value is 40.0
100 percentile value is 552.0
*****
99.1 percentile value is 42.0
99.2 percentile value is 44.0
99.3 percentile value is 47.0
99.4 percentile value is 50.0
99.5 percentile value is 55.0
99.6 percentile value is 61.0
99.7 percentile value is 70.0
99.8 percentile value is 84.0
99.9 percentile value is 112.0
100.0 percentile value is 552.0
```

```
In [9]: 1  print('99% of people having 40 or less followers')
```

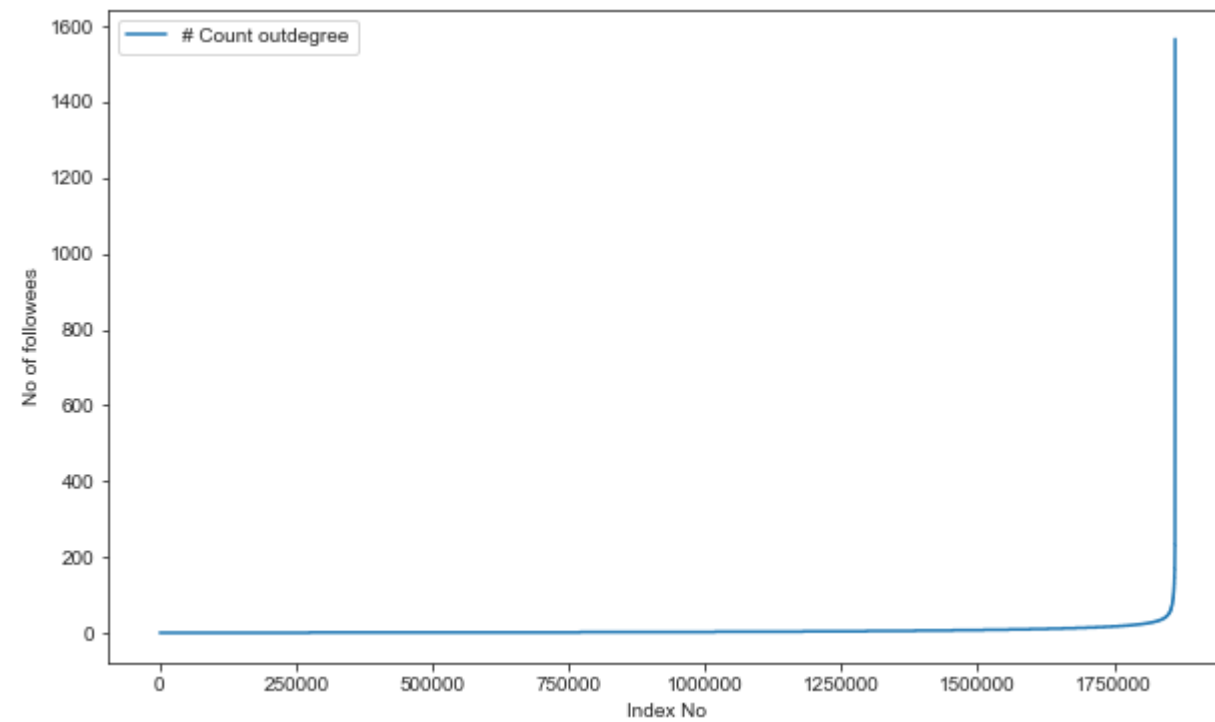
```
99% of people having 40 or less followers
```

```
In [10]: ▶ 1 %matplotlib inline
2 sns.set_style('ticks')
3 fig,ax = plt.subplots()
4 fig.set_size_inches(11.7,8.27)
5 sns.distplot(indegree_dist,color='#16A085')
6 plt.xlabel('PDF of indegree')
7 sns.despine()
```

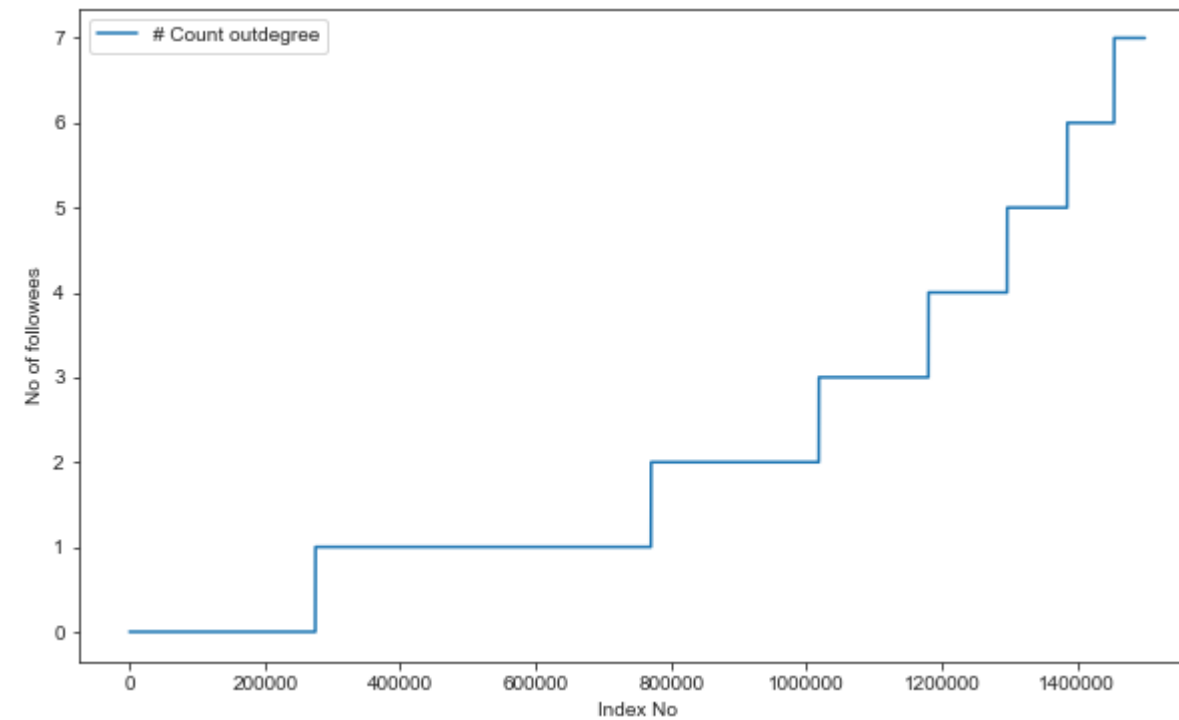


1.2 Number of Followees for each person

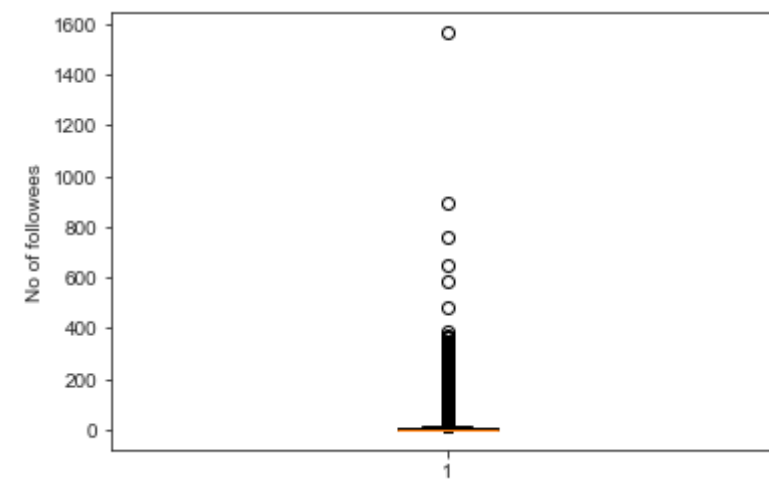
```
In [11]: 1  ## Number of followers of each person
2  outdegree_dist = list(dict(g.out_degree()).values())
3  outdegree_dist.sort()
4  plt.figure(figsize=(10,6))
5  plt.plot(outdegree_dist,label='# Count outdegree')
6  plt.xlabel('Index No')
7  plt.ylabel('No of followees')
8  plt.legend()
9  plt.show()
```




```
In [12]: 1  ## Lets zoom the portion where there is lot of followers
2  outdegree_dist = list(dict(g.out_degree()).values())
3  outdegree_dist.sort()
4  plt.figure(figsize=(10,6))
5  plt.plot(outdegree_dist[0:1500000],label='# Count outdegree')
6  plt.xlabel('Index No')
7  plt.ylabel('No of followees')
8  plt.legend()
9  plt.show()
```



```
In [13]: 1  plt.boxplot(outdegree_dist)
2  plt.ylabel('No of followees')
3  plt.show()
```



```
In [14]: 1  ### 90th to 100th percentile
2  for i in range(0,11):
3      print(90+i, 'percentile value is', np.percentile(outdegree_dist, 90+i))
4  print('*'*50)
5  for i in range(10,110,10):
6      print(99+(i/100), 'percentile value is', np.percentile(outdegree_dist, 99+(i/100)))
```

```
90 percentile value is 12.0
91 percentile value is 13.0
92 percentile value is 14.0
93 percentile value is 15.0
94 percentile value is 17.0
95 percentile value is 19.0
96 percentile value is 21.0
97 percentile value is 24.0
98 percentile value is 29.0
99 percentile value is 40.0
100 percentile value is 1566.0
*****
99.1 percentile value is 42.0
99.2 percentile value is 45.0
99.3 percentile value is 48.0
99.4 percentile value is 52.0
99.5 percentile value is 56.0
99.6 percentile value is 63.0
99.7 percentile value is 73.0
99.8 percentile value is 90.0
99.9 percentile value is 123.0
100.0 percentile value is 1566.0
```

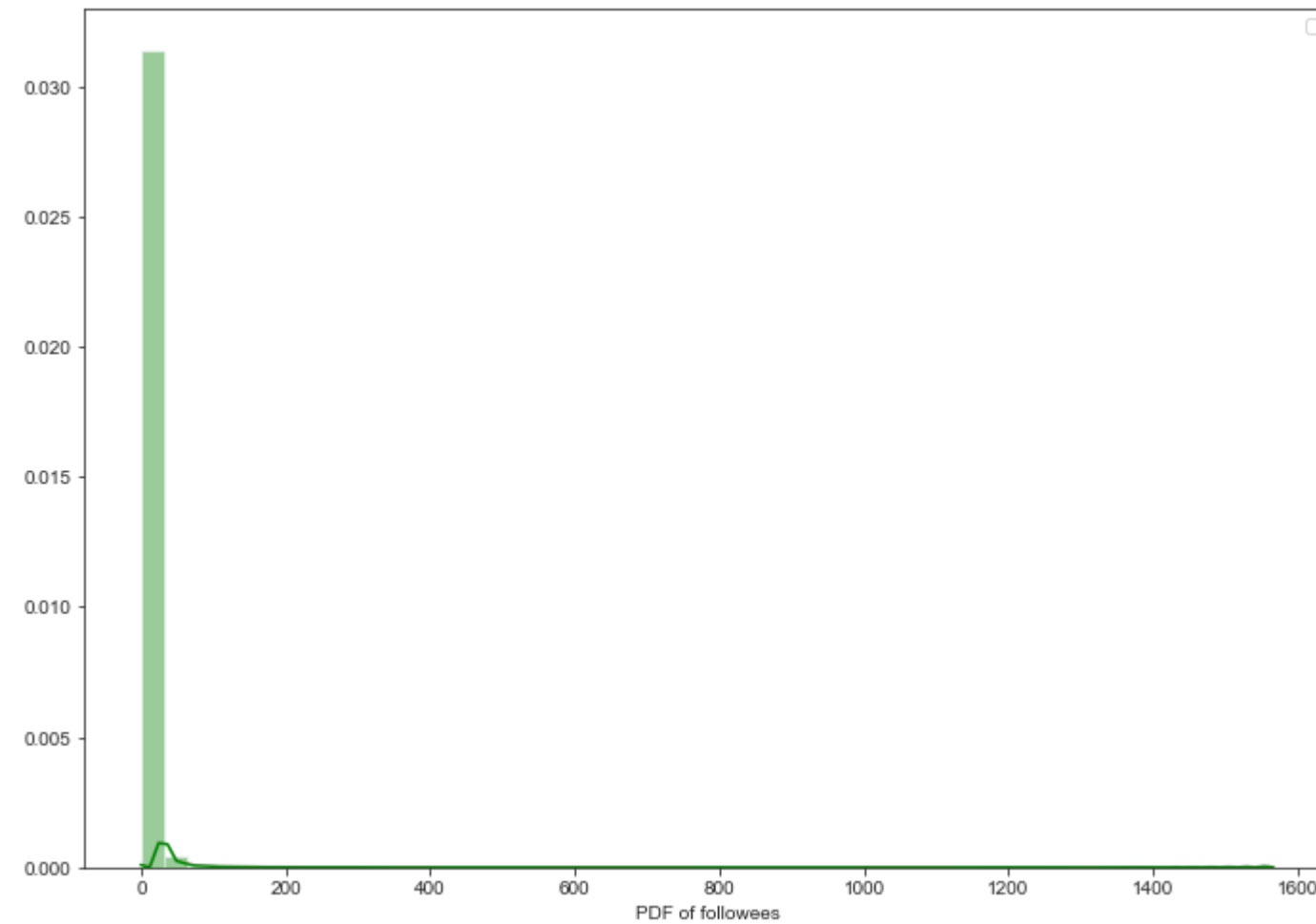
```
In [15]: 1  print('99% of people having 40 or less followees')
```

```
99% of people having 40 or less followees
```

```
In [16]: 1 sns.set_style('ticks')
2 fig,ax = plt.subplots()
3 fig.set_size_inches(11.7,8.27)
4 sns.distplot(outdegree_dist,color='green')
5 plt.xlabel('PDF of followees')
6 plt.legend()
```

No handles with labels found to put in legend.

Out[16]: <matplotlib.legend.Legend at 0x236342340f0>



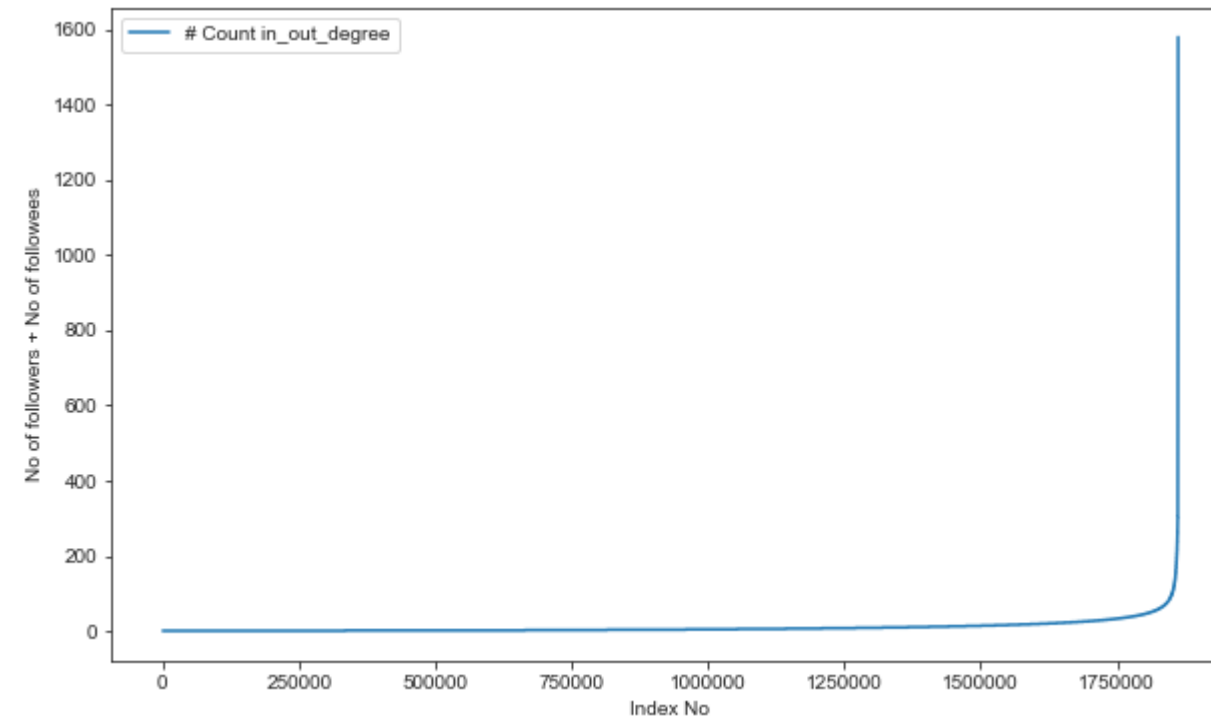
```
In [17]: ► 1 ### Number of people who are not following anyone are
2 print('Number of people who are not following anyone are ',sum(np.array(outdegree_dist)==0),
3       'and % is',sum(np.array(outdegree_dist)==0)*100 /len(outdegree_dist))
4 ### Number of people who are not having any followers
5 print('Number of people who are not having any followers ',sum(np.array(indegree_dist)==0),
6       'and % is',sum(np.array(indegree_dist)==0)*100 /len(indegree_dist))
7 ### Number of people who are having zero followers + zero followees
8 count=0
9 for i in g.nodes():
10     if len(list(g.predecessors(i)))==0 and len(list(g.successors(i)))==0:
11         count+=1
12 print('Number of people who are having zero followers + zero followees are',count)
```

Number of people who are not following anyone are 274512 and % is 14.7411154429
Number of people who are not having any followers 188043 and % is 10.0977865129
Number of people who are having zero followers + zero followees are 0

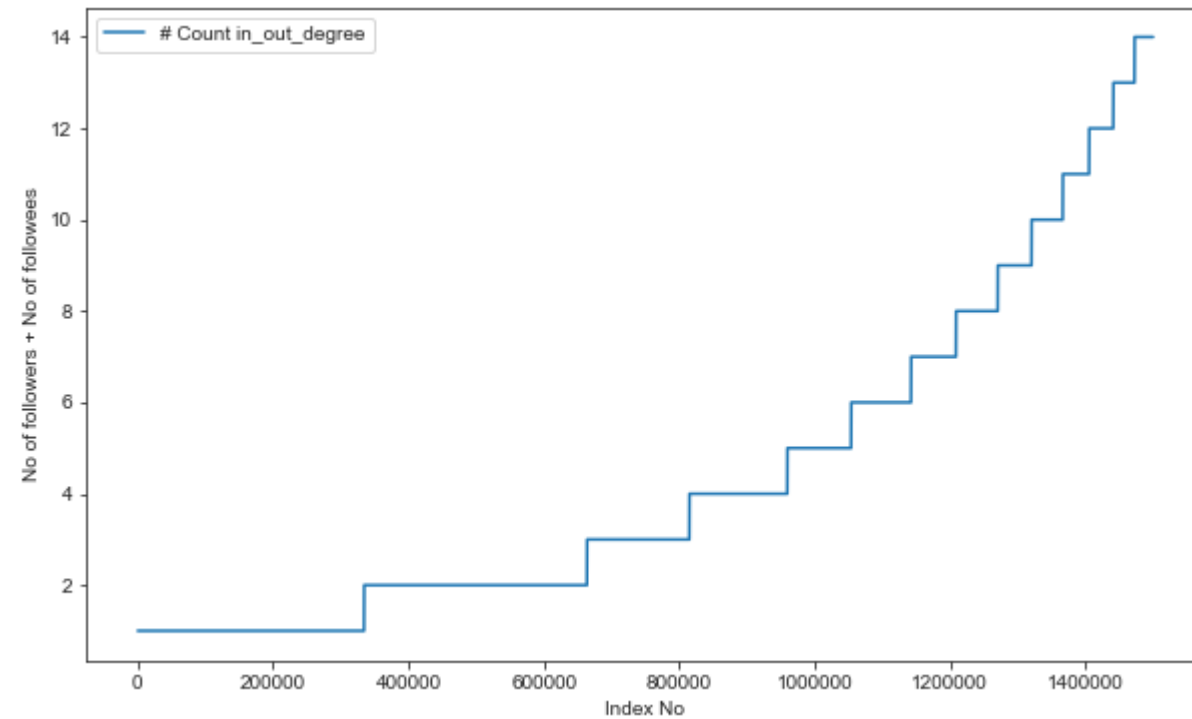
1.3 Both followers + following

```
In [18]: ► 1 from collections import Counter
2 dict_in = dict(g.in_degree())
3 dict_out = dict(g.out_degree())
4 d = Counter(dict_in) + Counter(dict_out)
5 in_out_degree_dict = np.array(list(d.values()))
```

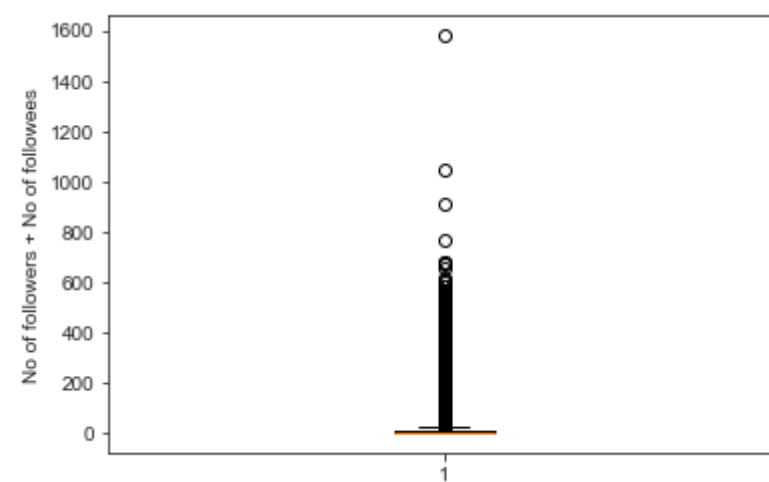
```
In [19]: 1  ## Number of followers of each person
2  in_out_degree_dict.sort()
3  plt.figure(figsize=(10,6))
4  plt.plot(in_out_degree_dict,label='# Count in_out_degree')
5  plt.xlabel('Index No')
6  plt.ylabel('No of followers + No of followees')
7  plt.legend()
8  plt.show()
```



```
In [20]: 1  ## Lets zoom the portion where there is lot of followers
2  in_out_degree_dict.sort()
3  plt.figure(figsize=(10,6))
4  plt.plot(in_out_degree_dict[0:1500000],label='# Count in_out_degree')
5  plt.xlabel('Index No')
6  plt.ylabel('No of followers + No of followees')
7  plt.legend()
8  plt.show()
```



```
In [21]: 1  plt.boxplot(in_out_degree_dict)
2  plt.ylabel('No of followers + No of followees')
3  plt.show()
```



```
In [22]: 1 ### 90th to 100th percentile
2 for i in range(0,11):
3     print(90+i, 'percentile value is', np.percentile(in_out_degree_dict, 90+i))
4 print('*'*50)
5 for i in range(10,110,10):
6     print(99+(i/100), 'percentile value is', np.percentile(in_out_degree_dict, 99+(i/100)))
```

```
90 percentile value is 24.0
91 percentile value is 26.0
92 percentile value is 28.0
93 percentile value is 31.0
94 percentile value is 33.0
95 percentile value is 37.0
96 percentile value is 41.0
97 percentile value is 48.0
98 percentile value is 58.0
99 percentile value is 79.0
100 percentile value is 1579.0
*****
99.1 percentile value is 83.0
99.2 percentile value is 87.0
99.3 percentile value is 93.0
99.4 percentile value is 99.0
99.5 percentile value is 108.0
99.6 percentile value is 120.0
99.7 percentile value is 138.0
99.8 percentile value is 168.0
99.9 percentile value is 221.0
100.0 percentile value is 1579.0
```

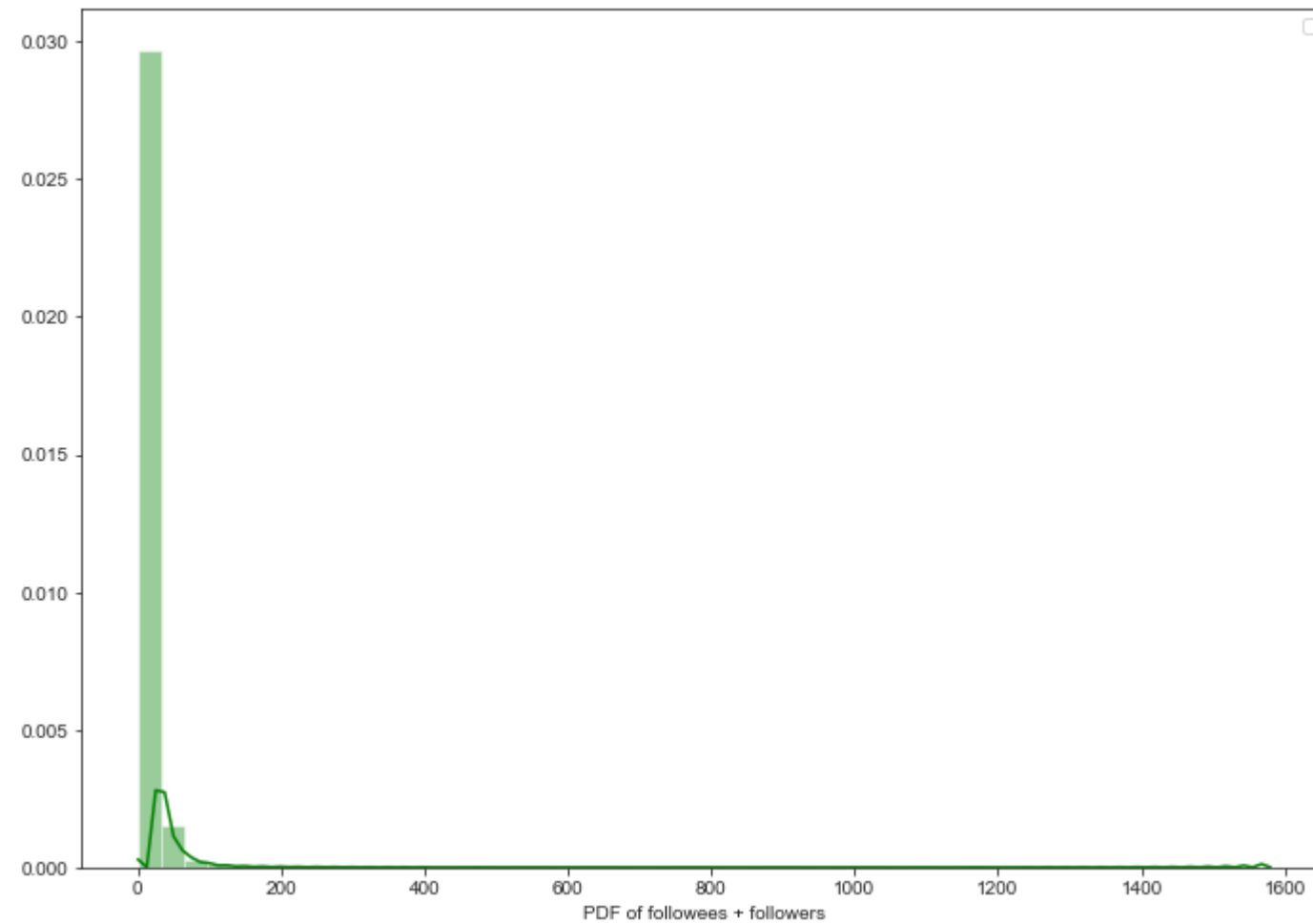
```
In [23]: 1 print('99% of people having 40 or less followees and followers combined')
```

```
99% of people having 40 or less followees and followers combined
```

```
In [24]: 1 sns.set_style('ticks')
2 fig,ax = plt.subplots()
3 fig.set_size_inches(11.7,8.27)
4 sns.distplot(in_out_degree_dict,color='green')
5 plt.xlabel('PDF of followees + followers')
6 plt.legend()
```

No handles with labels found to put in legend.

Out[24]: <matplotlib.legend.Legend at 0x236348e8f60>




```
In [25]: 1 print('Min of no of followers + following is',in_out_degree_dict.min())
2 print(np.sum(in_out_degree_dict==in_out_degree_dict.min()),' persons having minimum no of followers + following')
3 print('*'*50)
4 print('Max of no of followers + following is',in_out_degree_dict.max())
5 print(np.sum(in_out_degree_dict==in_out_degree_dict.max()),' persons having maximum no of followers + following')
6 print('*'*50)
7 print('No of persons having followers + following less than 10 are',np.sum(in_out_degree_dict<10))
```

```
Min of no of followers + following is 1
334291 persons having minimum no of followers + following
*****
Max of no of followers + following is 1579
1 persons having maximum no of followers + following
*****
No of persons having followers + following less than 10 are 1320326
```

```
In [26]: 1 print('No of Weakly connected components',len(list(nx.weakly_connected_components(g))))
2 count=0
3 for i in list(nx.weakly_connected_components(g)):
4     if len(i)==2:
5         count+=1
6 print('weakly connected components with 2 nodes',count)
```

```
No of Weakly connected components 45558
weakly connected components with 2 nodes 32195
```

2. Posing a problem as classification problem¶

2.1 Generating some edges which are not present in graph for supervised learning

- Generated Bad links from graph which are not in graph and whose shortest path is greater than 2

```
In [27]: ▶ 1 import random
2 import csv
3 if not os.path.isfile('missing_edges_final.p'):
4     r = csv.reader(open('train_woheader.csv','r'))
5     edges = dict()
6     for edge in r :
7         edges[(edge[0],edge[1])] = 1
8     missing_edges = set([])
9     while (len(missing_edges)<9437519):
10         a = random.randint(1,1862220)
11         b = random.randint(1,1862220)
12         tmp = edges.get((a,b),-1)
13         if tmp == -1 and a!=b:
14             try:
15                 if nx.shortest_path_length(g,source=a,target=b) > 2:
16                     missing_edges.add(a,b)
17             else:
18                 continue
19         except:
20             missing_edges.add((a,b))
21     pickle.dump(missing_edges,open('missing_edges_final.p','wb'))
22 else:
23     missing_edges = pickle.load(open('missing_edges_final.p','rb'))
```

```
In [28]: ▶ 1 missing_edges = pickle.load(open('missing_edges_final.p','rb'))
2 len(missing_edges)
```

Out[28]: 9437519

2.2 Training and Test data split:

- Removed edges from Graph and used as test data and after removing used that graph for creating features for Train and test data

```

In [29]: ▶ 1 from sklearn.model_selection import train_test_split
2 if (not os.path.isfile('train_pos_after_eda.csv')) and (not os.path.isfile('test_pos_after_eda.csv')):
3     #reading total data df
4     df_pos = pd.read_csv('train.csv')
5     df_neg = pd.DataFrame(list(missing_edges), columns=['source_node', 'destination_node'])
6     print('Number of nodes in graph with edges',df_pos.shape[0])
7     print('Number of nodes in graph without edges',df_neg.shape[0])
8     ## train test split is done separately for pos and neg because we need pos edges for
9     ## feature generation and creating graph
10    x_train_pos,x_test_pos,y_train_pos,y_test_pos = train_test_split(df_pos,np.ones(len(df_pos)))
11    x_train_neg,x_test_neg,y_train_neg,y_test_neg = train_test_split(df_neg,np.zeros(len(df_neg)))
12    print('='*60)
13    print("Number of nodes in the train data graph with edges", x_train_pos.shape[0],"=",y_train_pos.shape[0])
14    print("Number of nodes in the train data graph without edges", x_train_neg.shape[0],"=", y_train_neg.shape[0])
15    print('='*60)
16    print("Number of nodes in the test data graph with edges", x_test_pos.shape[0],"=",y_test_pos.shape[0])
17    print("Number of nodes in the test data graph without edges", x_test_neg.shape[0],"=",y_test_neg.shape[0])
18    #removing header and saving
19    x_train_pos.to_csv('train_pos_after_eda.csv',header=False, index=False)
20    x_test_pos.to_csv('test_pos_after_eda.csv',header=False, index=False)
21    x_train_neg.to_csv('train_neg_after_eda.csv',header=False, index=False)
22    x_test_neg.to_csv('test_neg_after_eda.csv',header=False, index=False)
23    pd.DataFrame(data=y_train_pos,columns=['indicator_link']).to_csv('y_tr_pos_after_eda.csv',header=False, index=False)
24    pd.DataFrame(data=y_test_pos,columns=['indicator_link']).to_csv('y_te_pos_after_eda.csv',header=False, index=False)
25    pd.DataFrame(data=y_train_neg,columns=['indicator_link']).to_csv('y_tr_neg_after_eda.csv',header=False, index=False)
26    pd.DataFrame(data=y_test_neg,columns=['indicator_link']).to_csv('y_te_neg_after_eda.csv',header=False, index=False)
27 else:
28     #Graph from Training data only
29     print('deleting .....')
30     del missing_edges

```

deleting

```

In [30]: 1 if (os.path.isfile('train_pos_after_eda.csv')) and (os.path.isfile('test_pos_after_eda.csv')):
2     train_graph=nx.read_edgelist('train_pos_after_eda.csv',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
3     test_graph=nx.read_edgelist('test_pos_after_eda.csv',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
4     print(nx.info(train_graph))
5     print(nx.info(test_graph))
6
7     # finding the unique nodes in the both train and test graphs
8     train_nodes_pos = set(train_graph.nodes())
9     test_nodes_pos = set(test_graph.nodes())
10
11     trY_teY = len(train_nodes_pos.intersection(test_nodes_pos))
12     trY_teN = len(train_nodes_pos - test_nodes_pos)
13     teY_trN = len(test_nodes_pos - train_nodes_pos)
14
15     print('no of people common in train and test -- ',trY_teY)
16     print('no of people present in train but not present in test -- ',trY_teN)
17
18     print('no of people present in test but not present in train -- ',teY_trN)
19     print(' % of people not there in Train but exist in Test in total Test data are {} %'.format(teY_trN/len(test_nodes_pos)*100))

```

Name:

Type: DiGraph

Number of nodes: 1755007

Number of edges: 7078139

Average in degree: 4.0331

Average out degree: 4.0331

Name:

Type: DiGraph

Number of nodes: 1252565

Number of edges: 2359380

Average in degree: 1.8836

Average out degree: 1.8836

no of people common in train and test -- 1145352

no of people present in train but not present in test -- 609655

no of people present in test but not present in train -- 107213

% of people not there in Train but exist in Test in total Test data are 8.559475955339643 %

```

In [34]: 1 #final train and test data sets
2 if (os.path.isfile('train_pos_after_eda.csv')) and \
3     (os.path.isfile('test_pos_after_eda.csv')) and \
4     (os.path.isfile('train_neg_after_eda.csv')) and \
5     (os.path.isfile('test_neg_after_eda.csv')) and \
6     (os.path.isfile('y_tr_pos_after_eda.csv')) and \
7     (os.path.isfile('y_te_pos_after_eda.csv')) and \
8     (os.path.isfile('y_tr_neg_after_eda.csv')) and \
9     (os.path.isfile('y_te_neg_after_eda.csv')) :
10
11 x_train_pos = pd.read_csv('train_pos_after_eda.csv', names=['source_node', 'destination_node'])
12 x_test_pos = pd.read_csv('test_pos_after_eda.csv', names=['source_node', 'destination_node'])
13 x_train_neg = pd.read_csv('train_neg_after_eda.csv', names=['source_node', 'destination_node'])
14 x_test_neg = pd.read_csv('test_neg_after_eda.csv', names=['source_node', 'destination_node'])
15 y_tr_pos = pd.read_csv('y_tr_pos_after_eda.csv', names=['source_node', 'destination_node'])
16 y_tr_neg = pd.read_csv('y_tr_neg_after_eda.csv', names=['source_node', 'destination_node'])
17 y_te_pos = pd.read_csv('y_te_pos_after_eda.csv', names=['source_node', 'destination_node'])
18 y_te_neg = pd.read_csv('y_te_neg_after_eda.csv', names=['source_node', 'destination_node'])
19 print('='*60)
20 print("Number of nodes in the train data graph with edges", x_train_pos.shape[0])
21 print("Number of nodes in the train data graph without edges", x_train_neg.shape[0])
22 print('='*60)
23 print("Number of nodes in the test data graph with edges", x_test_pos.shape[0])
24 print("Number of nodes in the test data graph without edges", x_test_neg.shape[0])
25
26 x_train = x_train_pos.append(x_train_neg,ignore_index=True)
27 y_train = y_tr_pos.append(y_tr_neg,ignore_index=True)
28 x_test = x_test_pos.append(x_test_neg,ignore_index=True)
29 y_test = y_te_pos.append(y_te_neg,ignore_index=True)
30
31 x_train.to_csv('train_after_eda.csv',header=False,index=False)
32 x_test.to_csv('test_after_eda.csv',header=False,index=False)
33 y_train.to_csv('train_y.csv',header=False,index=False)
34 y_test.to_csv('test_y.csv',header=False,index=False)
35

```

```

=====
Number of nodes in the train data graph with edges 7078139
Number of nodes in the train data graph without edges 7078139
=====
Number of nodes in the test data graph with edges 2359380
Number of nodes in the test data graph without edges 2359380

```

```

In [35]: 1 print("Data points in train data",x_train.shape)
2 print("Data points in test data",x_test.shape)
3 print("Shape of traget variable in train",y_train.shape)
4 print("Shape of traget variable in test", y_test.shape)

```

```

Data points in train data (14156278, 2)
Data points in test data (4718760, 2)
Shape of traget variable in train (14156278, 2)
Shape of traget variable in test (4718760, 2)

```

FB featurization

1. Reading Data

```
In [36]: 1 train_graph=nx.read_edgelist('train_pos_after_eda.csv',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
2 print(nx.info(train_graph))
```

Name:
Type: DiGraph
Number of nodes: 1755007
Number of edges: 7078139
Average in degree: 4.0331
Average out degree: 4.0331

1 # 2. Similarity measures

2.1 Jaccard Distance:

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

```
In [37]: 1 #for followees
2 def jaccard_for_followees(a,b):
3     try:
4         if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.successors(b))) == 0:
5             return 0
6         sim = (len(set(train_graph.successors(a)).intersection(set(train_graph.successors(b))))) /\
7              (len(set(train_graph.successors(a)).union(set(train_graph.successors(b)))))
8     except:
9         return 0
10    return sim
```

```
In [38]: 1 #one test case
2 print(jaccard_for_followees(273084,1505602))
```

0.0

```
In [39]: 1 #node 1635354 not in graph
2 print(jaccard_for_followees(273084,1505602))
```

0.0

```
In [40]: 1 #for followers
2 def jaccard_for_followers(a,b):
3     try:
4         if len(set(train_graph.predecessors(a))) == 0 | len(set(g.predecessors(b))) == 0:
5             return 0
6         sim = (len(set(train_graph.predecessors(a)).intersection(set(train_graph.predecessors(b))))) /\
7              (len(set(train_graph.predecessors(a)).union(set(train_graph.predecessors(b)))))
8         return sim
9     except:
10    return 0
```

In [41]: 1 `print(jaccard_for_followers(273084,470294))`

0.0

In [42]: 1 `#node 1635354 not in graph`
2 `print(jaccard_for_followees(669354,1635354))`

0

2.2 Cosine distance

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

In [43]: 1 `#for followees`
2 `def cosine_for_followees(a,b):`
3 `try:`
4 `if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.successors(b))) == 0:`
5 `return 0`
6 `sim = (len(set(train_graph.successors(a)).intersection(set(train_graph.successors(b)))))/\`
7 `(math.sqrt(len(set(train_graph.successors(a)))*len((set(train_graph.successors(b)))))`
8 `return sim`
9 `except:`
10 `return 0`

In [44]: 1 `print(cosine_for_followees(273084,1505602))`

0.0

In [45]: 1 `print(cosine_for_followees(273084,1635354))`

0

In [46]: 1 `def cosine_for_followers(a,b):`
2 `try:`
3
4 `if len(set(train_graph.predecessors(a))) == 0 | len(set(train_graph.predecessors(b))) == 0:`
5 `return 0`
6 `sim = (len(set(train_graph.predecessors(a)).intersection(set(train_graph.predecessors(b))))) /\`
7 `(math.sqrt(len(set(train_graph.predecessors(a)))*(len(set(train_graph.predecessors(b)))))`
8 `return sim`
9 `except:`
10 `return 0`

In [47]: 1 `print(cosine_for_followers(2,470294))`

0.0

In [48]: `1 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)

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 [49]: `1 if not os.path.isfile('page_rank.p'):
2 pr = nx.pagerank(train_graph, alpha=0.85)
3 pickle.dump(pr,open('page_rank.p','wb'))
4 else:
5 pr = pickle.load(open('page_rank.p','rb'))`

In [50]: `1 print('min',pr[min(pr, key=pr.get)])
2 print('max',pr[max(pr, key=pr.get)])
3 print('mean',float(sum(pr.values())) / len(pr))`

min 1.710937236609678e-07
max 2.713872959435951e-05
mean 5.697982970986667e-07

In [51]: `1 #for imputing to nodes which are not there in Train data
2 mean_pr = float(sum(pr.values())) / len(pr)
3 print(mean_pr)`

5.697982970986667e-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 [52]: 1 #if has direct edge then deleting that edge and calculating shortest path
2 def compute_shortest_path_length(a,b):
3     p=-1
4     try:
5         if train_graph.has_edge(a,b):
6             train_graph.remove_edge(a,b)
7             p = nx.shortest_path_length(train_graph,source=a,target=b)
8             train_graph.add_edge(a,b)
9             return p
10        else:
11            p= nx.shortest_path_length(train_graph,source=a,target=b)
12            return p
13    except:
14        return -1
```

```
In [53]: 1 #testing
2 compute_shortest_path_length(1550756,583691)
```

Out[53]: 3

```
In [54]: 1 #testing
2 compute_shortest_path_length(669354,1635354)
```

Out[54]: -1

4.2 Checking for same community

```

In [55]: ▶ 1 #getting weekly connected edges from graph
2 wcc=list(nx.weakly_connected_components(train_graph))
3 def belongs_to_same_wcc(a,b):
4     index = []
5     if train_graph.has_edge(b,a):
6         return 1
7     if train_graph.has_edge(a,b):
8         for i in wcc:
9             if a in i:
10                index=i
11                break
12            if (b in index):
13                train_graph.remove_edge(a,b)
14                if compute_shortest_path_length(a,b)==-1:
15                    train_graph.add_edge(a,b)
16                    return 0
17                else:
18                    train_graph.add_edge(a,b)
19                    return 1
20            else:
21                return 0
22     else:
23         for i in wcc:
24             if a in i:
25                 index=i
26                 break
27             if(b in index):
28                 return 1
29             else:
30                 return 0

```

```

In [56]: ▶ 1 belongs_to_same_wcc(861, 1659750)

```

```

Out[56]: 0

```

```

In [57]: ▶ 1 belongs_to_same_wcc(669354,1635354)

```

```

Out[57]: 0

```

4.3 Adamic/Adar Index:

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

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

```
In [58]: 1 #adar index
2 def calc_adar_in(a,b):
3     sum=0
4     try:
5         n=list(set(train_graph.successors(a)).intersection(set(train_graph.successors(b))))
6         if len(n)!=0:
7             for i in n:
8                 sum=sum+(1/np.log10(len(list(train_graph.predecessors(i)))))
9         return sum
10    else:
11        return 0
12    except:
13        return 0
```

```
In [59]: 1 calc_adar_in(1,189226)
```

```
Out[59]: 0
```

```
In [60]: 1 calc_adar_in(669354,1635354)
```

```
Out[60]: 0
```

4.4 If person was following back:

```
In [61]: 1 def follows_back(a,b):
2         if train_graph.has_edge(b,a):
3             return 1
4         else:
5             return 0
```

```
In [62]: 1 follows_back(1,189226)
```

```
Out[62]: 1
```

```
In [63]: 1 follows_back(669354,1635354)
```

```
Out[63]: 0
```

4.5 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

λ

The parameter

controls the initial centrality and

$$\beta$$

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

```
In [64]: 1 if not os.path.isfile('katz.p'):
2         katz = nx.katz.katz_centrality(train_graph,alpha=0.005,beta=1)
3         pickle.dump(katz,open('katz.p','wb'))
4     else:
5         katz = pickle.load(open('katz.p','rb'))
```

```
In [65]: 1 print('min',katz[min(katz, key=katz.get)])
2         print('max',katz[max(katz, key=katz.get)])
3         print('mean',float(sum(katz.values())) / len(katz))
```

```
min 0.0007377861743833067
max 0.0030406227082377937
mean 0.0007540061180221656
```

```
In [66]: 1 mean_katz = float(sum(katz.values())) / len(katz)
2         print(mean_katz)
```

```
0.0007540061180221656
```

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)

```
In [67]: 1 if not os.path.isfile('hits.p'):
2         hits = nx.hits(train_graph, max_iter=100, tol=1e-08, nstart=None, normalized=True)
3         pickle.dump(hits,open('hits.p','wb'))
4     else:
5         hits = pickle.load(open('hits.p','rb'))
```

```
In [68]: 1 print('min',hits[0][min(hits[0], key=hits[0].get)])
2         print('max',hits[0][max(hits[0], key=hits[0].get)])
3         print('mean',float(sum(hits[0].values())) / len(hits[0]))
```

```
min 0.0
max 0.00520119619968776
mean 5.697982971018867e-07
```

5. Featurization

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

```
In [69]: 1 import random
2 if os.path.isfile('train_after_eda.csv'):
3     filename = "train_after_eda.csv"
4     # you uncomment this line, if you dont know the lentgh of the file name
5     # here we have hardcoded the number of lines as 15100030
6     n_train = sum(1 for line in open(filename)) #number of records in file (excludes header)
7     s = 100000 #desired sample size
8     skip_train = sorted(random.sample(range(1,n_train+1),n_train-s))
9     #https://stackoverflow.com/a/22259008/4084039
```

```
In [70]: 1 if os.path.isfile('train_after_eda.csv'):
2     filename = "test_after_eda.csv"
3     # you uncomment this line, if you dont know the lentgh of the file name
4     # here we have hardcoded the number of lines as 3775008
5     n_test = sum(1 for line in open(filename)) #number of records in file (excludes header)
6     s = 50000 #desired sample size
7     skip_test = sorted(random.sample(range(1,n_test+1),n_test-s))
8     #https://stackoverflow.com/a/22259008/4084039
```

```
In [71]: 1 print("Number of rows in the train data file:", n_train)
2 print("Number of rows we are going to elimiate in train data are",len(skip_train))
3 print("Number of rows in the test data file:", n_test)
4 print("Number of rows we are going to elimiate in test data are",len(skip_test))
```

Number of rows in the train data file: 14156278
 Number of rows we are going to elimiate in train data are 14056278
 Number of rows in the test data file: 4718760
 Number of rows we are going to elimiate in test data are 4668760

```
In [72]: 1 df_final_train = pd.read_csv('train_after_eda.csv', skiprows=skip_train, names=['source_node', 'destination_node'])
2 df_final_train['indicator_link'] = pd.read_csv('train.y.csv', skiprows=skip_train, names=['indicator_link'],index_col=False)
3 print("Our train matrix size ",df_final_train.shape)
4 df_final_train.head(2)
```

Our train matrix size (100001, 3)

```
Out[72]:
```

	source_node	destination_node	indicator_link
0	1055841	514429	1.0
1	898694	385431	1.0

```
In [73]: 1 df_final_test = pd.read_csv('test_after_eda.csv', skiprows=skip_test, names=['source_node', 'destination_node'])
2 df_final_test['indicator_link'] = pd.read_csv('test.y.csv', skiprows=skip_test, names=['indicator_link'],index_col=False)
3 print("Our test matrix size ",df_final_test.shape)
4 df_final_test.head(2)
```

Our test matrix size (50001, 3)

```
Out[73]:
```

	source_node	destination_node	indicator_link
0	1099069	36203	1.0
1	589270	72825	1.0

```
In [74]: 1
2 #Importing Libraries
3 # please do go through this python notebook:
4 import warnings
5 warnings.filterwarnings("ignore")
6
7 import csv
8 import pandas as pd#pandas to create small dataframes
9 import datetime #Convert to unix time
10 import time #Convert to unix time
11 # if numpy is not installed already : pip3 install numpy
12 import numpy as np#Do arithmetic operations on arrays
13 # matplotlib: used to plot graphs
14 import matplotlib
15 import matplotlib.pyplot as plt
16 import seaborn as sns#Plots
17 from matplotlib import rcParams#Size of plots
18 from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
19 import math
20 import pickle
21 import os
22 # to install xgboost: pip3 install xgboost
23 import xgboost as xgb
24
25 import warnings
26 import networkx as nx
27 import pdb
28 import pickle
29 from pandas import HDFStore, DataFrame
30 from pandas import read_hdf
31 from scipy.sparse.linalg import svds, eigs
32 import gc
33 from tqdm import tqdm
34 from sklearn.ensemble import RandomForestClassifier
35 from sklearn.metrics import f1_score
```

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 [75]:

```
1 if not os.path.isfile('storage_sample_stage1.h5'):  
2     #mapping jaccrd followers to train and test data  
3     df_final_train['jaccard_followers'] = df_final_train.apply(lambda row:  
4         jaccard_for_followers(row['source_node'],row['destination_node']),axis=1)  
5     df_final_test['jaccard_followers'] = df_final_test.apply(lambda row:  
6         jaccard_for_followers(row['source_node'],row['destination_node']),axis=1)  
7  
8     #mapping jaccrd followees to train and test data  
9     df_final_train['jaccard_followees'] = df_final_train.apply(lambda row:  
10        jaccard_for_followees(row['source_node'],row['destination_node']),axis=1)  
11    df_final_test['jaccard_followees'] = df_final_test.apply(lambda row:  
12        jaccard_for_followees(row['source_node'],row['destination_node']),axis=1)  
13  
14  
15    #mapping jaccrd followers to train and test data  
16    df_final_train['cosine_followers'] = df_final_train.apply(lambda row:  
17        cosine_for_followers(row['source_node'],row['destination_node']),axis=1)  
18    df_final_test['cosine_followers'] = df_final_test.apply(lambda row:  
19        cosine_for_followers(row['source_node'],row['destination_node']),axis=1)  
20  
21    #mapping jaccrd followees to train and test data  
22    df_final_train['cosine_followees'] = df_final_train.apply(lambda row:  
23        cosine_for_followees(row['source_node'],row['destination_node']),axis=1)  
24    df_final_test['cosine_followees'] = df_final_test.apply(lambda row:  
25        cosine_for_followees(row['source_node'],row['destination_node']),axis=1)
```

```

In [76]: 1 def compute_features_stage1(df_final):
2         #calculating no of followers followees for source and destination
3         #calculating intersection of followers and followees for source and destination
4         num_followers_s=[]
5         num_followees_s=[]
6         num_followers_d=[]
7         num_followees_d=[]
8         inter_followers=[]
9         inter_followees=[]
10        for i,row in df_final.iterrows():
11            try:
12                s1=set(train_graph.predecessors(row['source_node']))
13                s2=set(train_graph.successors(row['source_node']))
14            except:
15                s1 = set()
16                s2 = set()
17            try:
18                d1=set(train_graph.predecessors(row['destination_node']))
19                d2=set(train_graph.successors(row['destination_node']))
20            except:
21                d1 = set()
22                d2 = set()
23            num_followers_s.append(len(s1))
24            num_followees_s.append(len(s2))
25
26            num_followers_d.append(len(d1))
27            num_followees_d.append(len(d2))
28
29            inter_followers.append(len(s1.intersection(d1)))
30            inter_followees.append(len(s2.intersection(d2)))
31
32        return num_followers_s, num_followers_d, num_followees_s, num_followees_d, inter_followers, inter_followees

```

```

In [77]: 1 if not os.path.isfile('storage_sample_stage1.h5'):
2         df_final_train['num_followers_s'], df_final_train['num_followers_d'], \
3         df_final_train['num_followees_s'], df_final_train['num_followees_d'], \
4         df_final_train['inter_followers'], df_final_train['inter_followees']= compute_features_stage1(df_final_train)
5
6         df_final_test['num_followers_s'], df_final_test['num_followers_d'], \
7         df_final_test['num_followees_s'], df_final_test['num_followees_d'], \
8         df_final_test['inter_followers'], df_final_test['inter_followees']= compute_features_stage1(df_final_test)
9
10        hdf = pd.HDFStore('storage_sample_stage1.h5')
11        hdf.put('train_df',df_final_train, format='table', data_columns=True)
12        hdf.put('test_df',df_final_test, format='table', data_columns=True)
13        hdf.close()
14    else:
15        df_final_train = pd.read_hdf('storage_sample_stage1.h5', 'train_df',mode='r')
16        df_final_test = pd.read_hdf('storage_sample_stage1.h5', 'test_df',mode='r')

```

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 [78]: 1 if not os.path.isfile('storage_sample_stage2.h5'):
2         #mapping adar index on train
3         df_final_train['adar_index'] = df_final_train.apply(lambda row: calc_adar_in(row['source_node'],row['destination_node']),axis=1)
4         #mapping adar index on test
5         df_final_test['adar_index'] = df_final_test.apply(lambda row: calc_adar_in(row['source_node'],row['destination_node']),axis=1)
6
7         #-----
8         #mapping followback or not on train
9         df_final_train['follows_back'] = df_final_train.apply(lambda row: follows_back(row['source_node'],row['destination_node']),axis=1)
10
11        #mapping followback or not on test
12        df_final_test['follows_back'] = df_final_test.apply(lambda row: follows_back(row['source_node'],row['destination_node']),axis=1)
13
14        #-----
15        #mapping same component of wcc or not on train
16        df_final_train['same_comp'] = df_final_train.apply(lambda row: belongs_to_same_wcc(row['source_node'],row['destination_node']),axis=1)
17
18        ##mapping same component of wcc or not on train
19        df_final_test['same_comp'] = df_final_test.apply(lambda row: belongs_to_same_wcc(row['source_node'],row['destination_node']),axis=1)
20
21        #-----
22        #mapping shortest path on train
23        df_final_train['shortest_path'] = df_final_train.apply(lambda row: compute_shortest_path_length(row['source_node'],row['destination_node']),axis=1)
24        #mapping shortest path on test
25        df_final_test['shortest_path'] = df_final_test.apply(lambda row: compute_shortest_path_length(row['source_node'],row['destination_node']),axis=1)
26
27        hdf = pd.HDFStore('storage_sample_stage2.h5')
28        hdf.put('train_df',df_final_train, format='table', data_columns=True)
29        hdf.put('test_df',df_final_test, format='table', data_columns=True)
30        hdf.close()
31    else:
32        df_final_train = pd.read_hdf('storage_sample_stage2.h5', 'train_df',mode='r')
33        df_final_test = pd.read_hdf('storage_sample_stage2.h5', 'test_df',mode='r')
```

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
 - 2*weight of incoming edges + weight of outgoing edges
 - weight of incoming edges + 2*weight of outgoing edges
2. Page Ranking of source
3. Page Ranking of dest
4. katz of source
5. katz of dest
6. hubs of source

7. hubs of dest
8. authorities_s of source
9. authorities_s of dest

Weight Features

In order to determine the similarity of nodes, an edge weight value was calculated between nodes. Edge weight decreases as the neighbor count goes up. Intuitively, consider one million people following a celebrity on a social network then chances are most of them never met each other or the celebrity. On the other hand, if a user has 30 contacts in his/her social network, the chances are higher that many of them know each other.

credit - Graph-based Features for Supervised Link Prediction William Cukierski, Benjamin Hamner, Bo Yang

$$W = \frac{1}{\sqrt{1 + |X|}}$$

it is directed graph so calculated Weighted in and Weighted out differently

```
In [79]: 1 #weight for source and destination of each link
2 weight_in = {}
3 weight_out = {}
4 for i in train_graph.nodes():
5     s1=set(train_graph.predecessors(i))
6     w_in = 1.0/(np.sqrt(1+len(s1)))
7     weight_in[i]=w_in
8
9     s2=set(train_graph.successors(i))
10    w_out = 1.0/(np.sqrt(1+len(s2)))
11    weight_out[i]=w_out
12
13 #for imputing with mean
14 mean_weight_in = np.mean(list(weight_in.values()))
15 mean_weight_out = np.mean(list(weight_out.values()))
```

```
In [80]: 1 if not os.path.isfile('data/fea_sample/storage_sample_stage3.h5'):
2     #mapping to pandas train
3     df_final_train['weight_in'] = df_final_train.destination_node.apply(lambda x: weight_in.get(x,mean_weight_in))
4     df_final_train['weight_out'] = df_final_train.source_node.apply(lambda x: weight_out.get(x,mean_weight_out))
5
6     #mapping to pandas test
7     df_final_test['weight_in'] = df_final_test.destination_node.apply(lambda x: weight_in.get(x,mean_weight_in))
8     df_final_test['weight_out'] = df_final_test.source_node.apply(lambda x: weight_out.get(x,mean_weight_out))
9
10
11    #some features engineerings on the in and out weights
12    df_final_train['weight_f1'] = df_final_train.weight_in + df_final_train.weight_out
13    df_final_train['weight_f2'] = df_final_train.weight_in * df_final_train.weight_out
14    df_final_train['weight_f3'] = (2*df_final_train.weight_in + 1*df_final_train.weight_out)
15    df_final_train['weight_f4'] = (1*df_final_train.weight_in + 2*df_final_train.weight_out)
16
17    #some features engineerings on the in and out weights
18    df_final_test['weight_f1'] = df_final_test.weight_in + df_final_test.weight_out
19    df_final_test['weight_f2'] = df_final_test.weight_in * df_final_test.weight_out
20    df_final_test['weight_f3'] = (2*df_final_test.weight_in + 1*df_final_test.weight_out)
21    df_final_test['weight_f4'] = (1*df_final_test.weight_in + 2*df_final_test.weight_out)
```

```

In [81]: 1 if not os.path.isfile('storage_sample_stage3.h5'):
2
3     #page rank for source and destination in Train and Test
4     #if anything not there in train graph then adding mean page rank
5     df_final_train['page_rank_s'] = df_final_train.source_node.apply(lambda x: pr.get(x, mean_pr))
6     df_final_train['page_rank_d'] = df_final_train.destination_node.apply(lambda x: pr.get(x, mean_pr))
7
8     df_final_test['page_rank_s'] = df_final_test.source_node.apply(lambda x: pr.get(x, mean_pr))
9     df_final_test['page_rank_d'] = df_final_test.destination_node.apply(lambda x: pr.get(x, mean_pr))
10    #=====
11
12    #Katz centrality score for source and destination in Train and test
13    #if anything not there in train graph then adding mean katz score
14    df_final_train['katz_s'] = df_final_train.source_node.apply(lambda x: katz.get(x, mean_katz))
15    df_final_train['katz_d'] = df_final_train.destination_node.apply(lambda x: katz.get(x, mean_katz))
16
17    df_final_test['katz_s'] = df_final_test.source_node.apply(lambda x: katz.get(x, mean_katz))
18    df_final_test['katz_d'] = df_final_test.destination_node.apply(lambda x: katz.get(x, mean_katz))
19    #=====
20
21    #Hits algorithm score for source and destination in Train and test
22    #if anything not there in train graph then adding 0
23    df_final_train['hubs_s'] = df_final_train.source_node.apply(lambda x: hits[0].get(x, 0))
24    df_final_train['hubs_d'] = df_final_train.destination_node.apply(lambda x: hits[0].get(x, 0))
25
26    df_final_test['hubs_s'] = df_final_test.source_node.apply(lambda x: hits[0].get(x, 0))
27    df_final_test['hubs_d'] = df_final_test.destination_node.apply(lambda x: hits[0].get(x, 0))
28    #=====
29
30    #Hits algorithm score for source and destination in Train and Test
31    #if anything not there in train graph then adding 0
32    df_final_train['authorities_s'] = df_final_train.source_node.apply(lambda x: hits[1].get(x, 0))
33    df_final_train['authorities_d'] = df_final_train.destination_node.apply(lambda x: hits[1].get(x, 0))
34
35    df_final_test['authorities_s'] = df_final_test.source_node.apply(lambda x: hits[1].get(x, 0))
36    df_final_test['authorities_d'] = df_final_test.destination_node.apply(lambda x: hits[1].get(x, 0))
37    #=====
38
39    hdf = pd.HDFStore('storage_sample_stage3.h5')
40    hdf.put('train_df', df_final_train, format='table', data_columns=True)
41    hdf.put('test_df', df_final_test, format='table', data_columns=True)
42    hdf.close()
43 else:
44     df_final_train = pd.read_hdf('storage_sample_stage3.h5', 'train_df', mode='r')
45     df_final_test = pd.read_hdf('storage_sample_stage3.h5', 'test_df', mode='r')

```

```

In [82]: 1 if not os.path.isfile('storage_sample_stage3.h5'):
2
3     #page rank for source and destination in Train and Test
4     #if anything not there in train graph then adding mean page rank
5     df_final_train['page_rank_s'] = df_final_train.source_node.apply(lambda x: pr.get(x, mean_pr))
6     df_final_train['page_rank_d'] = df_final_train.destination_node.apply(lambda x: pr.get(x, mean_pr))
7
8     df_final_test['page_rank_s'] = df_final_test.source_node.apply(lambda x: pr.get(x, mean_pr))
9     df_final_test['page_rank_d'] = df_final_test.destination_node.apply(lambda x: pr.get(x, mean_pr))
10    #=====
11
12    #Katz centrality score for source and destination in Train and test
13    #if anything not there in train graph then adding mean katz score
14    df_final_train['katz_s'] = df_final_train.source_node.apply(lambda x: katz.get(x, mean_katz))
15    df_final_train['katz_d'] = df_final_train.destination_node.apply(lambda x: katz.get(x, mean_katz))
16
17    df_final_test['katz_s'] = df_final_test.source_node.apply(lambda x: katz.get(x, mean_katz))
18    df_final_test['katz_d'] = df_final_test.destination_node.apply(lambda x: katz.get(x, mean_katz))
19    #=====
20
21    #Hits algorithm score for source and destination in Train and test
22    #if anything not there in train graph then adding 0
23    df_final_train['hubs_s'] = df_final_train.source_node.apply(lambda x: hits[0].get(x, 0))
24    df_final_train['hubs_d'] = df_final_train.destination_node.apply(lambda x: hits[0].get(x, 0))
25
26    df_final_test['hubs_s'] = df_final_test.source_node.apply(lambda x: hits[0].get(x, 0))
27    df_final_test['hubs_d'] = df_final_test.destination_node.apply(lambda x: hits[0].get(x, 0))
28    #=====
29
30    #Hits algorithm score for source and destination in Train and Test
31    #if anything not there in train graph then adding 0
32    df_final_train['authorities_s'] = df_final_train.source_node.apply(lambda x: hits[1].get(x, 0))
33    df_final_train['authorities_d'] = df_final_train.destination_node.apply(lambda x: hits[1].get(x, 0))
34
35    df_final_test['authorities_s'] = df_final_test.source_node.apply(lambda x: hits[1].get(x, 0))
36    df_final_test['authorities_d'] = df_final_test.destination_node.apply(lambda x: hits[1].get(x, 0))
37    #=====
38
39    hdf = pd.HDFStore('storage_sample_stage3.h5')
40    hdf.put('train_df', df_final_train, format='table', data_columns=True)
41    hdf.put('test_df', df_final_test, format='table', data_columns=True)
42    hdf.close()
43 else:
44     df_final_train = pd.read_hdf('storage_sample_stage3.h5', 'train_df', mode='r')
45     df_final_test = pd.read_hdf('storage_sample_stage3.h5', 'test_df', mode='r')

```

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 [83]: 1 def svd(x, S):
2         try:
3             z = sadj_dict[x]
4             return S[z]
5         except:
6             return [0,0,0,0,0,0]
```

```
In [84]: 1 #for svd features to get feature vector creating a dict node val and inedx in svd vector
2 sadj_col = sorted(train_graph.nodes())
3 sadj_dict = { val:idx for idx,val in enumerate(sadj_col)}
```

```
In [85]: 1 Adj = nx.adjacency_matrix(train_graph,nodelist=sorted(train_graph.nodes())).asfptype()
```

```
In [86]: 1 U, s, V = svds(Adj, k = 6)
2 print('Adjacency matrix Shape',Adj.shape)
3 print('U Shape',U.shape)
4 print('V Shape',V.shape)
5 print('s Shape',s.shape)
```

```
Adjacency matrix Shape (1755007, 1755007)
U Shape (1755007, 6)
V Shape (6, 1755007)
s Shape (6,)
```

```
In [94]: 1 if not os.path.isfile('storage_sample_stage4.h5'):
2         #=====
3
4         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_s_6']] = \
5         df_final_train.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)
6
7         df_final_train[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6']] = \
8         df_final_train.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)
9         #=====
10
11        df_final_train[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6',]] = \
12        df_final_train.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
13
14        df_final_train[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6']] = \
15        df_final_train.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
16        #=====
17
18        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_s_6']] = \
19        df_final_test.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)
20
21        df_final_test[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6']] = \
22        df_final_test.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)
23
24        #=====
25
26        df_final_test[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6',]] = \
27        df_final_test.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
28
29        df_final_test[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6']] = \
30        df_final_test.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
31        #=====
32
33        hdf = HDFStore('storage_sample_stage4.h5')
34        hdf.put('train_df',df_final_train, format='table', data_columns=True)
35        hdf.put('test_df',df_final_test, format='table', data_columns=True)
36        hdf.close()
37    else:
38        df_final_train = pd.read_hdf('storage_sample_stage4.h5', 'train_df',mode='r')
39        df_final_test = pd.read_hdf('storage_sample_stage4.h5', 'test_df',mode='r')
```

```
In [95]: 1 df_final_test.head(2)
```

Out[95]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_3	svd_v_s_4	svd_v_s_5	sv
0	1099069	36203	1.0	0.000000	0.000000	0.000000	0.000000	6	0	8	...	-1.502736e-07	1.956240e-14	4.281981e-13	8.7
1	198329	1522280	1.0	0.571429	0.333333	0.163299	0.503953	24	20	18	...	-1.023363e-14	5.172288e-14	2.810538e-16	4.3

2 rows × 55 columns

5.6 Adding new Feature: Preferential attachment

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

Preferential attachment for followers

```
In [96]: 1 ## number of followers of source * number of followers of destination train
2 num_f_dest_tr = np.array(df_final_train['num_followers_d'])
3 num_f_source_tr = np.array(df_final_train['num_followers_s'])
4 pref_followers = []
5 for i in range(len(num_f_source_tr)):
6     pref_followers.append(num_f_source_tr[i]*num_f_dest_tr[i])
7 df_final_train['preferential_attach_followers'] = pref_followers
8 df_final_train.head(2)
```

Out[96]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_4	svd_v_s_5	svd_v_s_6	...
0	1055841	514429	1.0	0.489796	0.224299	0.078657	0.383816	76	70	46	...	1.915737e-11	1.541002e-14	2.221265e-06	1.6
1	1281867	313127	1.0	0.000000	0.000000	0.000000	0.000000	1	1	2	...	3.765042e-15	1.091447e-13	1.725994e-14	0.0

2 rows × 56 columns

```
In [97]: 1 ## number of followers of source_te * number of followers of destination test
2 num_f_dest_te = np.array(df_final_test['num_followers_d'])
3 num_f_source_te = np.array(df_final_test['num_followers_s'])
4 pref_followers = []
5 for i in range(len(num_f_source_te)):
6     pref_followers.append(num_f_source_te[i]*num_f_dest_te[i])
7 df_final_test['preferential_attach_followers'] = pref_followers
8 df_final_test.head(2)
```

Out[97]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_4	svd_v_s_5	svd_v_s_6	...
0	1099069	36203	1.0	0.000000	0.000000	0.000000	0.000000	6	0	8	...	1.956240e-14	4.281981e-13	8.794428e-16	0.0
1	198329	1522280	1.0	0.571429	0.333333	0.163299	0.503953	24	20	18	...	5.172288e-14	2.810538e-16	4.385595e-17	3.4

2 rows × 56 columns

Preferential attachment for followees


```
In [98]: 1 ## number of followees of source * number of followees of destination train
2 num_f_dest_tr = np.array(df_final_train['num_followees_d'])
3 num_f_source_tr = np.array(df_final_train['num_followees_s'])
4 pref_followees = []
5 for i in range(len(num_f_source_tr)):
6     pref_followees.append(num_f_source_tr[i]*num_f_dest_tr[i])
7 df_final_train['preferential_attach_followees'] = pref_followees
8 df_final_train.head(2)
```

Out[98]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_5	svd_v_s_6	svd_v_d_1
0	1055841	514429	1.0	0.489796	0.224299	0.078657	0.383816	76	70	46	...	1.541002e-14	2.221265e-06	1.629140e-12
1	1281867	313127	1.0	0.000000	0.000000	0.000000	0.000000	1	1	2	...	1.091447e-13	1.725994e-14	0.000000e+00

2 rows × 57 columns

```
In [99]: 1 ## number of followees of source_te * number of followees of dest_teination test
2 num_f_dest_te = np.array(df_final_test['num_followees_d'])
3 num_f_source_te = np.array(df_final_test['num_followees_s'])
4 pref_followees = []
5 for i in range(len(num_f_source_te)):
6     pref_followees.append(num_f_source_te[i]*num_f_dest_te[i])
7 df_final_test['preferential_attach_followees'] = pref_followees
8 df_final_test.head(2)
```

Out[99]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_5	svd_v_s_6	svd_v_d_1
0	1099069	36203	1.0	0.000000	0.000000	0.000000	0.000000	6	0	8	...	4.281981e-13	8.794428e-16	0.000000e+00
1	198329	1522280	1.0	0.571429	0.333333	0.163299	0.503953	24	20	18	...	2.810538e-16	4.385595e-17	3.469511e-16

2 rows × 57 columns

5.7 Adding new Feature: SVD_dot

- * SVD dot is the dot product between source svd and destination svd
- * Dot product of columns a and b in low-rank approximation


```
In [101]: 1 ### now we need to get the dot product of all the source and destination svd
2 svd_dot = []
3
4 for i in tqdm(range(len(s1))):
5     ## a has all the source svd and d has all the destination svd
6     a,d = [],[]
7     a.append(s1[i])
8     a.append(s2[i])
9     a.append(s3[i])
10    a.append(s4[i])
11    a.append(s5[i])
12    a.append(s6[i])
13    a.append(s7[i])
14    a.append(s8[i])
15    a.append(s9[i])
16    a.append(s10[i])
17    a.append(s11[i])
18    a.append(s12[i])
19    d.append(d1[i])
20    d.append(d2[i])
21    d.append(d3[i])
22    d.append(d4[i])
23    d.append(d5[i])
24    d.append(d6[i])
25    d.append(d7[i])
26    d.append(d8[i])
27    d.append(d9[i])
28    d.append(d10[i])
29    d.append(d11[i])
30    d.append(d12[i])
31    svd_dot.append(np.dot(a,d))
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 100001/100001 [00:18<00:00, 5268.75it/s]
```

```

In [103]: 1  ### now we need to get the dot product of all the source and destination svd
          2  svd_dot_test =[]
          3
          4  for i in tqdm(range(len(s1))):
          5      ## a has all the source svd and d has all the destination svd
          6      a,d =[],[]
          7      a.append(s1[i])
          8      a.append(s2[i])
          9      a.append(s3[i])
         10      a.append(s4[i])
         11      a.append(s5[i])
         12      a.append(s6[i])
         13      a.append(s7[i])
         14      a.append(s8[i])
         15      a.append(s9[i])
         16      a.append(s10[i])
         17      a.append(s11[i])
         18      a.append(s12[i])
         19      d.append(d1[i])
         20      d.append(d2[i])
         21      d.append(d3[i])
         22      d.append(d4[i])
         23      d.append(d5[i])
         24      d.append(d6[i])
         25      d.append(d7[i])
         26      d.append(d8[i])
         27      d.append(d9[i])
         28      d.append(d10[i])
         29      d.append(d11[i])
         30      d.append(d12[i])
         31      svd_dot_test.append(np.dot(a,d))

```

```
In [104]: 1 ### Lets append the features in train and test dataset
          2 df_final_train['SVD_dot'] = svd_dot
          3 df_final_test['SVD_dot'] = svd_dot_test
```

In [105]:

▶

```
1 print('*'*50)
2 df_final_train.head(2)
```

Out[105]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_6	svd_v_d_1	svd_v_d_2
0	1055841	514429	1.0	0.489796	0.224299	0.078657	0.383816	76	70	46	...	2.221265e-06	1.629140e-12	-2.823192e-14
1	1281867	313127	1.0	0.000000	0.000000	0.000000	0.000000	1	1	2	...	1.725994e-14	0.000000e+00	0.000000e+00

2 rows × 58 columns

In [106]:

▶

```
1 print('*'*50)
2 df_final_test.head(2)
```

Out[106]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	num_followers_d	num_followees_s	...	svd_v_s_6	svd_v_d_1	svd_v_d_2
0	1099069	36203	1.0	0.000000	0.000000	0.000000	0.000000	6	0	8	...	8.794428e-16	0.000000e+00	0.000000e+00
1	198329	1522280	1.0	0.571429	0.333333	0.163299	0.503953	24	20	18	...	4.385595e-17	3.469511e-16	-1.294069e-15

2 rows × 58 columns

In [107]:

▶

```
1 hdf = HDFStore('storage_sample_stage5.h6')
2 hdf.put('train_df',df_final_train, format='table', data_columns=True)
3 hdf.put('test_df',df_final_test, format='table', data_columns=True)
4 hdf.close()
```

6. FB_Models

6.1 Applying Random Forest Classifier

In [108]:

▶

```
1 #reading
2 from pandas import read_hdf
3 df_final_train = read_hdf('storage_sample_stage5.h6', 'train_df',mode='r')
4 df_final_test = read_hdf('storage_sample_stage5.h6', 'test_df',mode='r')
```

```
In [109]: 1 df_final_train.columns
```

```
Out[109]: 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_attach_followers', 'preferential_attach_followees',  
               'SVD_dot'],  
              dtype='object')
```

```
In [110]: 1 y_train = df_final_train.indicator_link  
          2 y_test = df_final_test.indicator_link
```

```
In [111]: 1 df_final_train.drop(['source_node', 'destination_node', 'indicator_link'],axis=1,inplace=True)  
          2 df_final_test.drop(['source_node', 'destination_node', 'indicator_link'],axis=1,inplace=True)
```

```

In [112]: 1 estimators = [10,50,100,250,450]
2 train_scores = []
3 test_scores = []
4 for i in estimators:
5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
6                                 max_depth=5, 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=i, n_jobs=-1, random_state=25, verbose=0, warm_start=False)
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('Estimators = ', i, 'Train Score', train_sc, 'test Score', test_sc)
16 plt.plot(estimators, train_scores, label='Train Score')
17 plt.plot(estimators, test_scores, label='Test Score')
18 plt.xlabel('Estimators')
19 plt.ylabel('Score')
20 plt.title('Estimators vs score at depth of 5')

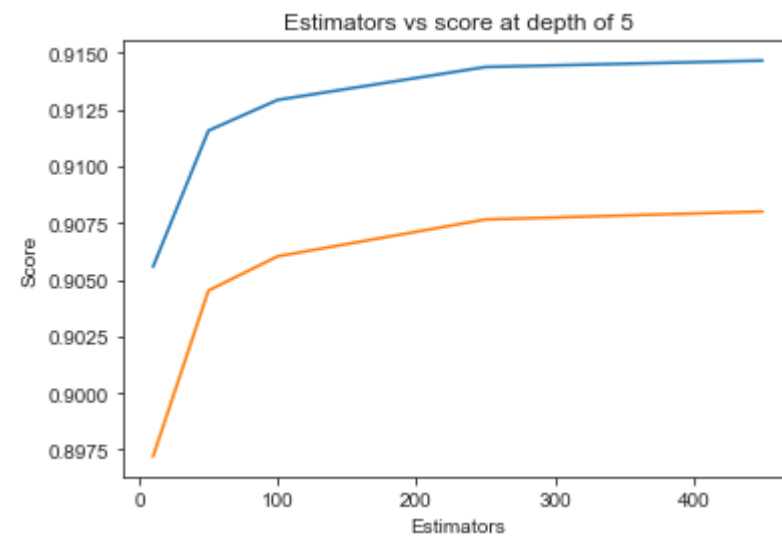
```

```

Estimators = 10 Train Score 0.905570995452 test Score 0.897197842868
Estimators = 50 Train Score 0.91157618673 test Score 0.904524305775
Estimators = 100 Train Score 0.912931976793 test Score 0.906034464576
Estimators = 250 Train Score 0.914383165075 test Score 0.907657848026
Estimators = 450 Train Score 0.914664788318 test Score 0.908008612319

```

Out[112]: Text(0.5, 1.0, 'Estimators vs score at depth of 5')



```

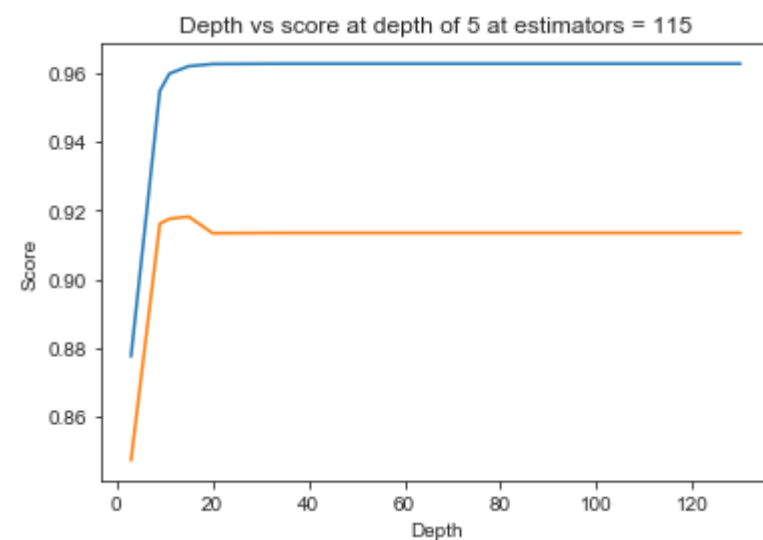
In [113]: 1 depths = [3,9,11,15,20,35,50,70,130]
2 train_scores = []
3 test_scores = []
4 for i in depths:
5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
6                                 max_depth=i, max_features='auto', max_leaf_nodes=None,
7                                 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_state=25, verbose=0, warm_start=False)
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 0.877465845335 test Score 0.847200882907
depth = 9 Train Score 0.954768594276 test Score 0.916144701521
depth = 11 Train Score 0.959819191475 test Score 0.917544604928
depth = 15 Train Score 0.961977806789 test Score 0.918204573416
depth = 20 Train Score 0.962622442843 test Score 0.913369393926
depth = 35 Train Score 0.962707787641 test Score 0.913461947958
depth = 50 Train Score 0.962707787641 test Score 0.913461947958
depth = 70 Train Score 0.962707787641 test Score 0.913461947958
depth = 130 Train Score 0.962707787641 test Score 0.913461947958

```



```
In [116]: 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,return_train_score=True,param_distributions=param_dist,
16                               n_iter=5,cv=10,scoring='f1',random_state=25,verbose=1)
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'])
```

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

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

[Parallel(n_jobs=1)]: Done 50 out of 50 | elapsed: 5.6min finished

mean test scores [0.96072737 0.96063258 0.95916736 0.96037088 0.96205632]

mean train scores [0.96141384 0.96121795 0.95958134 0.96097284 0.96272726]

```
In [139]: 1 print(rf_random.best_estimator_)
```

```
RandomForestClassifier(max_depth=14, min_samples_leaf=28, min_samples_split=111,
                       n_estimators=121, n_jobs=-1, random_state=25)
```

```
In [137]: 1 clf = xgb.XGBClassifier(max_depth=14, min_samples_leaf=28, min_samples_split=111,
2                               n_estimators=121, n_jobs=-1, random_state=25, scoring='f1',
3                               oob_score=False, verbose=0, warm_start=False)
```

```
In [138]: 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)
```

[11:50:39] WARNING: C:\Users\Administrator\workspace\xgboost-win64_release_1.1.0\src\learner.cc:480:
Parameters: { min_samples_leaf, min_samples_split, oob_score, scoring, verbose, warm_start } might not be used.

This may not be accurate due to some parameters are only used in language bindings but passed down to XGBoost core. Or some parameters are not used but slip through this verification. Please open an issue if you find above cases.

In [120]:

```
1 from sklearn.metrics import f1_score
2 print('Train f1 score',f1_score(y_train,y_train_pred))
3 print('Test f1 score',f1_score(y_test,y_test_pred))
```

Train f1 score 0.962789370721

Test f1 score 0.918308067423

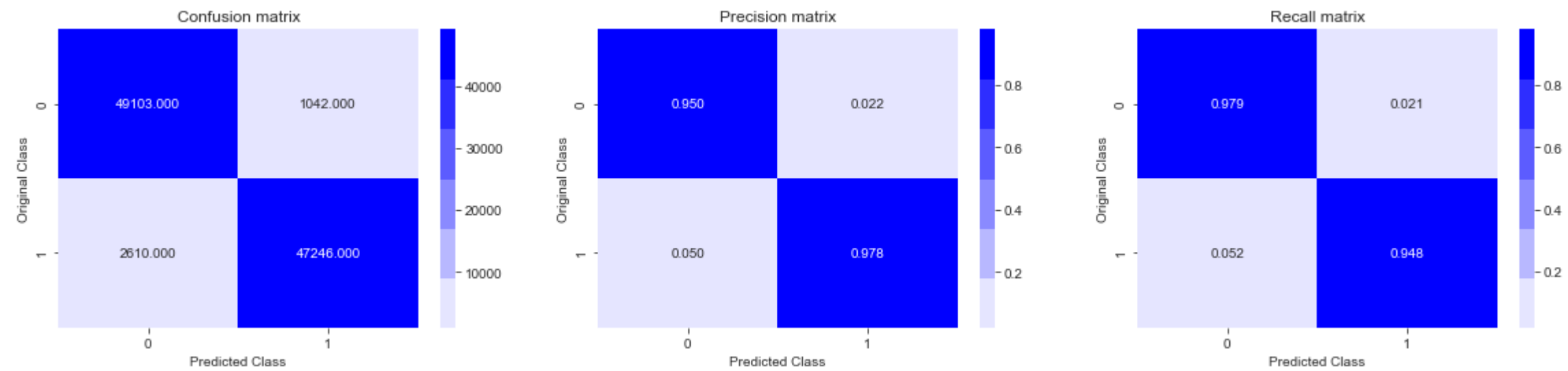
In [121]:

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

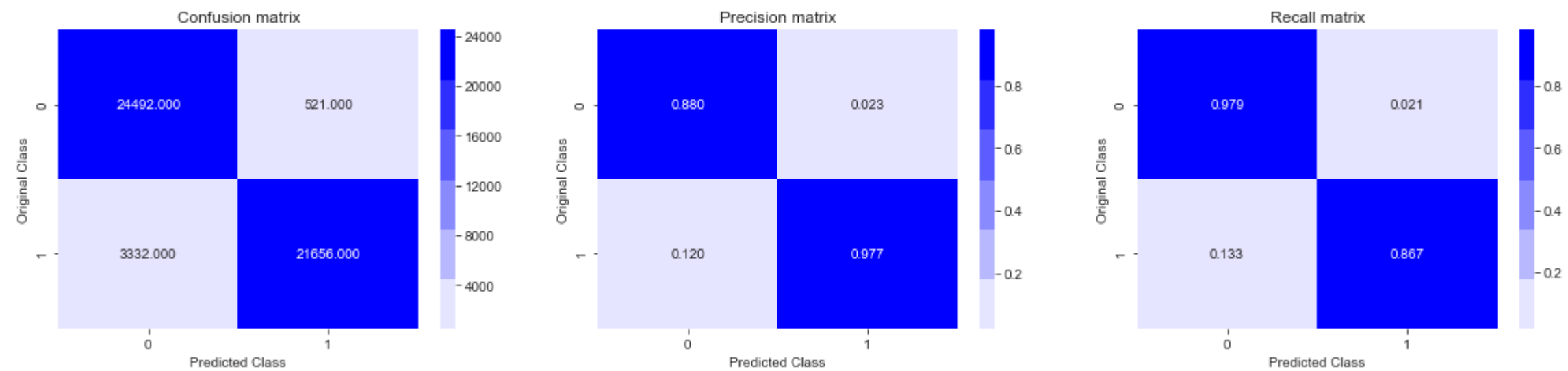


```
In [122]: 1 print('Train confusion_matrix')
2 plot_confusion_matrix(y_train,y_train_pred)
3 print('Test confusion_matrix')
4 plot_confusion_matrix(y_test,y_test_pred)
```

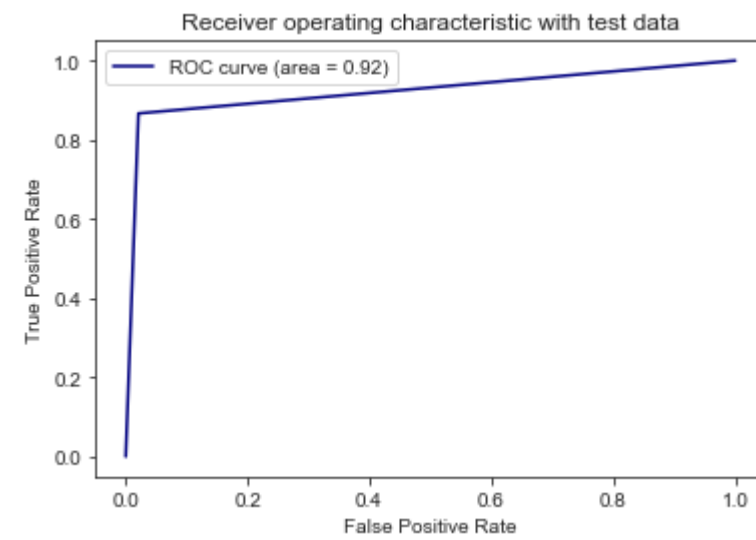
Train confusion_matrix



Test confusion_matrix



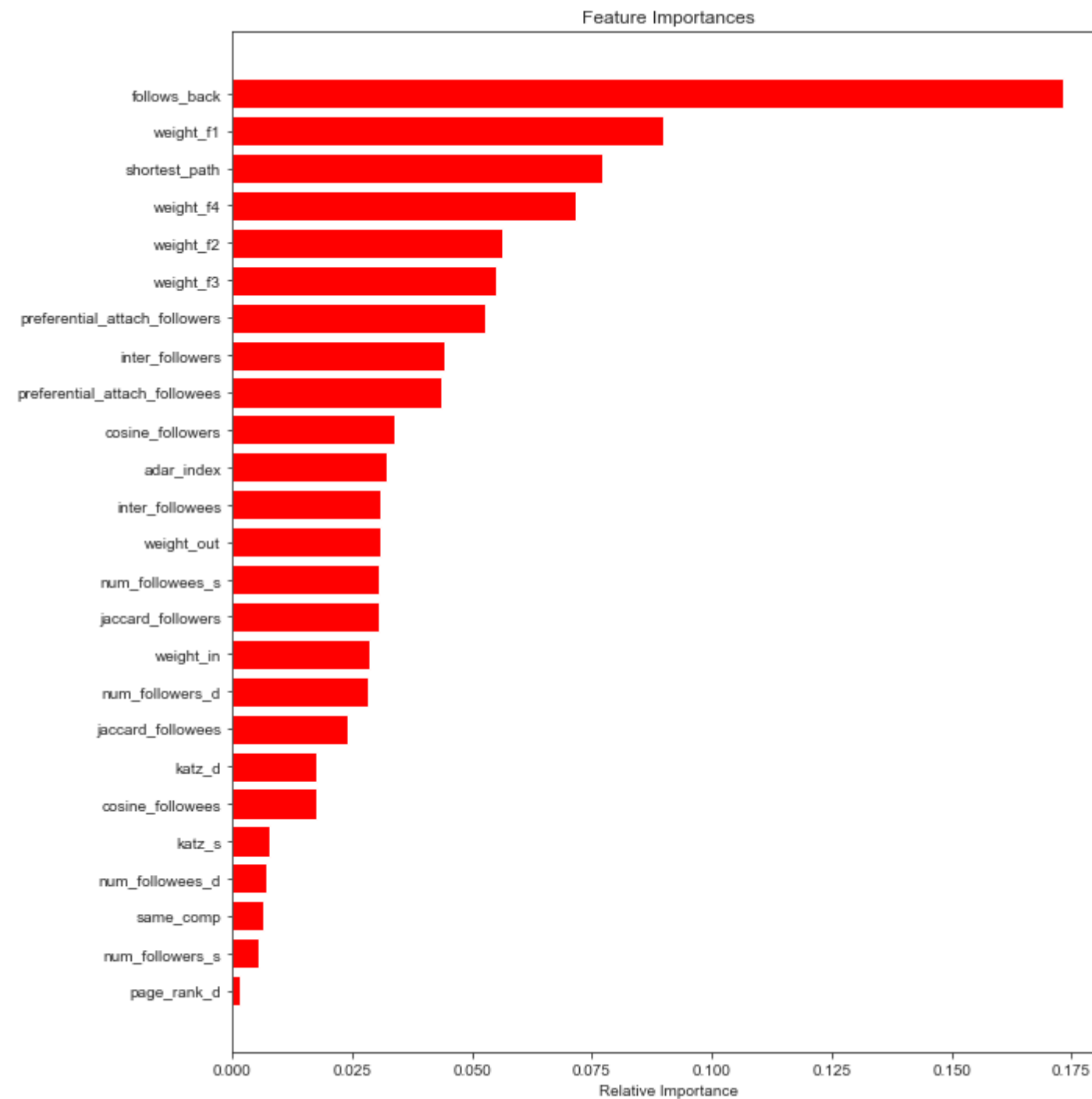
```
In [123]: 1 from sklearn.metrics import roc_curve, auc
2 fpr, tpr, ths = roc_curve(y_test, y_test_pred)
3 auc_sc = auc(fpr, tpr)
4 plt.plot(fpr, tpr, color='navy', label='ROC curve (area = %0.2f)' % auc_sc)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('Receiver operating characteristic with test data')
8 plt.legend()
9 plt.show()
```



```

In [124]: 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()

```



6.2 Applying XGBoost

```
In [152]: 1  ## ref : https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/
2  parameters = { 'max_depth':sp_randint(10,15),
3                 'n_estimators':sp_randint(105,250),
4                 'min_child_weight':range(1,6,2),
5                 'learning_rate': [0.001, 0.01, 0.1, 0.2, 0.3],
6                 'gamma':[0.1,0.2,0.3,0.4,0.5],
7                 'subsample':[0.5, 0.6, 0.7, 0.8, 0.9],
8                 'colsample_bytree':[0.5, 0.6, 0.7, 0.8, 0.9] ,
9                 'reg_alpha':[0.001, 0.005, 0.01, 0.05]
10         }
11  xgb_cl = xgb.XGBClassifier()
12  rs_xgb = RandomizedSearchCV(xgb_cl, parameters,scoring='f1',n_iter=20,verbose=10,
13                             cv=3,refit=False, random_state=42)
14  rs_xgb.fit(df_final_train,y_train)
15
```

Fitting 3 folds for each of 20 candidates, totalling 60 fits

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

[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=14, min_child_weight=1, n_estimators=207, reg_alpha=0.005, subsample=0.7

[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=14, min_child_weight=1, n_estimators=207, reg_alpha=0.005, subsample=0.7, score=0.981, total= 43.4s

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 43.3s remaining: 0.0s

[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=14, min_child_weight=1, n_estimators=207, reg_alpha=0.005, subsample=0.7

[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=14, min_child_weight=1, n_estimators=207, reg_alpha=0.005, subsample=0.7, score=0.982, total= 45.7s

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 1.5min remaining: 0.0s

[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=14, min_child_weight=1, n_estimators=207, reg_alpha=0.005, subsample=0.7

[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=14, min_child_weight=1, n_estimators=207, reg_alpha=0.005, subsample=0.7, score=0.982, total= 45.7s

[Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 2.2min remaining: 0.0s

[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0, max_depth=13, min_child_weight=5, n_estimators=157, reg_alpha=0.005, subsample=0.8

[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0, max_depth=13, min_child_weight=5, n_estimators=157, reg_alpha=0.005, subsample=0.8, score=0.000, total= 30.3s

[Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 2.8min remaining: 0.0s

[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0, max_depth=13, min_child_weight=5, n_estimators=157, reg_alpha=0.005, subsample=0.8

[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0, max_depth=13, min_child_weight=5, n_estimators=157, reg_alpha=0.005, subsample=0.8, score=0.000, total= 30.2s

[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 3.3min remaining: 0.0s

[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0, max_depth=13, min_child_weight=5, n_estimators=157, reg_alpha=0.005, subsample=0.8

[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0, max_depth=13, min_child_weight=5, n_estimators=157, reg_alpha=0.005, subsample=0.8, score=0.000, total= 30.9s

[Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 3.8min remaining: 0.0s

[CV] colsample_bytree=0.6, gamma=0.4, learning_rate=0, max_depth=10, min_child_weight=3, n_estimators=126, reg_alpha=0.001, subsample=0.8

[CV] colsample_bytree=0.6, gamma=0.4, learning_rate=0, max_depth=10, min_child_weight=3, n_estimators=126, reg_alpha=0.001, subsample=0.8, score=0.000, total= 19.8s

[Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 4.1min remaining: 0.0s

[CV] colsample_bytree=0.6, gamma=0.4, learning_rate=0, max_depth=10, min_child_weight=3, n_estimators=126, reg_alpha=0.001, subsample=0.8

[CV] colsample_bytree=0.6, gamma=0.4, learning_rate=0, max_depth=10, min_child_weight=3, n_estimators=126, reg_alpha=0.001, subsample=0.8, score=0.000, total= 19.7s

[Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 4.4min remaining: 0.0s

[CV] colsample_bytree=0.6, gamma=0.4, learning_rate=0, max_depth=10, min_child_weight=3, n_estimators=126, reg_alpha=0.001, subsample=0.8

[CV] colsample_bytree=0.6, gamma=0.4, learning_rate=0, max_depth=10, min_child_weight=3, n_estimators=126, reg_alpha=0.001, subsample=0.8, score=0.000, total= 19.8s

[illegible]

```
[CV] colsample_bytree=0.7, gamma=0.3, learning_rate=0.001, max_depth=12, min_child_weight=1, n_estimators=145, reg_alpha=0.05, subsample=0.5, score=0.976, total= 22.4s
[CV] colsample_bytree=0.8, gamma=0.2, learning_rate=0.001, max_depth=14, min_child_weight=5, n_estimators=208, reg_alpha=0.005, subsample=0.7
[CV] colsample_bytree=0.8, gamma=0.2, learning_rate=0.001, max_depth=14, min_child_weight=5, n_estimators=208, reg_alpha=0.005, subsample=0.7, score=0.975, total= 42.9s
[CV] colsample_bytree=0.8, gamma=0.2, learning_rate=0.001, max_depth=14, min_child_weight=5, n_estimators=208, reg_alpha=0.005, subsample=0.7
[CV] colsample_bytree=0.8, gamma=0.2, learning_rate=0.001, max_depth=14, min_child_weight=5, n_estimators=208, reg_alpha=0.005, subsample=0.7, score=0.975, total= 44.4s
[CV] colsample_bytree=0.8, gamma=0.2, learning_rate=0.001, max_depth=14, min_child_weight=5, n_estimators=208, reg_alpha=0.005, subsample=0.7
[CV] colsample_bytree=0.8, gamma=0.2, learning_rate=0.001, max_depth=14, min_child_weight=5, n_estimators=208, reg_alpha=0.005, subsample=0.7, score=0.977, total= 42.6s
[CV] colsample_bytree=0.7, gamma=0.1, learning_rate=0.1, max_depth=14, min_child_weight=5, n_estimators=235, reg_alpha=0.001, subsample=0.9
[CV] colsample_bytree=0.7, gamma=0.1, learning_rate=0.1, max_depth=14, min_child_weight=5, n_estimators=235, reg_alpha=0.001, subsample=0.9, score=0.982, total= 45.6s
[CV] colsample_bytree=0.7, gamma=0.1, learning_rate=0.1, max_depth=14, min_child_weight=5, n_estimators=235, reg_alpha=0.001, subsample=0.9
[CV] colsample_bytree=0.7, gamma=0.1, learning_rate=0.1, max_depth=14, min_child_weight=5, n_estimators=235, reg_alpha=0.001, subsample=0.9, score=0.982, total= 49.4s
[CV] colsample_bytree=0.7, gamma=0.1, learning_rate=0.1, max_depth=14, min_child_weight=5, n_estimators=235, reg_alpha=0.001, subsample=0.9
[CV] colsample_bytree=0.7, gamma=0.1, learning_rate=0.1, max_depth=14, min_child_weight=5, n_estimators=235, reg_alpha=0.001, subsample=0.9, score=0.982, total= 45.8s
[CV] colsample_bytree=0.6, gamma=0.3, learning_rate=0.001, max_depth=11, min_child_weight=3, n_estimators=228, reg_alpha=0.001, subsample=0.7
[CV] colsample_bytree=0.6, gamma=0.3, learning_rate=0.001, max_depth=11, min_child_weight=3, n_estimators=228, reg_alpha=0.001, subsample=0.7, score=0.974, total= 35.3s
[CV] colsample_bytree=0.6, gamma=0.3, learning_rate=0.001, max_depth=11, min_child_weight=3, n_estimators=228, reg_alpha=0.001, subsample=0.7
[CV] colsample_bytree=0.6, gamma=0.3, learning_rate=0.001, max_depth=11, min_child_weight=3, n_estimators=228, reg_alpha=0.001, subsample=0.7, score=0.973, total= 35.7s
[CV] colsample_bytree=0.6, gamma=0.3, learning_rate=0.001, max_depth=11, min_child_weight=3, n_estimators=228, reg_alpha=0.001, subsample=0.7
[CV] colsample_bytree=0.6, gamma=0.3, learning_rate=0.001, max_depth=11, min_child_weight=3, n_estimators=228, reg_alpha=0.001, subsample=0.7, score=0.974, total= 36.0s
[CV] colsample_bytree=0.5, gamma=0.4, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=247, reg_alpha=0.01, subsample=0.9
[CV] colsample_bytree=0.5, gamma=0.4, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=247, reg_alpha=0.01, subsample=0.9, score=0.000, total= 42.4s
[CV] colsample_bytree=0.5, gamma=0.4, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=247, reg_alpha=0.01, subsample=0.9
[CV] colsample_bytree=0.5, gamma=0.4, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=247, reg_alpha=0.01, subsample=0.9, score=0.000, total= 42.1s
[CV] colsample_bytree=0.5, gamma=0.4, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=247, reg_alpha=0.01, subsample=0.9
[CV] colsample_bytree=0.5, gamma=0.4, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=247, reg_alpha=0.01, subsample=0.9, score=0.000, total= 41.9s
[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=12, min_child_weight=3, n_estimators=132, reg_alpha=0.005, subsample=0.6
[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=12, min_child_weight=3, n_estimators=132, reg_alpha=0.005, subsample=0.6, score=0.981, total= 24.1s
[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=12, min_child_weight=3, n_estimators=132, reg_alpha=0.005, subsample=0.6
[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=12, min_child_weight=3, n_estimators=132, reg_alpha=0.005, subsample=0.6, score=0.981, total= 24.3s
[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=12, min_child_weight=3, n_estimators=132, reg_alpha=0.005, subsample=0.6
[CV] colsample_bytree=0.8, gamma=0.5, learning_rate=0.1, max_depth=12, min_child_weight=3, n_estimators=132, reg_alpha=0.005, subsample=0.6, score=0.981, total= 25.4s
[CV] colsample_bytree=0.9, gamma=0.1, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=179, reg_alpha=0.05, subsample=0.6
[CV] colsample_bytree=0.9, gamma=0.1, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=179, reg_alpha=0.05, subsample=0.6, score=0.000, total= 38.6s
[CV] colsample_bytree=0.9, gamma=0.1, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=179, reg_alpha=0.05, subsample=0.6
[CV] colsample_bytree=0.9, gamma=0.1, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=179, reg_alpha=0.05, subsample=0.6, score=0.000, total= 38.3s
[CV] colsample_bytree=0.9, gamma=0.1, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=179, reg_alpha=0.05, subsample=0.6
[CV] colsample_bytree=0.9, gamma=0.1, learning_rate=0, max_depth=13, min_child_weight=3, n_estimators=179, reg_alpha=0.05, subsample=0.6, score=0.000, total= 38.6s
[CV] colsample_bytree=0.8, gamma=0.1, learning_rate=3, max_depth=10, min_child_weight=1, n_estimators=225, reg_alpha=0.01, subsample=0.5
[CV] colsample_bytree=0.8, gamma=0.1, learning_rate=3, max_depth=10, min_child_weight=1, n_estimators=225, reg_alpha=0.01, subsample=0.5, score=0.000, total= 9.2s
[CV] colsample_bytree=0.8, gamma=0.1, learning_rate=3, max_depth=10, min_child_weight=1, n_estimators=225, reg_alpha=0.01, subsample=0.5
[CV] colsample_bytree=0.8, gamma=0.1, learning_rate=3, max_depth=10, min_child_weight=1, n_estimators=225, reg_alpha=0.01, subsample=0.5, score=0.000, total= 8.8s
[CV] colsample_bytree=0.8, gamma=0.1, learning_rate=3, max_depth=10, min_child_weight=1, n_estimators=225, reg_alpha=0.01, subsample=0.5
[CV] colsample_bytree=0.8, gamma=0.1, learning_rate=3, max_depth=10, min_child_weight=1, n_estimators=225, reg_alpha=0.01, subsample=0.5, score=0.000, total= 8.6s
```

[Parallel(n_jobs=1)]: Done 60 out of 60 | elapsed: 28.9min finished

```
Out[152]: RandomizedSearchCV(cv=3,
                             estimator=XGBClassifier(base_score=None, booster=None,
                                                       colsample_bylevel=None,
                                                       colsample_bynode=None,
                                                       colsample_bytree=None, gamma=None,
                                                       gpu_id=None, importance_type='gain',
                                                       interaction_constraints=None,
                                                       learning_rate=None,
                                                       max_delta_step=None, max_depth=None,
                                                       min_child_weight=None, missing=nan,
                                                       monotone_constraints=None,
                                                       n_estimators=100,...
                             'learning_rate': [0.001, 0.01, 0.1, 0.2,
                                                0, 3],
                             'max_depth': <scipy.stats._distn_infrastructure.rv_frozen object at 0x00000236EB5CFBA8>,
```

```

        'min_child_weight': range(1, 6, 2),
        'n_estimators': <scipy.stats._distn_infrastructure.rv_frozen object at 0x00000236EB5CF630>,
        'reg_alpha': [0.001, 0.005, 0.01, 0.05],
        'subsample': [0.5, 0.6, 0.7, 0.8, 0.9]},
    random_state=42, refit=False, scoring='f1', verbose=10)

```

In [153]: 1 `print(rs_xgb.best_params_)`

```
{'colsample_bytree': 0.7, 'gamma': 0.1, 'learning_rate': 0.1, 'max_depth': 14, 'min_child_weight': 5, 'n_estimators': 235, 'reg_alpha': 0.001, 'subsample': 0.9}
```

In [154]: 1 `clf = xgb.XGBClassifier(colsample_bytree= 0.7, gamma= 0.1, learning_rate= 0.1, max_depth= 14,`
 2 `min_child_weight= 5, n_estimators= 235, reg_alpha= 0.001,`
 3 `subsample= 0.9,n_jobs=-1, random_state=25,)`

In [155]: 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 [156]: 1 `from sklearn.metrics import f1_score`
 2 `print('Train f1 score',f1_score(y_train,y_train_pred))`
 3 `print('Test f1 score',f1_score(y_test,y_test_pred))`

```

Train f1 score 1.0
Test f1 score 0.919521571978

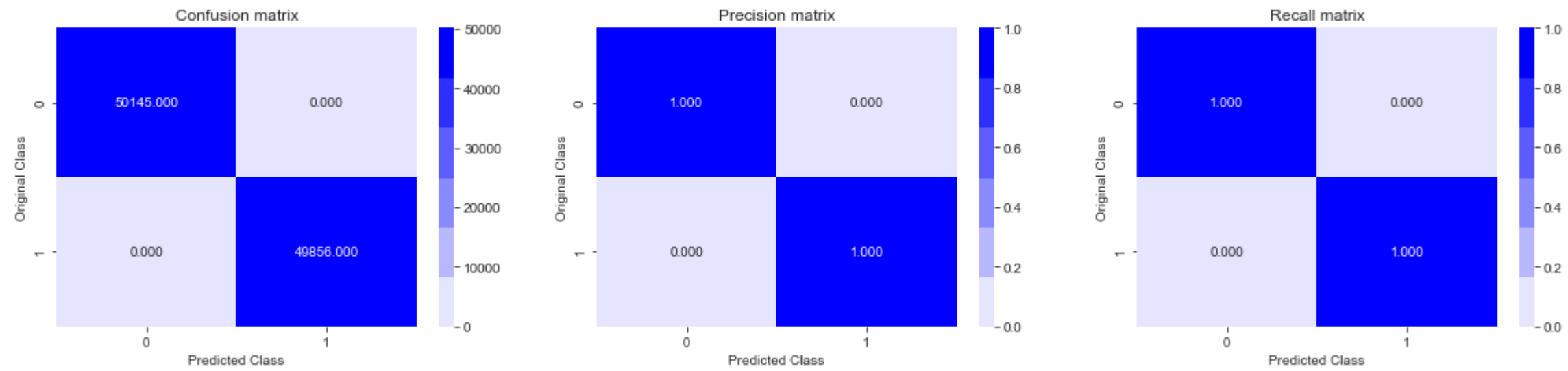
```



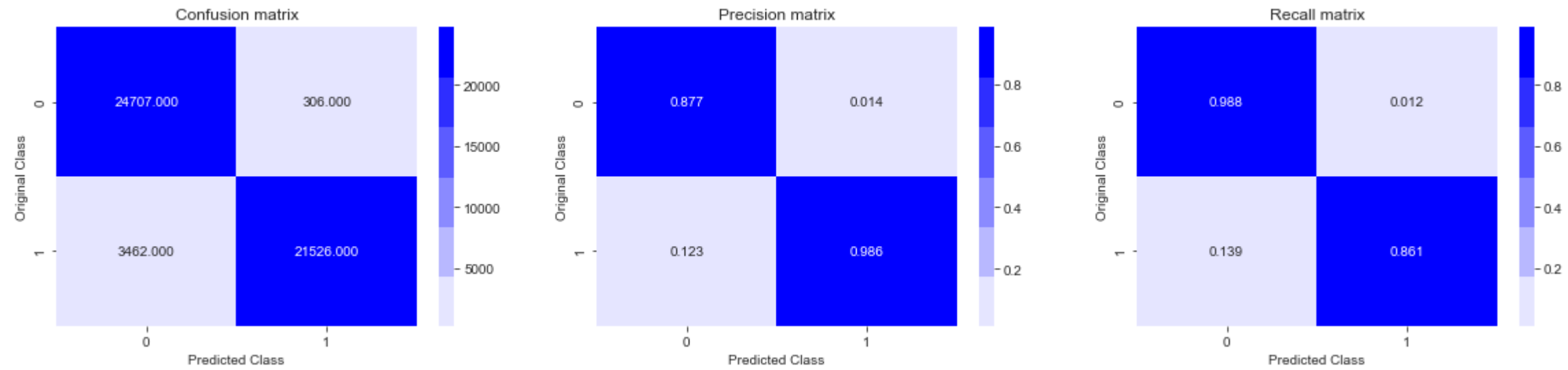
```
In [157]: 1 def plot_confusion_matrix(test_y, predict_y):
2         C = confusion_matrix(test_y, predict_y)
3
4         A = ((C.T)/(C.sum(axis=1))).T
5
6         B = (C/C.sum(axis=0))
7         plt.figure(figsize=(20,4))
8
9         labels = [0,1]
10        # representing A in heatmap format
11        cmap=sns.light_palette("blue")
12        plt.subplot(1, 3, 1)
13        sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
14        plt.xlabel('Predicted Class')
15        plt.ylabel('Original Class')
16        plt.title("Confusion matrix")
17
18        plt.subplot(1, 3, 2)
19        sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
20        plt.xlabel('Predicted Class')
21        plt.ylabel('Original Class')
22        plt.title("Precision matrix")
23
24        plt.subplot(1, 3, 3)
25        # representing B in heatmap format
26        sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
27        plt.xlabel('Predicted Class')
28        plt.ylabel('Original Class')
29        plt.title("Recall matrix")
30
31        plt.show()
```

```
In [158]: 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)
```

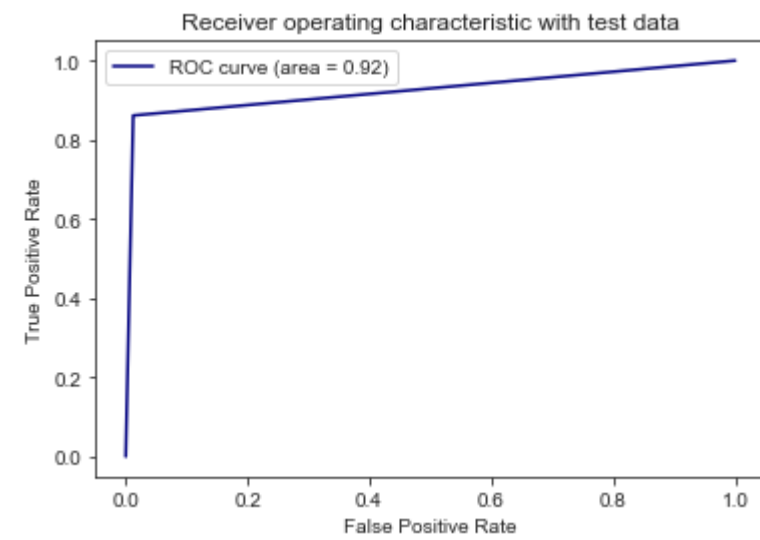
Train confusion_matrix



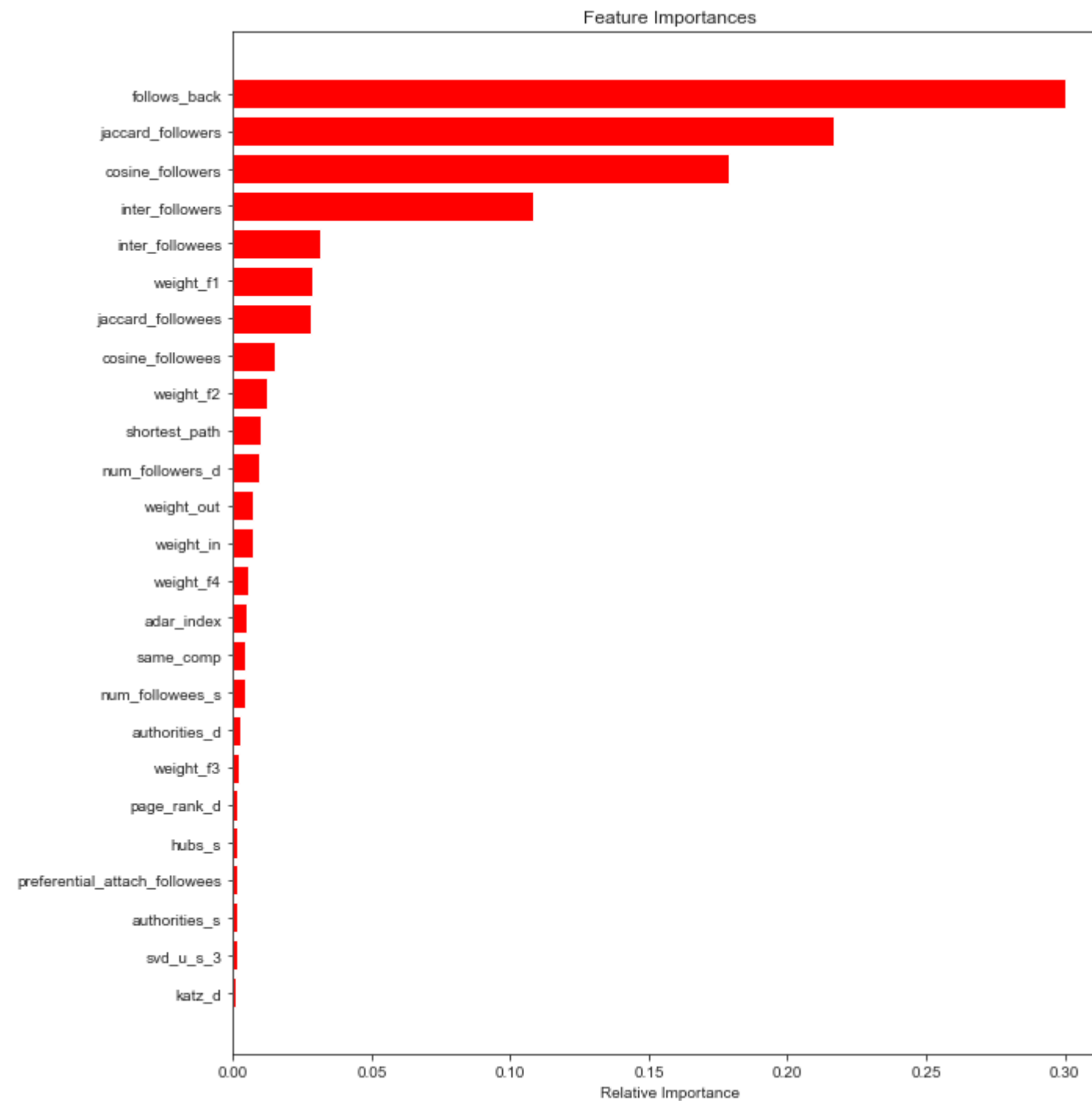
Test confusion_matrix



```
In [159]: 1 fpr,tpr,ths = roc_curve(y_test,y_test_pred)
2 auc_sc = auc(fpr, tpr)
3 plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
4 plt.xlabel('False Positive Rate')
5 plt.ylabel('True Positive Rate')
6 plt.title('Receiver operating characteristic with test data')
7 plt.legend()
8 plt.show()
```



```
In [160]: 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()
```



Conclusion :

```
In [14]: 1  ##http://zetcode.com/python/prettytable/
2
3  from prettytable import PrettyTable
4
5  x = PrettyTable()
6  x.field_names = ["Model", "Hyper Parameter", " Train F1score" ,"Test F1score "]
7  x.add_row(["Random Forest Classifier",'{'max_depth':14,'n_estimators':121}','',0.96,0.918])
8  x.add_row(["XGBoost",'{'max_depth': 14,'n_estimators': 235}','', 1, 0.919])
9  print(x)
```

Model	Hyper Parameter	Train F1score	Test F1score
Random Forest Classifier	{'max_depth':14,'n_estimators':121}	0.96	0.918
XGBoost	{'max_depth': 14,'n_estimators': 235}	1	0.919

Steps:

1. Exploratory Data Analysis

- * 1.1 Number of Followers for each person
- * 1.2 Number of Followees for each person
- * 1.3 Both followers + following

2. Posing a problem as classification problem¶

- * 2.1 Generating some edges which are not present in graph for supervised learning
- * 2.2 Training and Test data split:

FB featurization

- * Reading Data

2. Similarity measures

- * 2.1 Jaccard Distance:
- * 2.2 Cosine distance

3. Ranking Measures

- * 3.1 Page Ranking

4. Other Graph Features

- * 4.1 Shortest path
- * 4.2 Checking for same community
- * 4.3 Adamic/Adar Index
- * 4.4 If person was following back:
- * 4.5 Katz Centrality
- * 4.6 Hits Score

5. Featurization

- * 5. 1 Reading a sample of Data from both train and test
- * 5.2 Adding a set of features:jaccard distance,cosine features,num_followers,inter_followers
- * 5.3 Adding a set of features : adar index,is following back,belongs to same weakly connect components,shortest path between source and destination

5.4 Adding new set of features: weight features

- * 5.5 Adding new set of features: SVD features
- * 5.6 Adding new Feature: Preferential attachment
- * 5.7 Adding new Feature: SVD_dot

6. FB_Models

6.1 Applying Random Forest Classifier:

- * hyperparameter tuning
- * Confusion matrix , precisiona nd recall matrix
- * Roc with test data
- * feature importances

6.1 Applying XGBoost:

- * hyperparameter tuning
- * Confusion matrix , precisiona nd recall matrix
- * Roc with test data
- * feature importances

In []:

▶

1