

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://kaggle2.blob.core.windows.net/forum-message-attachments/2594/supervised_link_prediction.pdf (https://kaggle2.blob.core.windows.net/forum-message-attachments/2594/supervised_link_prediction.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]:

```

#Importing Libraries
# please do go through this python notebook:
import warnings
warnings.filterwarnings("ignore")

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

import warnings
import networkx as nx
import pdb
import pickle

```

In [2]:

```

#reading graph
if not os.path.isfile('train_woheader.csv'):
    traincsv = pd.read_csv('train.csv')
    print(traincsv[traincsv.isna().any(1)])
    print(traincsv.info())
    print("Number of duplicate entries: ",sum(traincsv.duplicated()))
    traincsv.to_csv('train_woheader.csv',header=False,index=False)
    print("saved the graph into file")
else:
    g=nx.read_edgelist('train_woheader.csv',delimiter=',',create_using=nx.DiGraph(),nod
etype=int)
    print(nx.info(g))

```

Name:

Type: DiGraph

Number of nodes: 1862220

Number of edges: 9437519

Average in degree: 5.0679

Average out degree: 5.0679

Displaying a sub graph

In [3]:

```
if not os.path.isfile('train_woheader_sample.csv'):
    pd.read_csv('train.csv', nrows=50).to_csv('train_woheader_sample.csv', header=False,
index=False)

subgraph=nx.read_edgelist('train_woheader_sample.csv',delimiter=',',create_using=nx.DiG
raph(),nodetype=int)
# https://stackoverflow.com/questions/9402255/drawing-a-huge-graph-with-networkx-and-ma
tplotlib

pos=nx.spring_layout(subgraph)
nx.draw(subgraph,pos,node_color='#A0CBE2',edge_color='#00bb5e',width=1,edge_cmap=plt.cm
.Blues,with_labels=True)
plt.savefig("graph_sample.pdf")
print(nx.info(subgraph))
```

Name:

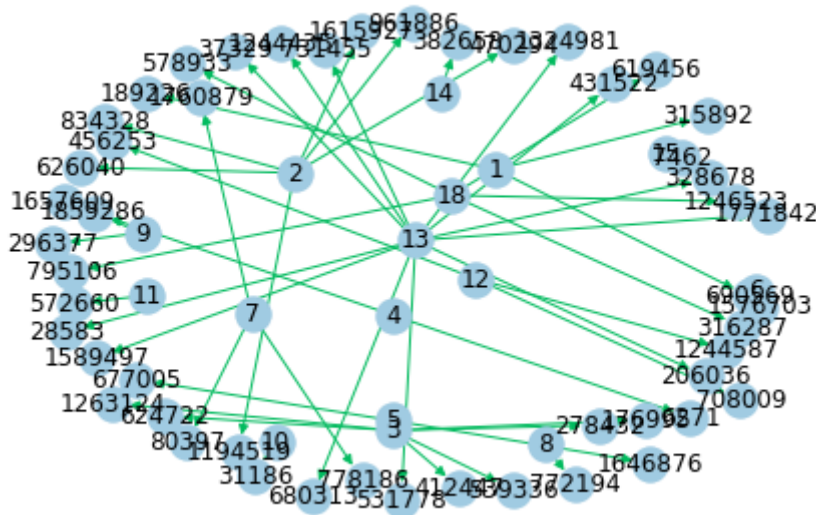
Type: DiGraph

Number of nodes: 66

Number of edges: 50

Average in degree: 0.7576

Average out degree: 0.7576



1. Exploratory Data Analysis

In [4]:

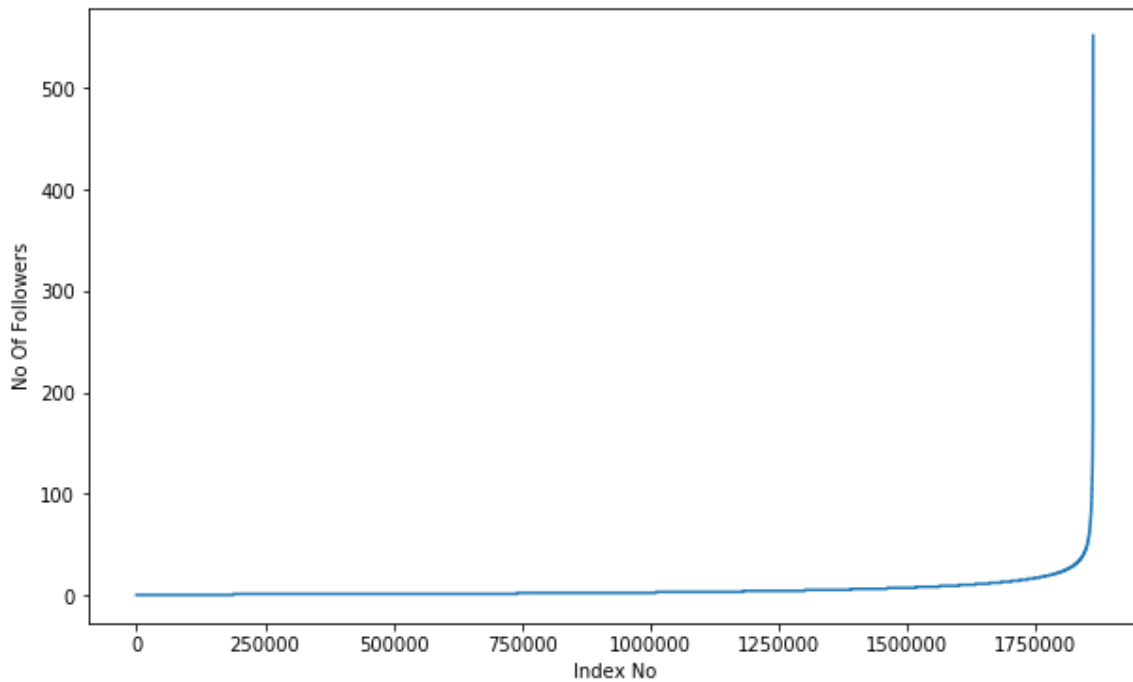
```
# No of Unique persons
print("The number of unique persons",len(g.nodes()))
```

The number of unique persons 1862220

1.1 No of followers for each person

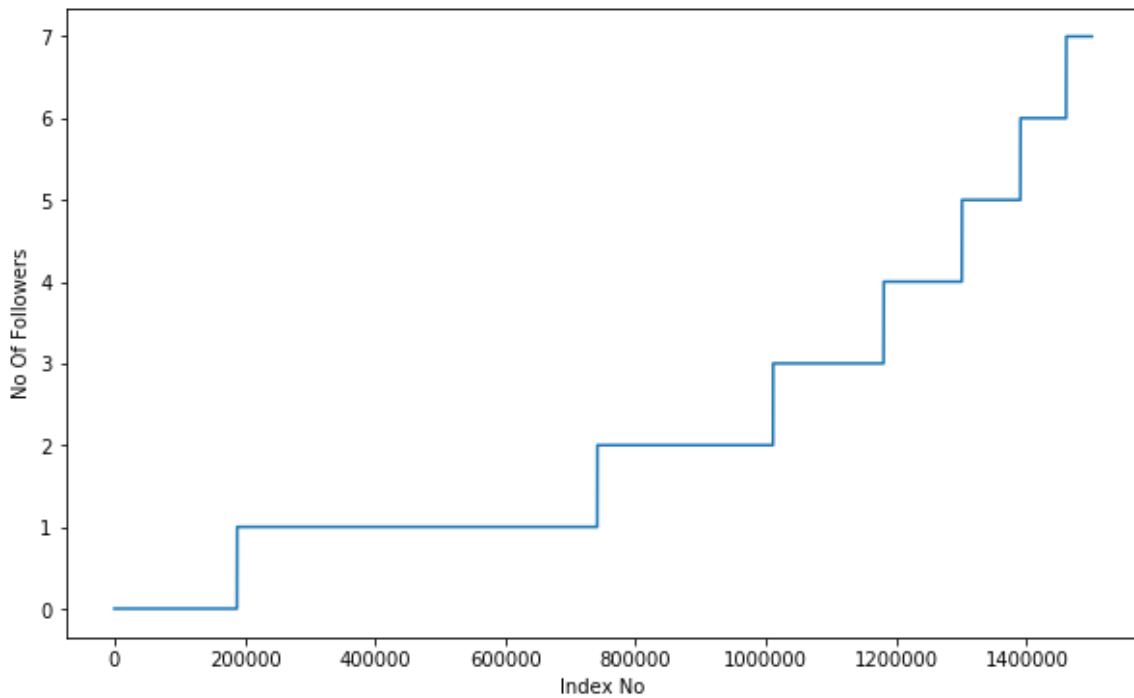
In [5]:

```
indegree_dist = list(dict(g.in_degree()).values())  
indegree_dist.sort()  
plt.figure(figsize=(10,6))  
plt.plot(indegree_dist)  
plt.xlabel('Index No')  
plt.ylabel('No Of Followers')  
plt.show()
```



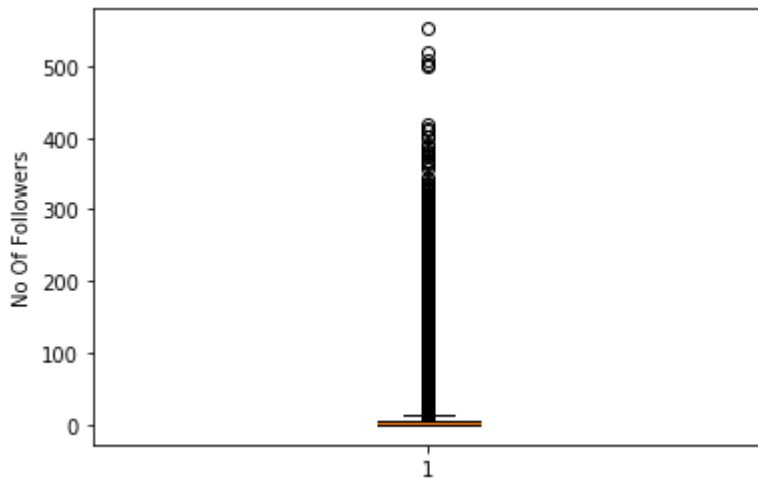
In [6]:

```
indegree_dist = list(dict(g.in_degree()).values())
indegree_dist.sort()
plt.figure(figsize=(10,6))
plt.plot(indegree_dist[0:1500000])
plt.xlabel('Index No')
plt.ylabel('No Of Followers')
plt.show()
```



In [7]:

```
plt.boxplot(indegree_dist)
plt.ylabel('No Of Followers')
plt.show()
```



In [8]:

```
### 90-100 percentile
for i in range(0,11):
    print(90+i, 'percentile value is', np.percentile(indegree_dist, 90+i))
```

```
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% of data having followers of 40 only.

In [9]:

```

### 99-100 percentile
for i in range(10,110,10):
    print(99+(i/100), 'percentile value is', np.percentile(indegree_dist, 99+(i/100)))

```

```

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

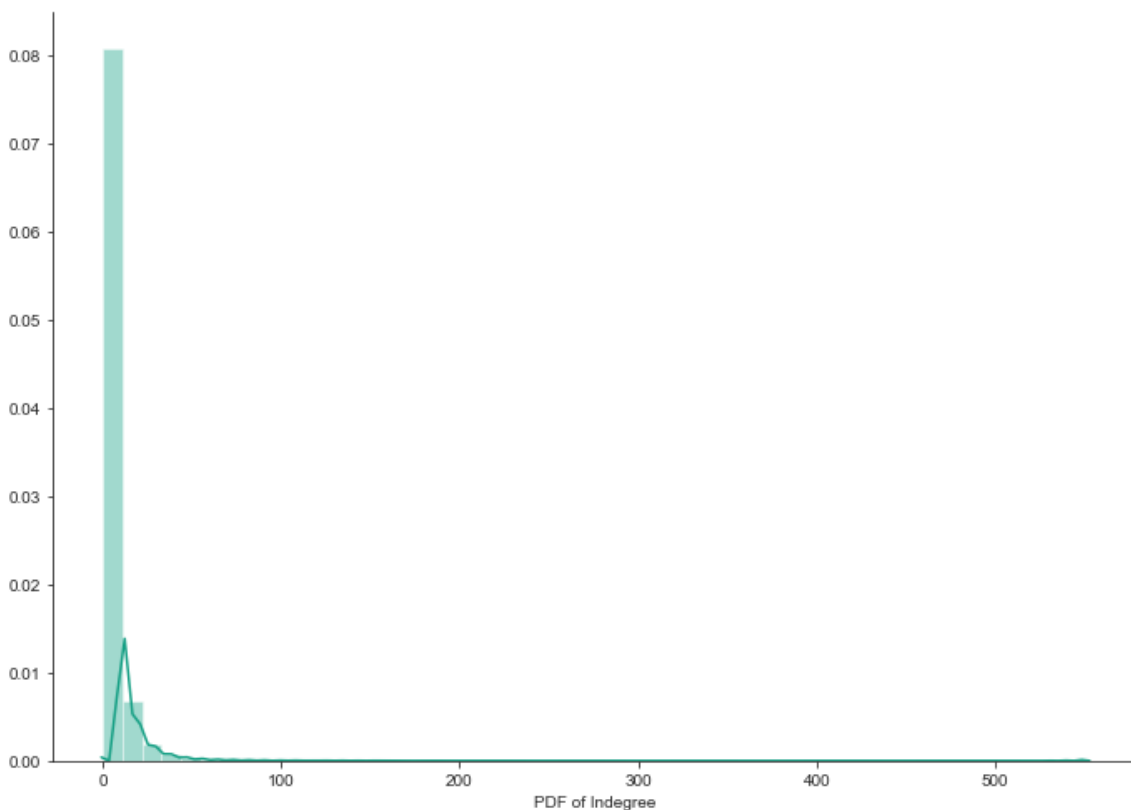
```

In [10]:

```

%matplotlib inline
sns.set_style('ticks')
fig, ax = plt.subplots()
fig.set_size_inches(11.7, 8.27)
sns.distplot(indegree_dist, color='#16A085')
plt.xlabel('PDF of Indegree')
sns.despine()
#plt.show()

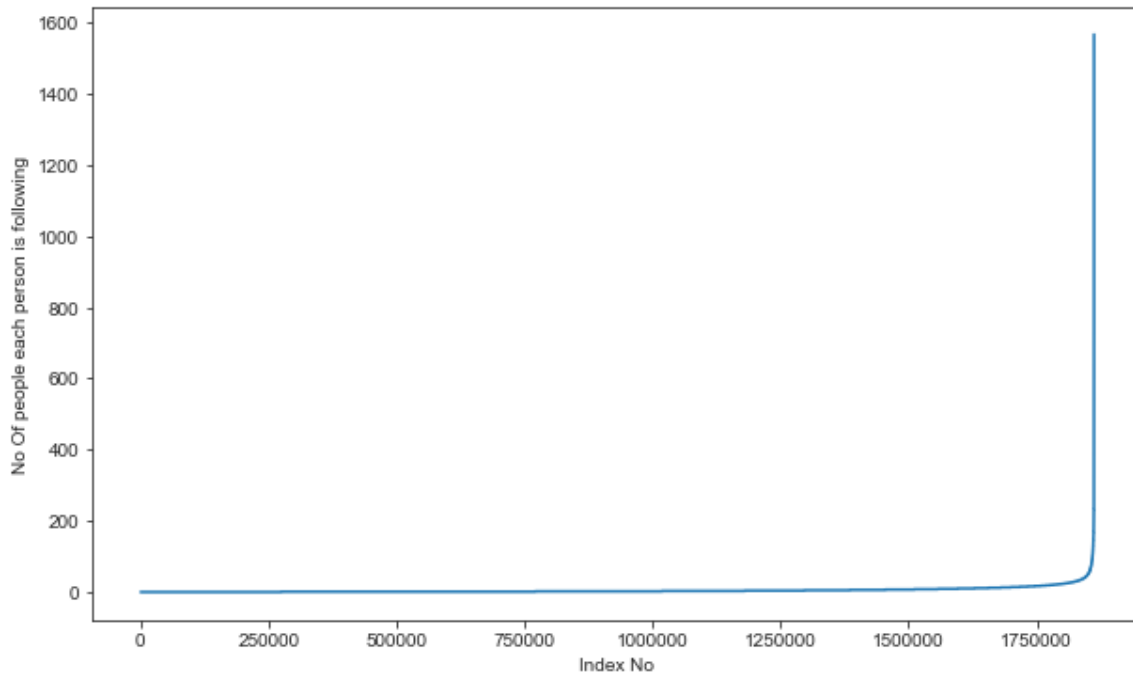
```



1.2 No of people each person is following

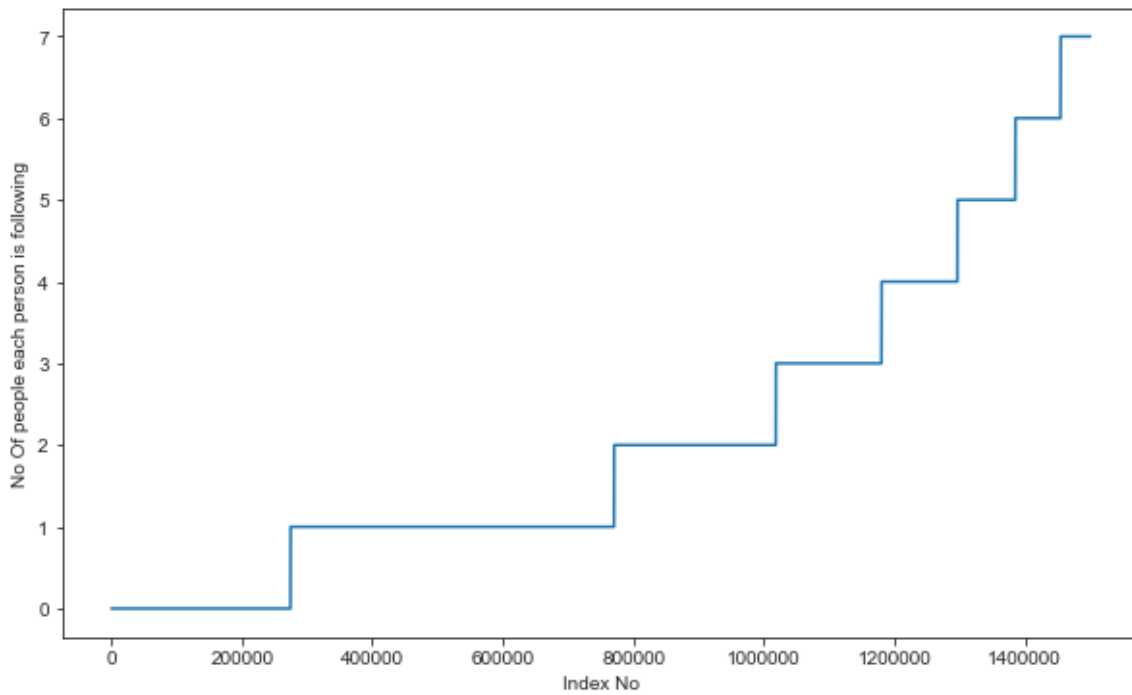
In [11]:

```
outdegree_dist = list(dict(g.out_degree()).values())
outdegree_dist.sort()
plt.figure(figsize=(10,6))
plt.plot(outdegree_dist)
plt.xlabel('Index No')
plt.ylabel('No Of people each person is following')
plt.show()
```



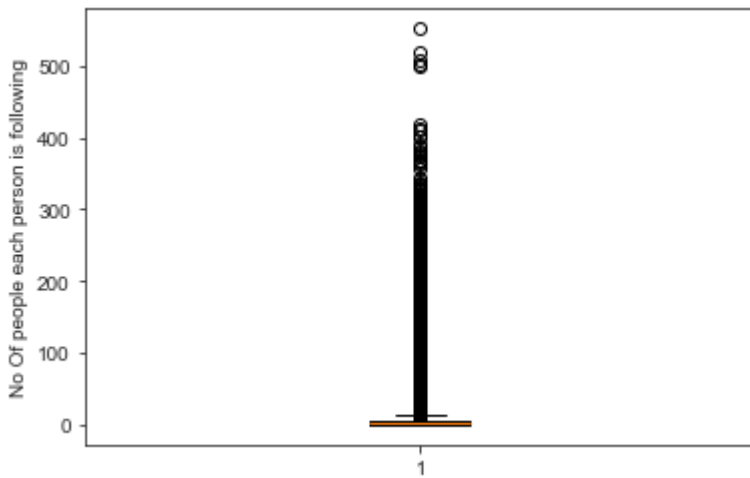
In [12]:

```
indegree_dist = list(dict(g.in_degree()).values())  
indegree_dist.sort()  
plt.figure(figsize=(10,6))  
plt.plot(outdegree_dist[0:1500000])  
plt.xlabel('Index No')  
plt.ylabel('No Of people each person is following')  
plt.show()
```



In [13]:

```
plt.boxplot(indegree_dist)
plt.ylabel('No Of people each person is following')
plt.show()
```



In [14]:

```
### 90-100 percentile
for i in range(0,11):
    print(90+i, 'percentile value is', np.percentile(outdegree_dist, 90+i))
```

```
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
```

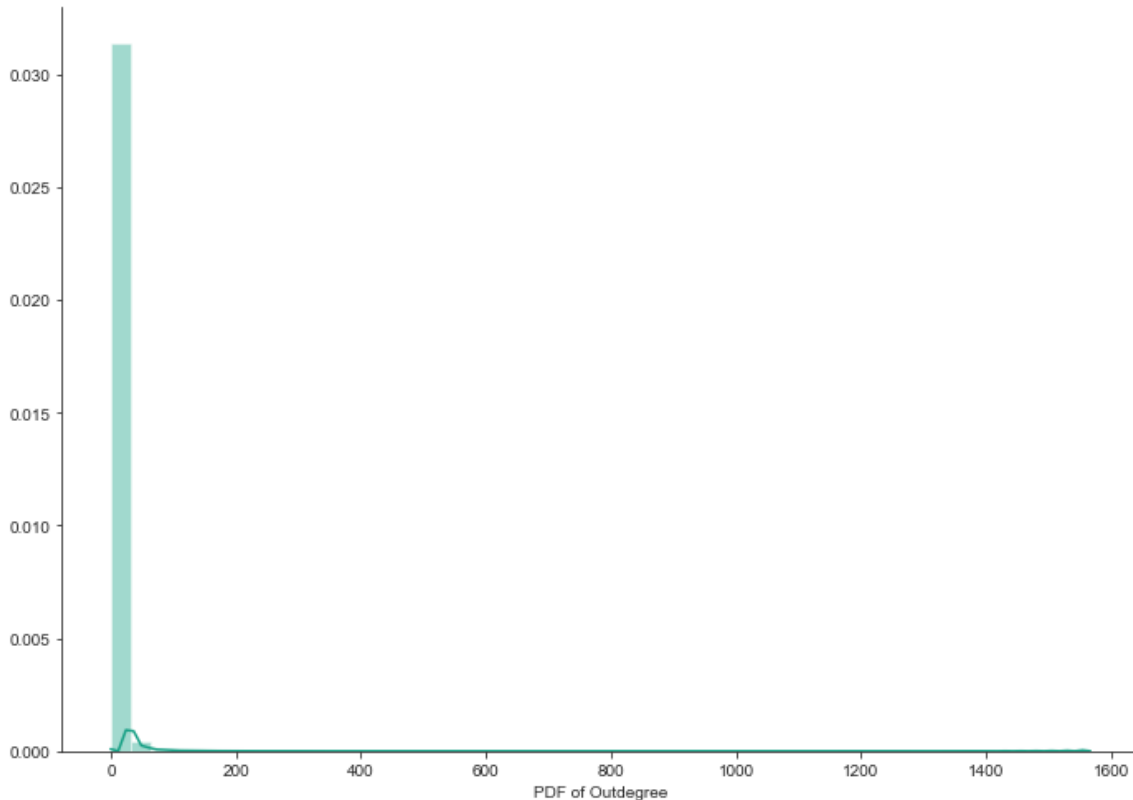
In [15]:

```
### 99-100 percentile
for i in range(10,110,10):
    print(99+(i/100), 'percentile value is', np.percentile(outdegree_dist, 99+(i/100)))
```

```
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 [16]:

```
sns.set_style('ticks')
fig, ax = plt.subplots()
fig.set_size_inches(11.7, 8.27)
sns.distplot(outdegree_dist, color='#16A085')
plt.xlabel('PDF of Outdegree')
sns.despine()
```



In [17]:

```
print('No of persons those are not following anyone are' ,sum(np.array(outdegree_dist)=
=0),'and % is',
      sum(np.array(outdegree_dist)==0)*100/len(outdegree_dist)
) )
```

No of persons those are not following anyone are 274512 and % is 14.741115
442858524

In [18]:

```
print('No of persons having zero followers are' ,sum(np.array(indegree_dist)==0),'and %
is',
      sum(np.array(indegree_dist)==0)*100/len(indegree_dist)
) )
```

No of persons having zero followers are 188043 and % is 10.097786512871734

In [19]:

```

count=0
for i in g.nodes():
    if len(list(g.predecessors(i)))==0 :
        if len(list(g.successors(i)))==0:
            count+=1
print('No of persons those are not not following anyone and also not having any followers are',count)

```

No of persons those are not not following anyone and also not having any followers are 0

1.3 both followers + following

In [20]:

```

from collections import Counter
dict_in = dict(g.in_degree())
dict_out = dict(g.out_degree())
d = Counter(dict_in) + Counter(dict_out)
in_out_degree = np.array(list(d.values()))

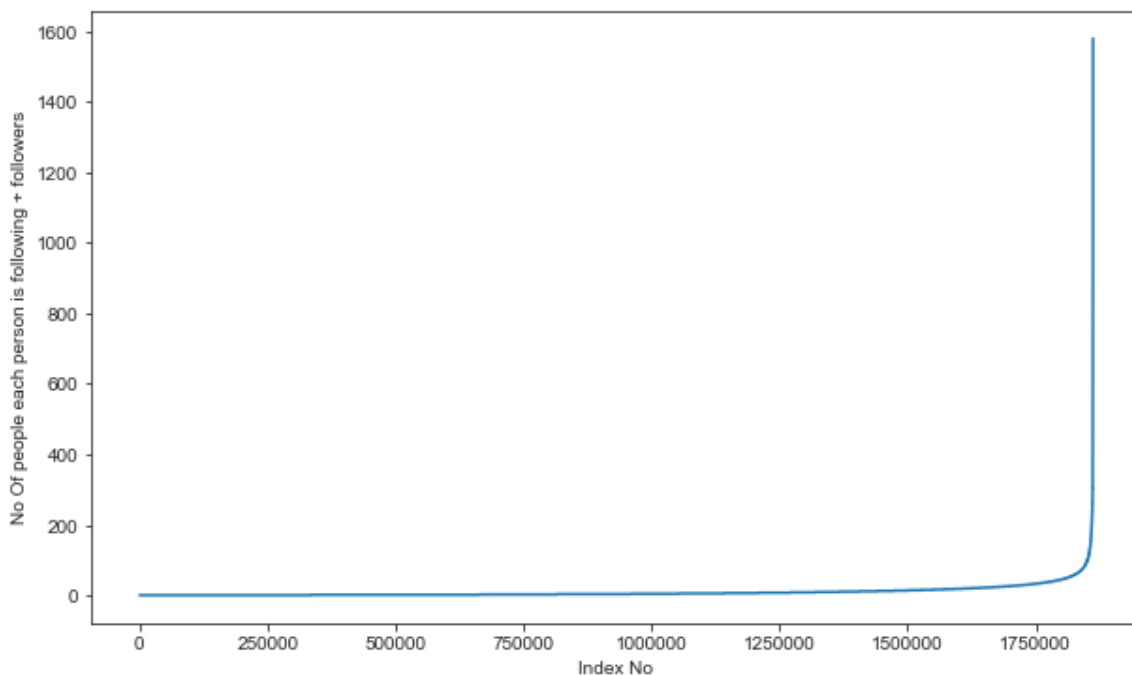
```

In [21]:

```

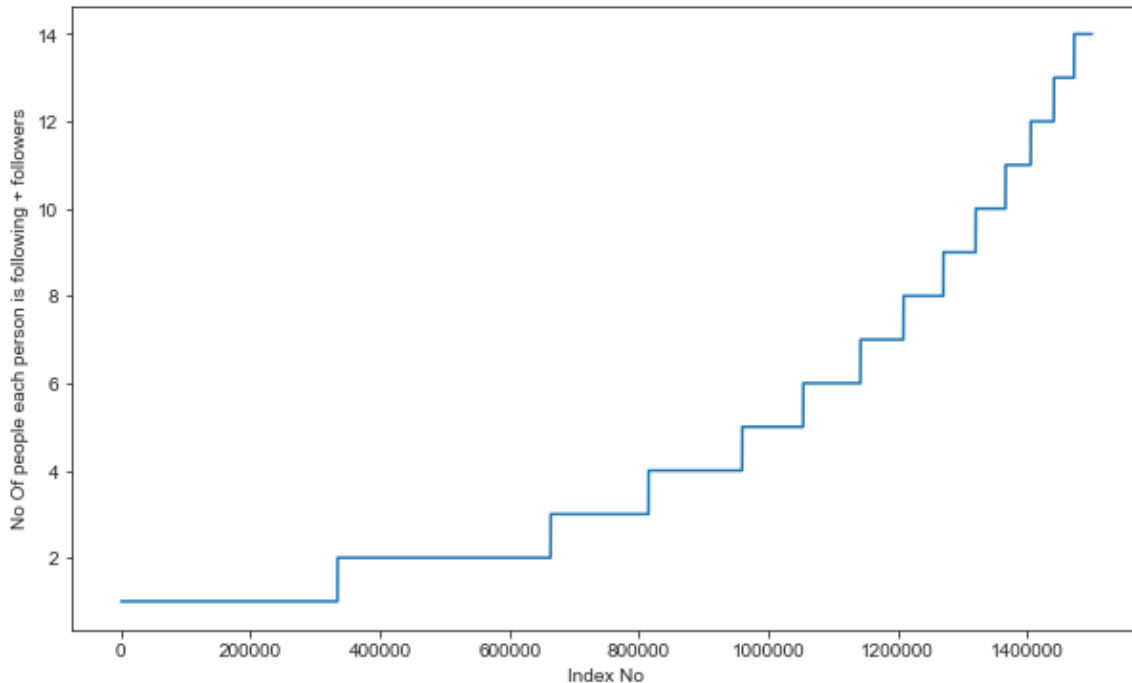
in_out_degree_sort = sorted(in_out_degree)
plt.figure(figsize=(10,6))
plt.plot(in_out_degree_sort)
plt.xlabel('Index No')
plt.ylabel('No Of people each person is following + followers')
plt.show()

```



In [22]:

```
in_out_degree_sort = sorted(in_out_degree)
plt.figure(figsize=(10,6))
plt.plot(in_out_degree_sort[0:1500000])
plt.xlabel('Index No')
plt.ylabel('No Of people each person is following + followers')
plt.show()
```



In [23]:

```
### 90-100 percentile
for i in range(0,11):
    print(90+i,'percentile value is',np.percentile(in_out_degree_sort,90+i))
```

```
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
```

In [24]:

```
### 99-100 percentile
for i in range(10,110,10):
    print(99+(i/100), 'percentile value is', np.percentile(in_out_degree_sort, 99+(i/100)))
```

```
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 [25]:

```
print('Min of no of followers + following is', in_out_degree.min())
print(np.sum(in_out_degree==in_out_degree.min()), ' persons having minimum no of followe
rs + following')
```

```
Min of no of followers + following is 1
334291 persons having minimum no of followers + following
```

In [26]:

```
print('Max of no of followers + following is', in_out_degree.max())
print(np.sum(in_out_degree==in_out_degree.max()), ' persons having maximum no of followe
rs + following')
```

```
Max of no of followers + following is 1579
1 persons having maximum no of followers + following
```

In [27]:

```
print('No of persons having followers + following less than 10 are', np.sum(in_out_degre
e<10))
```

```
No of persons having followers + following less than 10 are 1320326
```

In [28]:

```
print('No of weakly connected components', len(list(nx.weakly_connected_components(g))))
count=0
for i in list(nx.weakly_connected_components(g)):
    if len(i)==2:
        count+=1
print('weakly connected components wit 2 nodes', count)
```

```
No of weakly connected components 45558
weakly connected components wit 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 [29]:

```
%%time
###generating bad edges from given graph
import random
if not os.path.isfile('missing_edges_final.p'):
    #getting all set of edges
    r = csv.reader(open('train_woheader.csv','r'))
    edges = dict()
    for edge in r:
        edges[(edge[0], edge[1])] = 1

    missing_edges = set([])
    while (len(missing_edges)<9437519):
        a=random.randint(1, 1862220)
        b=random.randint(1, 1862220)
        tmp = edges.get((a,b),-1)
        if tmp == -1 and a!=b:
            try:
                if nx.shortest_path_length(g,source=a,target=b) > 2:
                    missing_edges.add((a,b))
            else:
                continue
        except:
            missing_edges.add((a,b))
        else:
            continue
    pickle.dump(missing_edges,open('missing_edges_final.p','wb'))
else:
    missing_edges = pickle.load(open('missing_edges_final.p','rb'))
```

Wall time: 3.18 s

In [30]:

```
len(missing_edges)
```

Out[30]:

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 [31]:

```
from sklearn.model_selection import train_test_split
missing_edges = pickle.load(open('missing_edges_final.p','rb'))
df_pos = pd.read_csv('train.csv')
df_neg = pd.DataFrame(list(missing_edges), columns=['source_node', 'destination_node'])

print("Number of nodes in the graph with edges", df_pos.shape[0])
print("Number of nodes in the graph without edges", df_neg.shape[0])

#Train test split
#Spiltted data into 80-20
#positive links and negative links seperatly because we need positive training data onl
y for creating graph and for feature generation
X_train_pos, X_test_pos, y_train_pos, y_test_pos = train_test_split(df_pos,np.ones(len
(df_pos)),test_size=0.2, random_state=9)
X_train_neg, X_test_neg, y_train_neg, y_test_neg = train_test_split(df_neg,np.zeros(le
n(df_neg)),test_size=0.2, random_state=9)

print('='*60)
print("Number of nodes in the train data graph with edges", X_train_pos.shape[0],"=",y_
train_pos.shape[0])
print("Number of nodes in the train data graph without edges", X_train_neg.shape[0],"="
, y_train_neg.shape[0])
print('='*60)
print("Number of nodes in the test data graph with edges", X_test_pos.shape[0],"=",y_te
st_pos.shape[0])
print("Number of nodes in the test data graph without edges", X_test_neg.shape[0],"=",y
_test_neg.shape[0])

#removing header and saving
X_train_pos.to_csv('train_pos_after_eda.csv',header=False, index=False)
X_test_pos.to_csv('test_pos_after_eda.csv',header=False, index=False)
X_train_neg.to_csv('train_neg_after_eda.csv',header=False, index=False)
X_test_neg.to_csv('test_neg_after_eda.csv',header=False, index=False)
```

Number of nodes in the graph with edges 9437519

Number of nodes in the graph without edges 9437519

=====

Number of nodes in the train data graph with edges 7550015 = 7550015

Number of nodes in the train data graph without edges 7550015 = 7550015

=====

Number of nodes in the test data graph with edges 1887504 = 1887504

Number of nodes in the test data graph without edges 1887504 = 1887504

In [32]:

```

if (os.path.isfile('train_pos_after_eda.csv')) and (os.path.isfile('test_pos_after_eda.csv')):
    train_graph=nx.read_edgelist('train_pos_after_eda.csv',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
    test_graph=nx.read_edgelist('test_pos_after_eda.csv',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
    print(nx.info(train_graph))
    print(nx.info(test_graph))

    # finding the unique nodes in the both train and test graphs
    train_nodes_pos = set(train_graph.nodes())
    test_nodes_pos = set(test_graph.nodes())

    trY_teY = len(train_nodes_pos.intersection(test_nodes_pos))
    trY_teN = len(train_nodes_pos - test_nodes_pos)
    teY_trN = len(test_nodes_pos - train_nodes_pos)

    print('no of people common in train and test -- ',trY_teY)
    print('no of people present in train but not present in test -- ',trY_teN)

    print('no of people present in test but not present in train -- ',teY_trN)
    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: 1780722

Number of edges: 7550015

Average in degree: 4.2399

Average out degree: 4.2399

Name:

Type: DiGraph

Number of nodes: 1144623

Number of edges: 1887504

Average in degree: 1.6490

Average out degree: 1.6490

no of people common in train and test -- 1063125

no of people present in train but not present in test -- 717597

no of people present in test but not present in train -- 81498

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

7.1200735962845405 %

we have a cold start problem here

In [33]:

```
#final train and test data sets
X_train_pos = pd.read_csv('train_pos_after_eda.csv', names=['source_node', 'destination_node'])
X_test_pos = pd.read_csv('test_pos_after_eda.csv', names=['source_node', 'destination_node'])
X_train_neg = pd.read_csv('train_neg_after_eda.csv', names=['source_node', 'destination_node'])
X_test_neg = pd.read_csv('test_neg_after_eda.csv', names=['source_node', 'destination_node'])

print('='*60)
print("Number of nodes in the train data graph with edges", X_train_pos.shape[0])
print("Number of nodes in the train data graph without edges", X_train_neg.shape[0])
print('='*60)
print("Number of nodes in the test data graph with edges", X_test_pos.shape[0])
print("Number of nodes in the test data graph without edges", X_test_neg.shape[0])

X_train = X_train_pos.append(X_train_neg, ignore_index=True)
y_train = np.concatenate((y_train_pos, y_train_neg))
X_test = X_test_pos.append(X_test_neg, ignore_index=True)
y_test = np.concatenate((y_test_pos, y_test_neg))

X_train.to_csv('train_after_eda.csv', header=False, index=False)
X_test.to_csv('test_after_eda.csv', header=False, index=False)
pd.DataFrame(y_train.astype(int)).to_csv('train_y.csv', header=False, index=False)
pd.DataFrame(y_test.astype(int)).to_csv('test_y.csv', header=False, index=False)
```

```
=====
Number of nodes in the train data graph with edges 7550015
Number of nodes in the train data graph without edges 7550015
=====
Number of nodes in the test data graph with edges 1887504
Number of nodes in the test data graph without edges 1887504
```

In [34]:

```
print("Data points in train data", X_train.shape)
print("Data points in test data", X_test.shape)
print("Shape of target variable in train", y_train.shape)
print("Shape of target variable in test", y_test.shape)
```

```
Data points in train data (15100030, 2)
Data points in test data (3775008, 2)
Shape of target variable in train (15100030,)
Shape of target variable in test (3775008,)
```

2.1 Jaccard Distance:

In [35]:

```
#for followees
def jaccard_for_followees(a,b):
    try:
        if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.successors(b))) == 0:
            return 0
        sim = (len(set(train_graph.successors(a)).intersection(set(train_graph.successors(b)))))/\
            (len(set(train_graph.successors(a)).union(set(train_graph.successors(b)))))
    except:
        return 0
    return sim
```

In [36]:

```
#one test case
print(jaccard_for_followees(273084,1505602))
```

0.0

In [37]:

```
#node 1635354 not in graph
print(jaccard_for_followees(273084,1505602))
```

0.0

In [38]:

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

In [39]:

```
print(jaccard_for_followers(273084,470294))
#node 1635354 not in graph
print(jaccard_for_followers(669354,1635354))
```

0.0

0

2.2 Cosine distance

In [40]:

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

In [41]:

```
print(cosine_for_followees(273084,1505602))
```

0.0

In [42]:

```
print(cosine_for_followees(273084,1635354))
```

0

In [43]:

```
def cosine_for_followers(a,b):
    try:
        if len(set(train_graph.predecessors(a))) == 0 | len(set(train_graph.predecessors(b))) == 0:
            return 0
        sim = (len(set(train_graph.predecessors(a)).intersection(set(train_graph.predecessors(b))))) /\
            (math.sqrt(len(set(train_graph.predecessors(a)))*len(set(train_graph.predecessors(b)))))
        return sim
    except:
        return 0
```

In [44]:

```
print(cosine_for_followers(2,470294))
print(cosine_for_followers(669354,1635354))
```

0.02886751345948129

0

3. Ranking Measures

3.1 Page Ranking

In [45]:

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

In [46]:

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

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

In [47]:

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

```
5.615699699389075e-07
```

4. Other Graph Features

In [48]:

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

In [49]:

```
#testing
compute_shortest_path_length(77697, 826021)
```

Out[49]:

```
10
```

In [50]:

```
#testing
compute_shortest_path_length(669354,1635354)
```

Out[50]:

-1

In [51]:

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

In [52]:

```
belongs_to_same_wcc(861, 1659750)
```

Out[52]:

0

In [53]:

```
belongs_to_same_wcc(669354,1635354)
```

Out[53]:

0

4.3 Adamic/Adar Index:

In [54]:

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

In [55]:

```
calc_adar_in(1,189226)
calc_adar_in(669354,1635354)
```

Out[55]:

0

4.4 Is person was following back:

In [56]:

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

In [57]:

```
follows_back(1,189226)
```

Out[57]:

1

In [58]:

```
follows_back(669354,1635354)
```

Out[58]:

0

4.5 Katz Centrality:

In [59]:

```
if not os.path.isfile('katz.p'):
    katz = nx.katz.katz_centrality(train_graph,alpha=0.005,beta=1)
    pickle.dump(katz,open('katz.p','wb'))
else:
    katz = pickle.load(open('katz.p','rb'))
```

In [60]:

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

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

In [61]:

```
mean_katz = float(sum(katz.values())) / len(katz)
print(mean_katz)
```

```
0.0007483800935562018
```

Hits Score

In [62]:

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

In [63]:

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

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

Calculating Preferential Attachment

In [62]:

```
#Preferential Attachment
def calc_pref_att(a,b):
    try:
        return len(set(train_graph.predecessors(a))) * len(set(train_graph.predecessors
(b)))
    except:
        return 0
```


In [63]:

```
#testing
calc_pref_att(1,189226)
```

Out[63]:

9

SVD Dot Features

In [74]:

```
#svd_dot_u
def svd_dot_u(node):
    try:
        s_node = node[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5',
                        'svd_u_s_6']]
        d_node = node[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5',
                        'svd_u_d_6']]
        return np.dot(s_node,d_node)
    except:
        return 0
```

In [75]:

```
#svd_dot_v
def svd_dot_v(node):
    try:
        s_node = node[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5',
                        'svd_v_s_6'],]
        d_node = node[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5',
                        'svd_v_d_6']]
        return np.dot(s_node,d_node)
    except:
        return 0
```

In [76]:

```
svd_dot_v(df_final_train.iloc[1])
```

Out[76]:

0

5. Featurization

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

In [68]:

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

In [69]:

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

In [70]:

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

Number of rows in the train data file: 15100028

Number of rows we are going to elimiate in train data are 15000028

Number of rows in the test data file: 3775006

Number of rows we are going to elimiate in test data are 3725006

In [71]:

```
df_final_train = pd.read_csv('train_after_eda.csv', skiprows=skip_train, names=['source
_node', 'destination_node'])
df_final_train['indicator_link'] = pd.read_csv('train_y.csv', skiprows=skip_train, name
s=['indicator_link'])
print("Our train matrix size ",df_final_train.shape)
df_final_train.head(2)
```

Our train matrix size (100002, 3)

Out[71]:

	source_node	destination_node	indicator_link
0	273084	1505602	1
1	1859230	521884	1

In [72]:

```
df_final_test = pd.read_csv('test_after_eda.csv', skiprows=skip_test, names=['source_node', 'destination_node'])
df_final_test['indicator_link'] = pd.read_csv('test_y.csv', skiprows=skip_test, names=['indicator_link'])
print("Our test matrix size ", df_final_test.shape)
df_final_test.head(2)
```

Our test matrix size (50002, 3)

Out[72]:

	source_node	destination_node	indicator_link
0	848424	784690	1
1	1341156	1679887	1

5.2 Adding a set of features

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

1.jaccard_followers 2.jaccard_followees 3.cosine_followers 4.cosine_followees 5.num_followers_s
6.num_followees_s 7.num_followers_d 8.num_followees_d 9.inter_followers 10.inter_followees

In [77]:

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

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

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

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

In [78]:

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

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

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

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

In [79]:

```
from pandas import HDFStore,DataFrame
from pandas import read_hdf
```

In [80]:

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

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

    hdf = HDFStore('storage_sample_stage1.h5')
    hdf.put('train_df', df_final_train, format='table', data_columns=True)
    hdf.put('test_df', df_final_test, format='table', data_columns=True)
    hdf.close()
else:
    df_final_train = read_hdf('storage_sample_stage1.h5', 'train_df', mode='r')
    df_final_test = 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 [81]:

```
if not os.path.isfile('storage_sample_stage2.h5'):
    #mapping adar index on train
    df_final_train['adar_index'] = df_final_train.apply(lambda row: calc_adar_in(row['source_node'],row['destination_node']),axis=1)
    #mapping adar index on test
    df_final_test['adar_index'] = df_final_test.apply(lambda row: calc_adar_in(row['source_node'],row['destination_node']),axis=1)

    #-----
    #mapping followback or not on train
    df_final_train['follows_back'] = df_final_train.apply(lambda row: follows_back(row['source_node'],row['destination_node']),axis=1)

    #mapping followback or not on test
    df_final_test['follows_back'] = df_final_test.apply(lambda row: follows_back(row['source_node'],row['destination_node']),axis=1)

    #-----
    #mapping same component of wcc or not on train
    df_final_train['same_comp'] = df_final_train.apply(lambda row: belongs_to_same_wcc(row['source_node'],row['destination_node']),axis=1)

    ##mapping same component of wcc or not on train
    df_final_test['same_comp'] = df_final_test.apply(lambda row: belongs_to_same_wcc(row['source_node'],row['destination_node']),axis=1)

    #-----
    #mapping shortest path on train
    df_final_train['shortest_path'] = df_final_train.apply(lambda row: compute_shortest_path_length(row['source_node'],row['destination_node']),axis=1)
    #mapping shortest path on test
    df_final_test['shortest_path'] = df_final_test.apply(lambda row: compute_shortest_path_length(row['source_node'],row['destination_node']),axis=1)

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

5.4 Adding new set of features

In [84]:

```

if not os.path.isfile('storage_sample_stage3.h5'):

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

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

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

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

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

    df_final_test['hubs_s'] = df_final_test.source_node.apply(lambda x: hits[0].get(x,0
))
    df_final_test['hubs_d'] = df_final_test.destination_node.apply(lambda x: hits[0].ge
t(x,0))
    #=====

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

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

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

```

```
df_final_train = read_hdf('storage_sample_stage3.h5', 'train_df', mode='r')
df_final_test = read_hdf('storage_sample_stage3.h5', 'test_df', mode='r')
```

5.5 Adding new set of features

In [85]:

```
#SVD features for both source and destination
def svd(x, S):
    try:
        z = sadj_dict[x]
        return S[z]
    except:
        return [0,0,0,0,0,0]
```

In [86]:

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

In [87]:

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

In [88]:

```
from scipy.sparse.linalg import svds, eigs
import gc
U, s, V = svds(Adj, k = 6)
print('Adjacency matrix Shape', Adj.shape)
print('U Shape', U.shape)
print('V Shape', V.shape)
print('s Shape', s.shape)
```

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

In [89]:

```

if not os.path.isfile('storage_sample_stage4.h5'):
    #=====
    =====

    df_final_train[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5', 'sv
d_u_s_6']] = \
        df_final_train.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)

    df_final_train[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5', 'sv
d_u_d_6']] = \
        df_final_train.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)
    #=====
    =====

    df_final_train[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'sv
d_v_s_6',]] = \
        df_final_train.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)

    df_final_train[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'sv
d_v_d_6']] = \
        df_final_train.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
    #=====
    =====

    df_final_test[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5', 'svd
_u_s_6']] = \
        df_final_test.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)

    df_final_test[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5', 'svd
_u_d_6']] = \
        df_final_test.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)

    #=====
    =====

    df_final_test[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd
_v_s_6',]] = \
        df_final_test.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)

    df_final_test[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd
_v_d_6']] = \
        df_final_test.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
    #=====
    =====

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

```

In [90]:

```

#reading
from pandas import read_hdf
df_final_train = read_hdf('storage_sample_stage4.h5', 'train_df', mode='r')
df_final_test = read_hdf('storage_sample_stage4.h5', 'test_df', mode='r')

```

In [91]:

```
df_final_train.columns
```

Out[91]:

```
Index(['source_node', 'destination_node', 'indicator_link',
      'jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_inde
x',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_
out',
      'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
      'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities
_s',
      'authorities_d', 'svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_
4',
      'svd_u_s_5', 'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
      'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2',
      'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1',
      'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6'],
      dtype='object')
```

In [92]:

```
%%time
s_node = df_final_train.loc[1][['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd
_v_s_5', 'svd_v_s_6',]]
```

Wall time: 68.8 ms

In [93]:

```
d_node = df_final_train.iloc[182][['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4',
'svd_v_s_5', 'svd_v_s_6',]]
```

In [94]:

```
type(s_node)
```

Out[94]:

pandas.core.series.Series

In [95]:

```
s_node[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6',]]
```

Out[95]:

```
svd_v_s_1    -9.996461e-10
svd_v_s_2     6.107418e-10
svd_v_s_3     2.482648e-09
svd_v_s_4     1.757569e-11
svd_v_s_5     1.154567e-09
svd_v_s_6     1.519087e-13
Name: 1, dtype: float64
```

In [97]:

```
d_node[['svd_v_s_1','svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6'],]
```

Out[97]:

```
svd_v_s_1    -6.066617e-13
svd_v_s_2     5.092009e-13
svd_v_s_3     5.362832e-12
svd_v_s_4     8.472017e-12
svd_v_s_5     1.397855e-12
svd_v_s_6     2.041652e-14
Name: 182, dtype: float64
```

In [98]:

```
%%time
sum_x = 0.0
for i in range(6):
    sum_x += s_node[i]*d_node[i]
```

Wall time: 0 ns

In [99]:

```
print(sum_x)
```

1.599428486610478e-20

In [100]:

```
%%time
np.dot(np.array(s_node),np.array(d_node))
```

Wall time: 0 ns

Out[100]:

1.5994284866104777e-20

In [101]:

```
%%time
np.dot(s_node,d_node)
```

Wall time: 0 ns

Out[101]:

1.5994284866104777e-20

In [102]:

```
df_final_test.iloc[1][['source_node','destination_node']]
```

Out[102]:

```
source_node      15078.0
destination_node  370241.0
Name: 1, dtype: float64
```

In [103]:

```
%%time
df_final_train['svd_dot_u'] = df_final_train.apply(lambda row:svd_dot_u(row),axis=1)
df_final_train['svd_dot_v'] = df_final_train.apply(lambda row:svd_dot_v(row),axis=1)

df_final_test['svd_dot_u'] = df_final_test.apply(lambda row:svd_dot_u(row),axis=1)
df_final_test['svd_dot_v'] = df_final_test.apply(lambda row:svd_dot_v(row),axis=1)
```

Wall time: 4min 2s

In [104]:

```
df_final_train['pref_att'] = df_final_train.apply(lambda row:
                                                    calc_pref_att(row['source_node'],row['desti
nation_node']),axis=1)
df_final_test['pref_att'] = df_final_test.apply(lambda row:
                                                    calc_pref_att(row['source_node'],row['desti
nation_node']),axis=1)
```

In [105]:

```
df_final_train.shape
```

Out[105]:

(100002, 57)

In [106]:

```
df_final_test.shape
```

Out[106]:

(50002, 57)

In [107]:

```
df_final_train.iloc[1]['svd_dot_u']
```

Out[107]:

1.678875635579273e-17

In [108]:

```
svd_dot_u(df_final_train.iloc[1])
```

Out[108]:

1.678875635579273e-17

Social network Graph Link Prediction - Facebook Challenge

In [115]:

```
from sklearn.ensemble import RandomForestClassifier
```

In [117]:

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

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

import warnings
import networkx as nx
import pdb
import pickle
from pandas import HDFStore, DataFrame
from pandas import read_hdf
from scipy.sparse.linalg import svds, eigs
import gc
from tqdm import tqdm
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
```

In [119]:

```
#reading
from pandas import read_hdf
df_final_train = read_hdf('storage_sample_stage4.h5', 'train_df', mode='r')
df_final_test = read_hdf('storage_sample_stage4.h5', 'test_df', mode='r')
```

In [121]:

```
type(df_final_train)
```

Out[121]:

```
pandas.core.frame.DataFrame
```

In [122]:

```
df_final_train.columns
```

Out[122]:

```
Index(['source_node', 'destination_node', 'indicator_link',
      'jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_inde
x',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_
out',
      'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
      'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities
_s',
      'authorities_d', 'svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_
4',
      'svd_u_s_5', 'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
      'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2',
      'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1',
      'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6'],
      dtype='object')
```

In [123]:

```
y_train = df_final_train.indicator_link
y_test = df_final_test.indicator_link
```

In [124]:

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


In [125]:

```

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

```

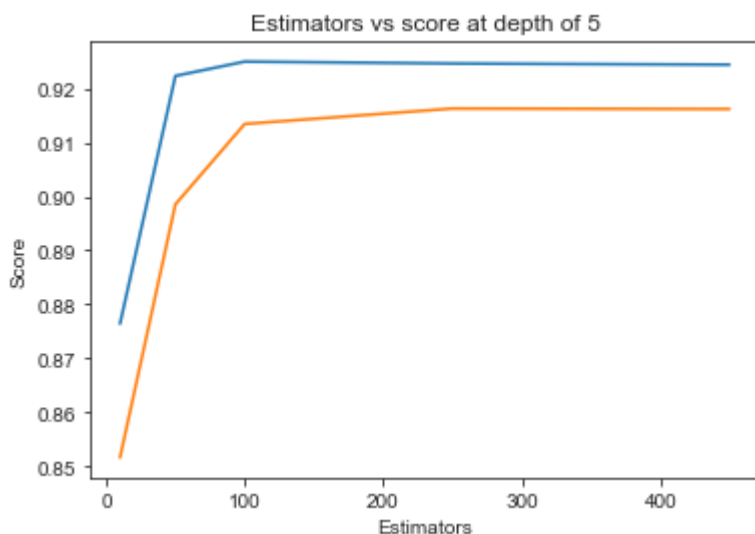
```

Estimators = 10 Train Score 0.8763989111178326 test Score 0.8515874423554
451
Estimators = 50 Train Score 0.9223470872251155 test Score 0.8985419722198
258
Estimators = 100 Train Score 0.9250028905683384 test Score 0.913442106830
607
Estimators = 250 Train Score 0.9246523388116308 test Score 0.916303492846
3616
Estimators = 450 Train Score 0.9244286375689212 test Score 0.916212027936
4

```

Out[125]:

Text(0.5, 1.0, 'Estimators vs score at depth of 5')



In [127]:

```

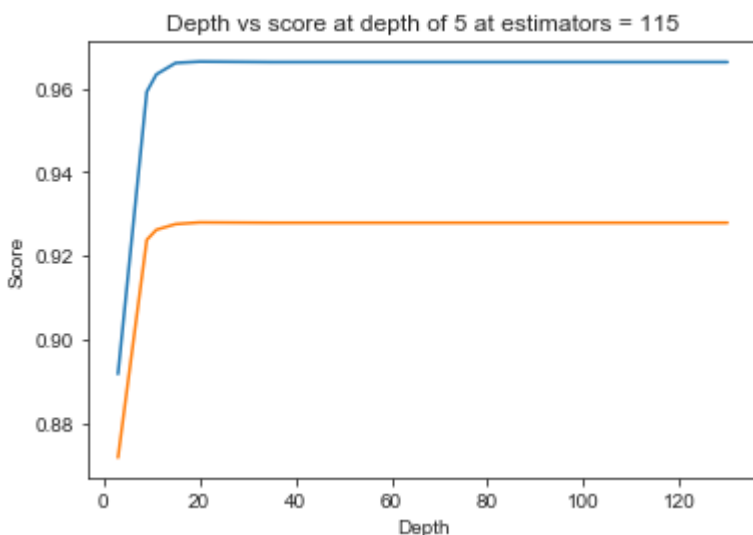
depths = [3,9,11,15,20,35,50,70,130]
train_scores = []
test_scores = []
for i in depths:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=i, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0,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

```

```

depth = 3 Train Score 0.8916639914392723 test Score 0.8716660477511523
depth = 9 Train Score 0.9591895446325827 test Score 0.9237609019357582
depth = 11 Train Score 0.9633332311120653 test Score 0.9261716514403203
depth = 15 Train Score 0.9660652614635439 test Score 0.9275227416966364
depth = 20 Train Score 0.9663885238051573 test Score 0.9279292619465658
depth = 35 Train Score 0.9662790816170552 test Score 0.9278429049049897
depth = 50 Train Score 0.9662790816170552 test Score 0.9278429049049897
depth = 70 Train Score 0.9662790816170552 test Score 0.9278429049049897
depth = 130 Train Score 0.9662790816170552 test Score 0.9278429049049897

```



In [128]:

```

from sklearn.metrics import f1_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint
from scipy.stats import uniform

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

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

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

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

```

```

mean test scores [0.96468376 0.96393932 0.96176657 0.96410511 0.96615727]
mean train scores [0.96527135 0.96456104 0.96208427 0.9647248  0.96716544]

```

In [129]:

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

```

RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini
i',
                        max_depth=14, max_features='auto', max_leaf_nodes=N
one,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=28, min_samples_split=111,
                        min_weight_fraction_leaf=0.0, n_estimators=121,
                        n_jobs=-1, oob_score=False, random_state=25, verbos
e=0,
                        warm_start=False)

```

In [130]:

```

clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                             max_depth=14, max_features='auto', max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=28, min_samples_split=111,
                             min_weight_fraction_leaf=0.0, n_estimators=121, n_jobs=-1,
                             oob_score=False, random_state=25, verbose=0, warm_start=False)

```

In [131]:

```

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

```

In [132]:

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

Train f1 score 0.9676348040464134

Test f1 score 0.9286348830322771

In [133]:

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

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

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

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

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

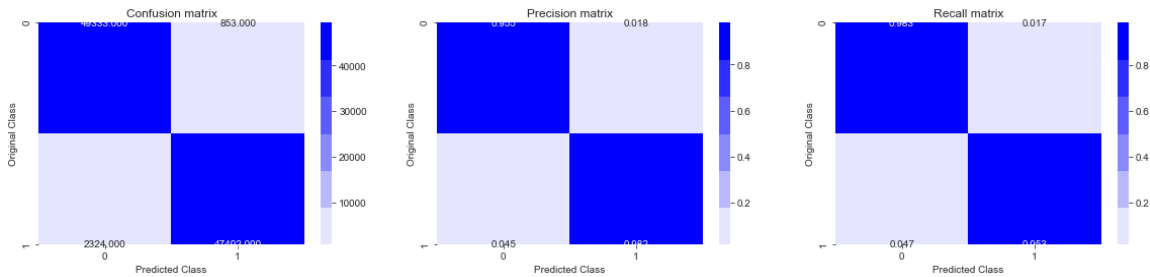
    plt.subplot(1, 3, 3)
    # representing B in heatmap format
    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
bels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Recall matrix")

    plt.show()
```

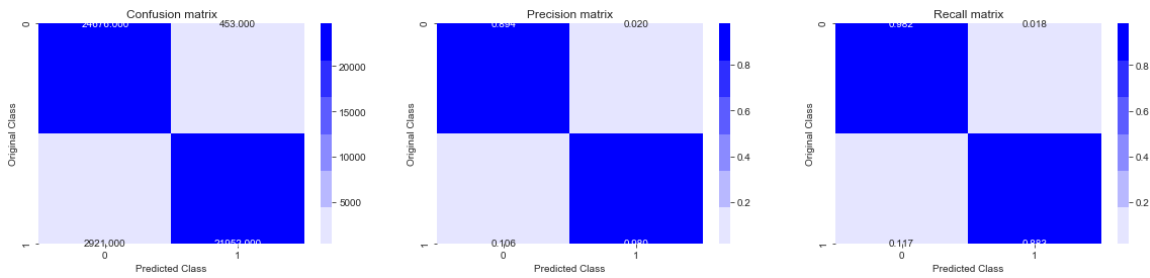
In [134]:

```
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

Train confusion_matrix

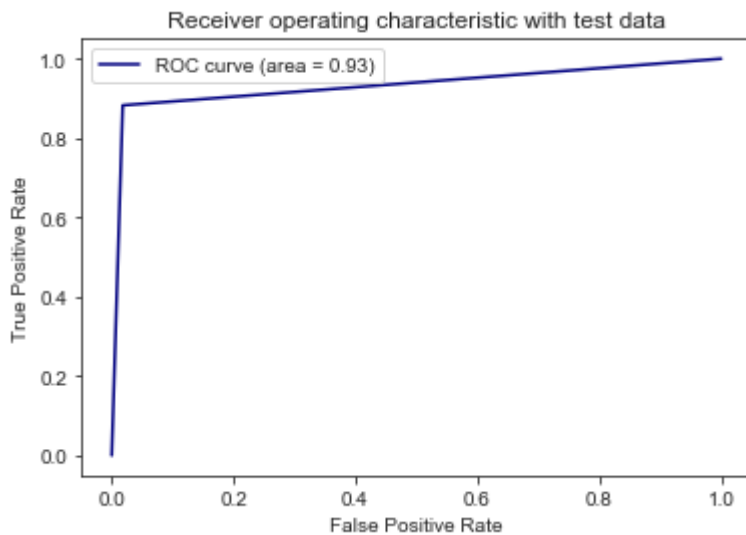


Test confusion_matrix



In [135]:

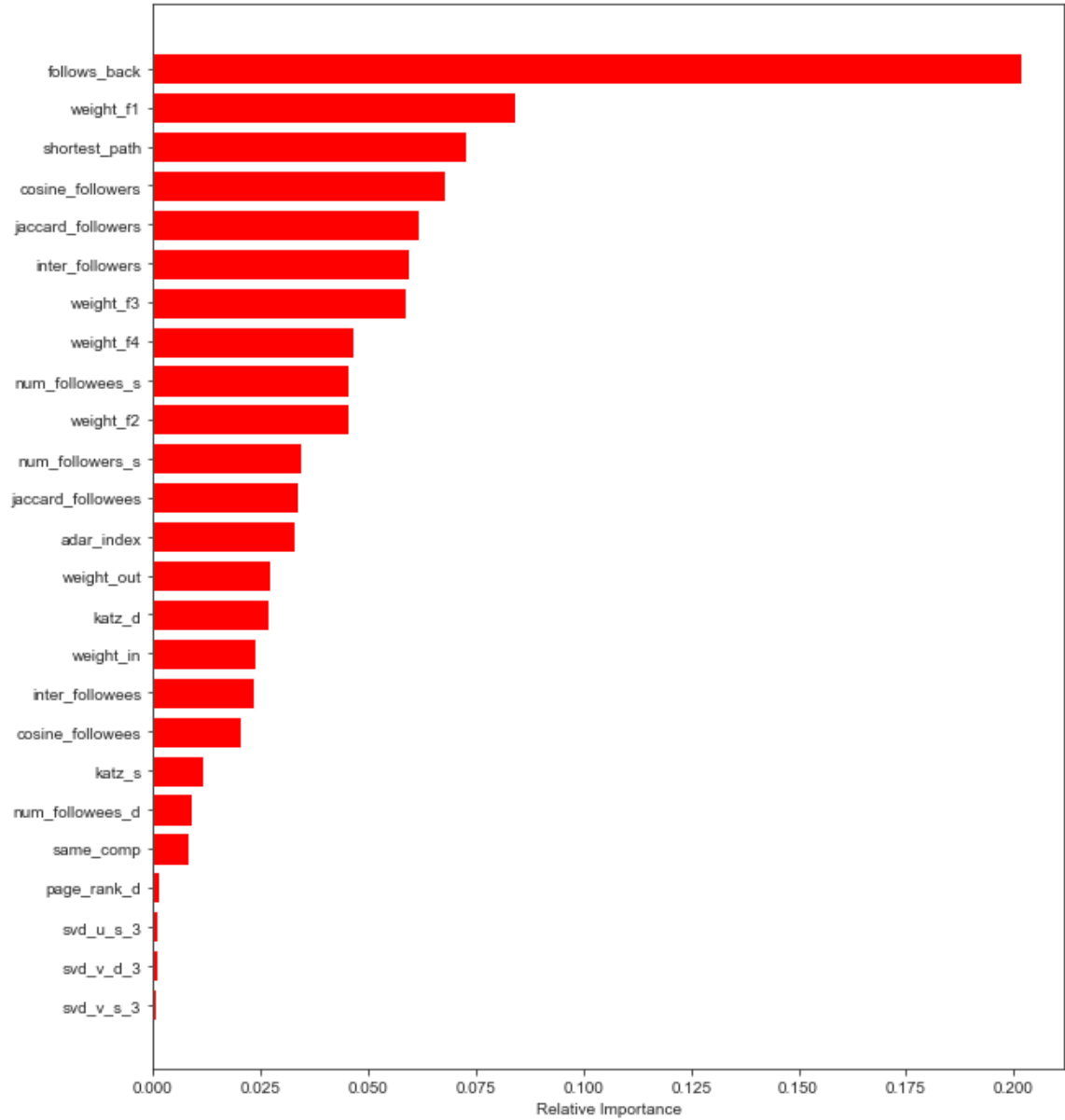
```
from sklearn.metrics import roc_curve, auc
fpr,tpr,ths = roc_curve(y_test,y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic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



In [136]:

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

Feature Importances



In [181]:

```
#reading
from pandas import read_hdf
df_final_train = read_hdf('storage_sample_stage5.h5', 'train_df', mode='r')
df_final_test = read_hdf('storage_sample_stage5.h5', 'test_df', mode='r')
```

In [182]:

```
df_final_train.columns
```

Out[182]:

```
Index(['source_node', 'destination_node', 'indicator_link',
      'jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_inde
x',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_
out',
      'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
      'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities
_s',
      'authorities_d', 'svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_
4',
      'svd_u_s_5', 'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
      'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2',
      'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1',
      'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6',
      'preferential_attachment_followers',
      'preferential_attachment_followees', 'svd_u_1_dot', 'svd_v_1_dot',
      'svd_u_2_dot', 'svd_v_2_dot', 'svd_u_3_dot', 'svd_v_3_dot',
      'svd_u_4_dot', 'svd_v_4_dot', 'svd_u_5_dot', 'svd_v_5_dot',
      'svd_u_6_dot', 'svd_v_6_dot', 'preferential_followers',
      'preferential_followees', 'svd_dot_1', 'svd_dot_2', 'svd_dot_3',
      'svd_dot_4', 'svd_dot_5', 'svd_dot_6'],
      dtype='object')
```

In [183]:

```
y_train = df_final_train.indicator_link
y_test = df_final_test.indicator_link
```

In [184]:

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

In [185]:

```
#taking depth and n_estimator as hyperparameter
max_depth = [1, 5, 10, 50, 100, 500]
n_estimators = [10, 100, 500]
```

In [137]:

```
df_final_train.columns
```

Out[137]:

```
Index(['jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_inde
x',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_
out',
      'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
      'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities
_s',
      'authorities_d', 'svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_
4',
      'svd_u_s_5', 'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
      'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2',
      'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1',
      'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6'],
      dtype='object')
```

In [138]:

```
df_final_test.columns
```

Out[138]:

```
Index(['jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_inde
x',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_
out',
      'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
      'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities
_s',
      'authorities_d', 'svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_
4',
      'svd_u_s_5', 'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
      'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2',
      'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1',
      'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6'],
      dtype='object')
```

In [139]:

```
from sklearn.metrics import f1_score
from sklearn.metrics import f1_score
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint as sp_randint
from scipy.stats import uniform
```

APPLYING XGBOOST

In [140]:

```
import xgboost as xgb
clf = xgb.XGBClassifier()
param_dist = {"n_estimators":sp_randint(105,125),
              "max_depth": sp_randint(2,10)
              }
model = RandomizedSearchCV(clf, param_distributions=param_dist,n_jobs=4,
                           n_iter=5,cv=3,scoring='f1',random_state=25,return_train_score = True)

model.fit(df_final_train,y_train)
print('mean test scores',model.cv_results_['mean_test_score'])
print('mean train scores',model.cv_results_['mean_train_score'])
```

```
mean test scores [0.97943378 0.98016113 0.97698755 0.97899135 0.97574467]
mean train scores [0.98340111 0.9865051  0.97815187 0.98236004 0.97616602]
```

In [141]:

model.cv_results_

Out[141]:

```
{'mean_fit_time': array([104.0736444 , 121.62355947,  72.88366477, 100.550
78999,
        51.50878231]),
  'std_fit_time': array([0.26882925, 2.7000768 , 0.38902393, 0.55008867, 3.
56157694]),
  'mean_score_time': array([0.51510811, 0.53787621, 0.3827463 , 0.46413136,
0.27016481]),
  'std_score_time': array([0.01124898, 0.01453507, 0.00294629, 0.02765494,
0.03342654]),
  'param_max_depth': masked_array(data=[6, 7, 4, 6, 3],
        mask=[False, False, False, False, False],
        fill_value='?',
        dtype=object),
  'param_n_estimators': masked_array(data=[120, 117, 113, 109, 110],
        mask=[False, False, False, False, False],
        fill_value='?',
        dtype=object),
  'params': [{'max_depth': 6, 'n_estimators': 120},
{'max_depth': 7, 'n_estimators': 117},
{'max_depth': 4, 'n_estimators': 113},
{'max_depth': 6, 'n_estimators': 109},
{'max_depth': 3, 'n_estimators': 110}],
  'split0_test_score': array([0.98024654, 0.98070085, 0.97830134, 0.9795694
7, 0.97696696]),
  'split1_test_score': array([0.97947767, 0.98004313, 0.97717507, 0.9793250
2, 0.97580473]),
  'split2_test_score': array([0.97857708, 0.97973938, 0.97548616, 0.9780795
2, 0.97446224]),
  'mean_test_score': array([0.97943378, 0.98016113, 0.97698755, 0.97899135,
0.97574467]),
  'std_test_score': array([0.00068226, 0.00040129, 0.00115692, 0.00065242,
0.00102343]),
  'rank_test_score': array([2, 1, 4, 3, 5]),
  'split0_train_score': array([0.98274106, 0.98565291, 0.97774795, 0.981746
83, 0.97569831]),
  'split1_train_score': array([0.98385871, 0.9868439 , 0.97830399, 0.982814
21, 0.97590104]),
  'split2_train_score': array([0.98360357, 0.9870185 , 0.97840367, 0.982519
08, 0.9768987 ]),
  'mean_train_score': array([0.98340111, 0.9865051 , 0.97815187, 0.9823600
4, 0.97616602]),
  'std_train_score': array([0.00047821, 0.0006068 , 0.0002885 , 0.00045003,
0.00052465])}
```

In [142]:

```

results = pd.DataFrame.from_dict(model.cv_results_)
results = results.sort_values(['param_max_depth', 'param_n_estimators'])

train_auc = results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']

results_score_sorted = results.sort_values(by=['mean_test_score'], ascending=False)
results_score_sorted.head()

```

Out[142]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max_depth	param_n
1	121.623559	2.700077	0.537876	0.014535	7	
0	104.073644	0.268829	0.515108	0.011249	6	
3	100.550790	0.550089	0.464131	0.027655	6	
2	72.883665	0.389024	0.382746	0.002946	4	
4	51.508782	3.561577	0.270165	0.033427	3	

In [143]:

```
print(model.best_estimator_)
```

```

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=1, gamma=0,
              learning_rate=0.1, max_delta_step=0, max_depth=7,
              min_child_weight=1, missing=None, n_estimators=117, n_jobs=
1,
              nthread=None, objective='binary:logistic', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=None, subsample=1, verbosity=1)

```

In [144]:

```

clf=xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
                      max_depth=7, min_child_weight=1, missing=None, n_estimators=117,
                      n_jobs=4, nthread=None, objective='binary:logistic', random_state=0,
                      reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                      silent=True, subsample=1)

```

In [145]:

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

In [146]:

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

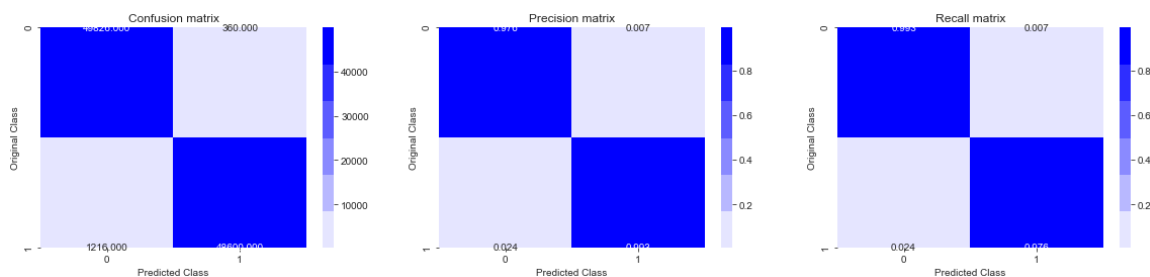
Train f1 score 0.9840447072163278

Test f1 score 0.9314194577352471

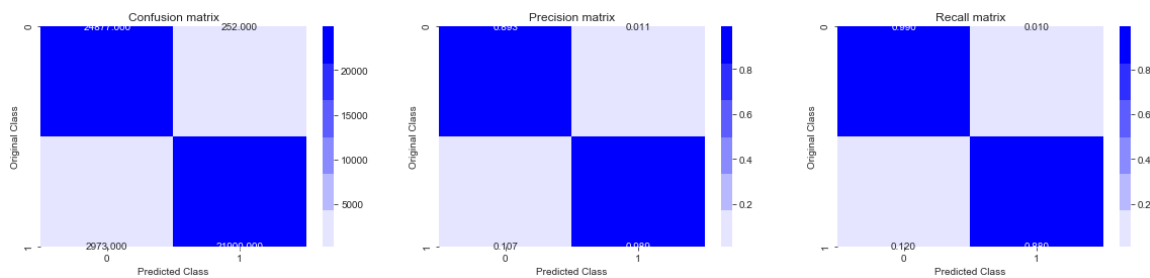
In [147]:

```
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

Train confusion_matrix

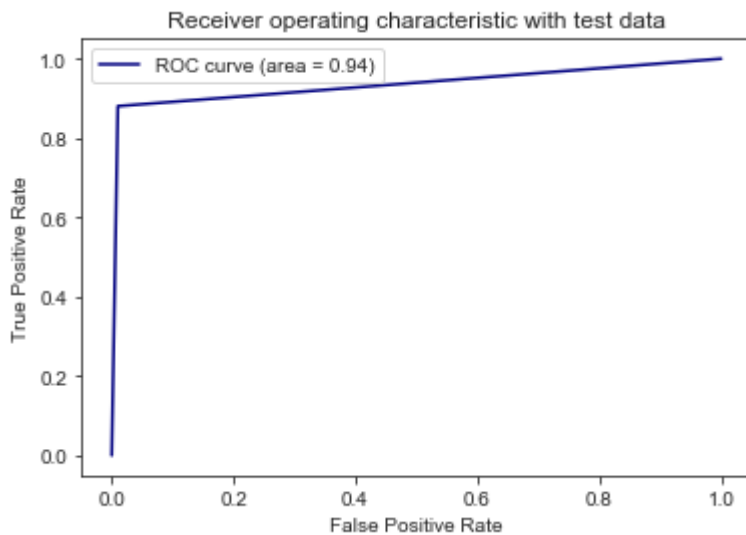


Test confusion_matrix



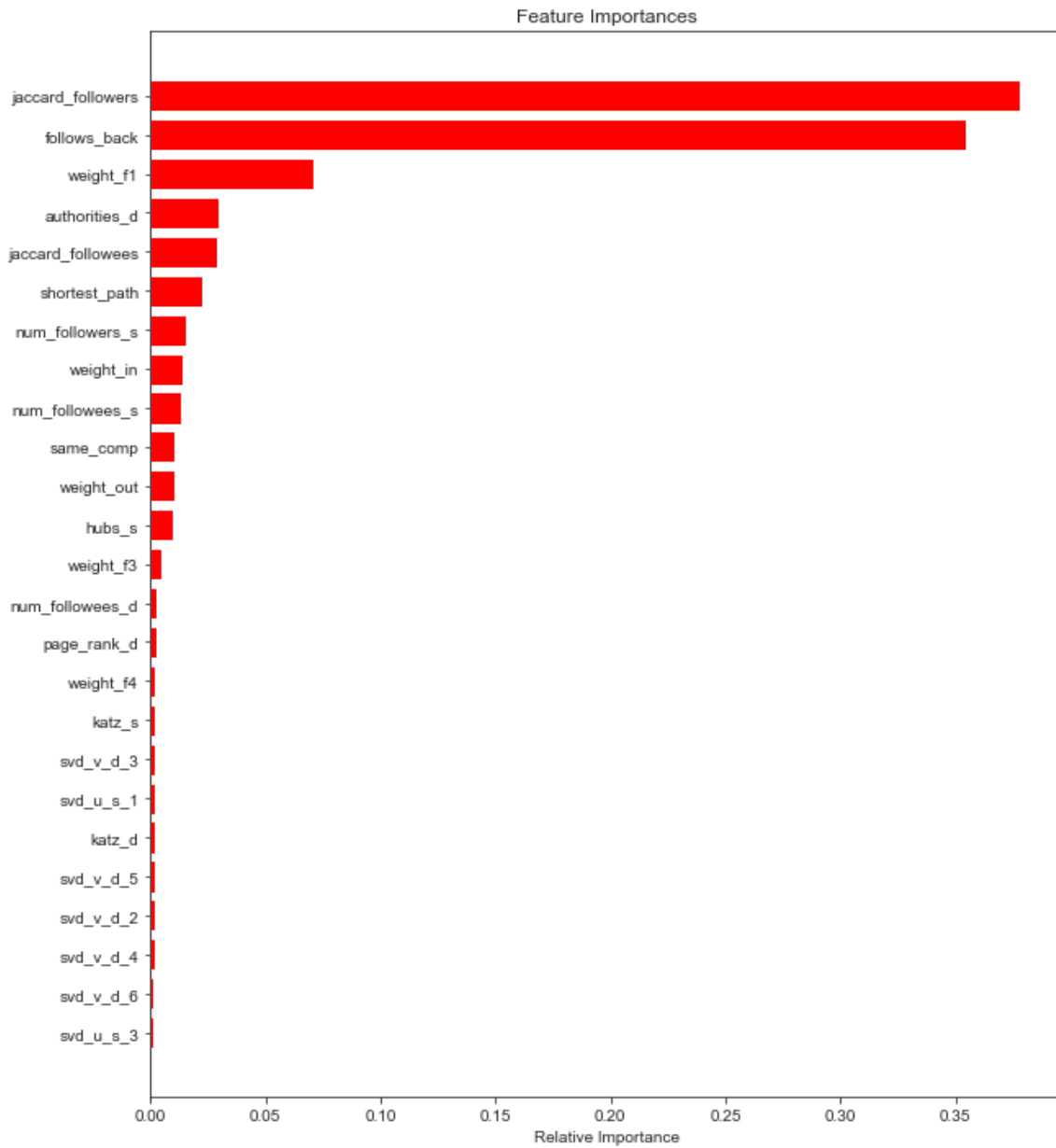
In [148]:

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



In [149]:

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

Observations: 1.XGBoost also performs very similar to Random Forest. 2.Two new added features Preferential attachment and svd_dot are also not very important as per XGBoost model , hence not much improvement in results .

In [150]:

```
# Please compare all models using Prettytable Library
# http://zetcode.com/python/prettytable/

from prettytable import PrettyTable
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                        max_depth=14, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=28, min_samples_split=111,
                        min_weight_fraction_leaf=0.0, n_estimators=121, n_jobs=-1,
                        oob_score=False, random_state=25, verbose=0, warm_start=False)

#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "F1-Score"]

x.add_row(["Previous Graph Based features", "Random Forest", "Max Depth:14 , Estimators : 111, min_samples_leaf:28, min_samples_split:111", 0.92])
x.add_row(["Previous Graph Based features + Two new features", "XGBoost", "Max Depth:7 , Estimators : 117", 0.93])

print(x)
```

- 1.This is a problem statement of Social network Graph Link Prediction. We have a given a dataset comprising of vertex and edges (excluding metadata). We have a graph which is directed.
- 2.We have created or posed this data as a classification task. For the mapping into supervised learning problem, we generated training samples of good and bad links from given directed graph and for each link.
- 3.We have to add new features as their is absence of metadata. Following are the list of engineered features, which are incorporated into the dataset:
- 4.Jaccard Distance, Preferential Attachment , Cosine distance (Otsuka-Ochiai coefficient) (both for followees and followers)
- 5.Ranking Measures, Shortest path, Adamic/Adar Index, person follow back, Katz Centrality, HITS Score, Weight Features
- 6.At last we have engineered another feature called svd_dot (Dot product between source node svd and destination node svd features)
- 7.Models used for machine learning here were Random Forest and Xgboost.
- 8.Performance can be compared from the above created table