

# 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://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]:

```

traincsv = pd.read_csv('data/train.csv')
traincsv

```

Out[2]:

	source_node	destination_node
0	1	690569
1	1	315892
2	1	189226
3	2	834328
4	2	1615927
...	...	...
9437514	1862219	1187308
9437515	1862219	563943
9437516	1862219	1044046
9437517	1862219	1022613
9437518	1862220	1748794

9437519 rows × 2 columns

In [3]:

```
#reading graph
if not os.path.isfile('data/after_eda/train_woheader.csv'):
    traincsv = pd.read_csv('data/train.csv')
    print(traincsv[traincsv.isna().any(1)])
    print(traincsv.info())
    print("Number of diplicate entries: ",sum(traincsv.duplicated()))
    traincsv.to_csv('data/after_eda/train_woheader.csv',header=False,index=False)
    print("saved the graph into file")
else:
    g=nx.read_edgelist('data/after_eda/train_woheader.csv',delimiter=',',create_using=nx.Di
    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 [4]:

```

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

subgraph=nx.read_edgelist('train_woheader_sample.csv',delimiter=',',create_using=nx.DiGraph)
# https://stackoverflow.com/questions/9402255/drawing-a-huge-graph-with-networkx-and-matplotlib

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

```

Name:

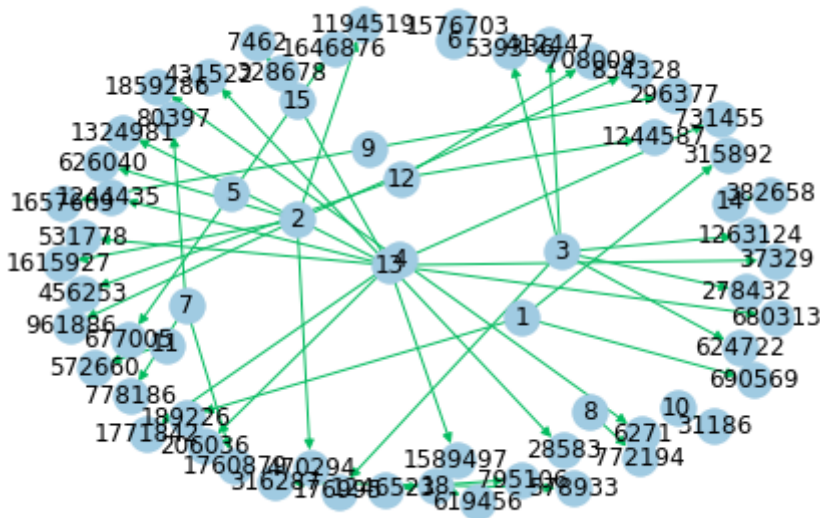
Type: DiGraph

Number of nodes: 66

Number of edges: 50

Average in degree: 0.7576

Average out degree: 0.7576



## 1. Exploratory Data Analysis

In [5]:

```

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

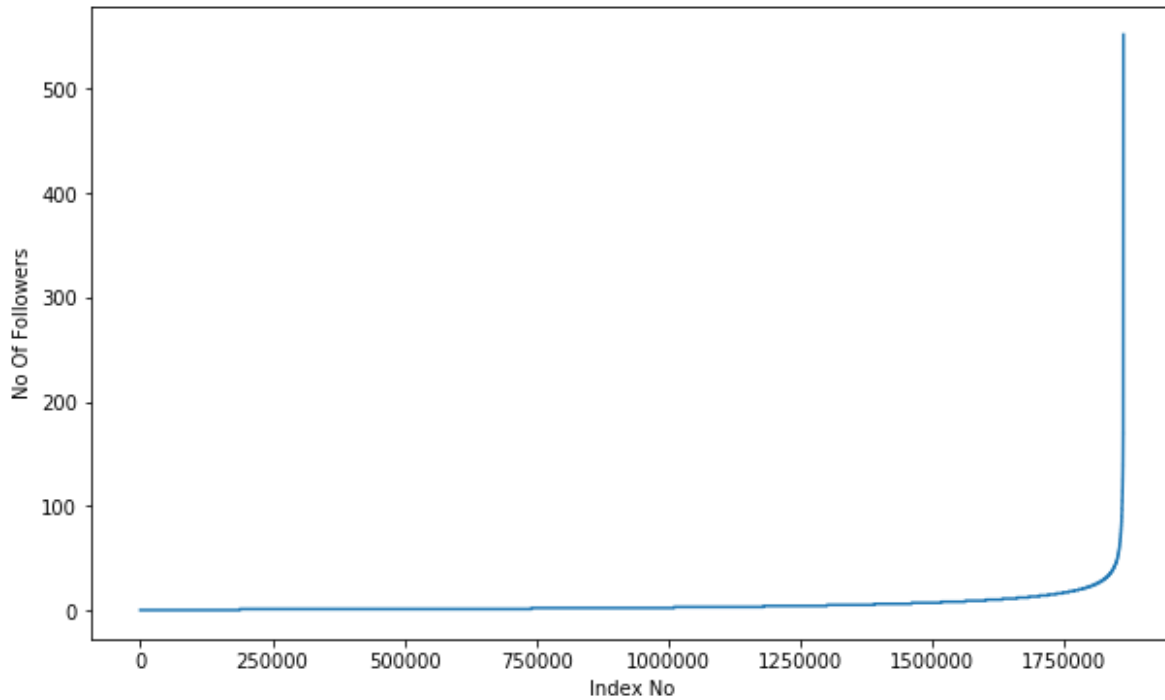
```

The number of unique persons 1862220

### 1.1 No of followers for each person

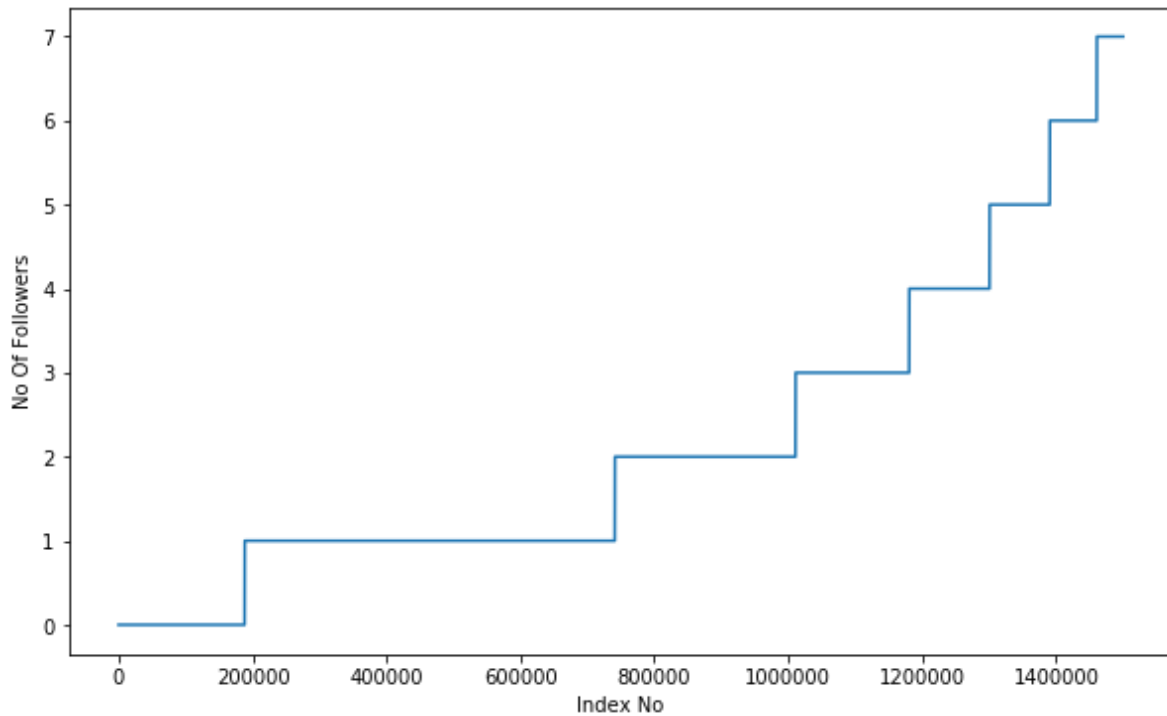
In [6]:

```
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()
```



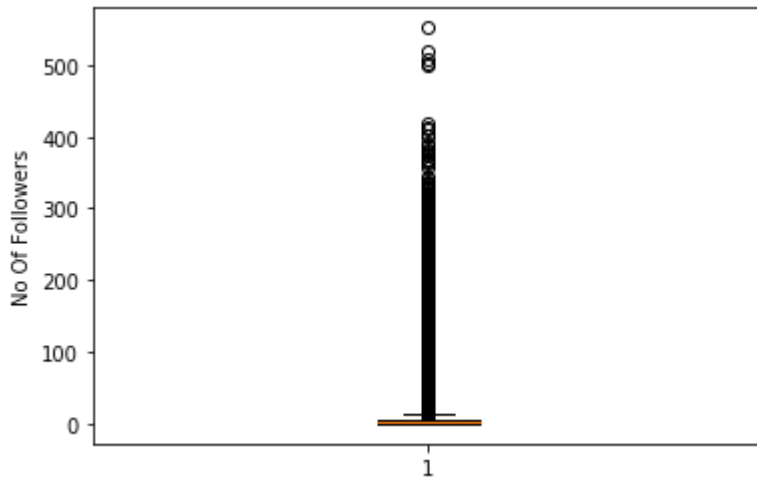
In [7]:

```
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 [8]:

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



In [9]:

```
### 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
```

1. Nearly the in-degree of 1800000 user is between (0-100).
2. Few users have more than 500 followers, they might be popular personalities.
3. 99% of data having followers of 40 only.

In [10]:

```

### 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

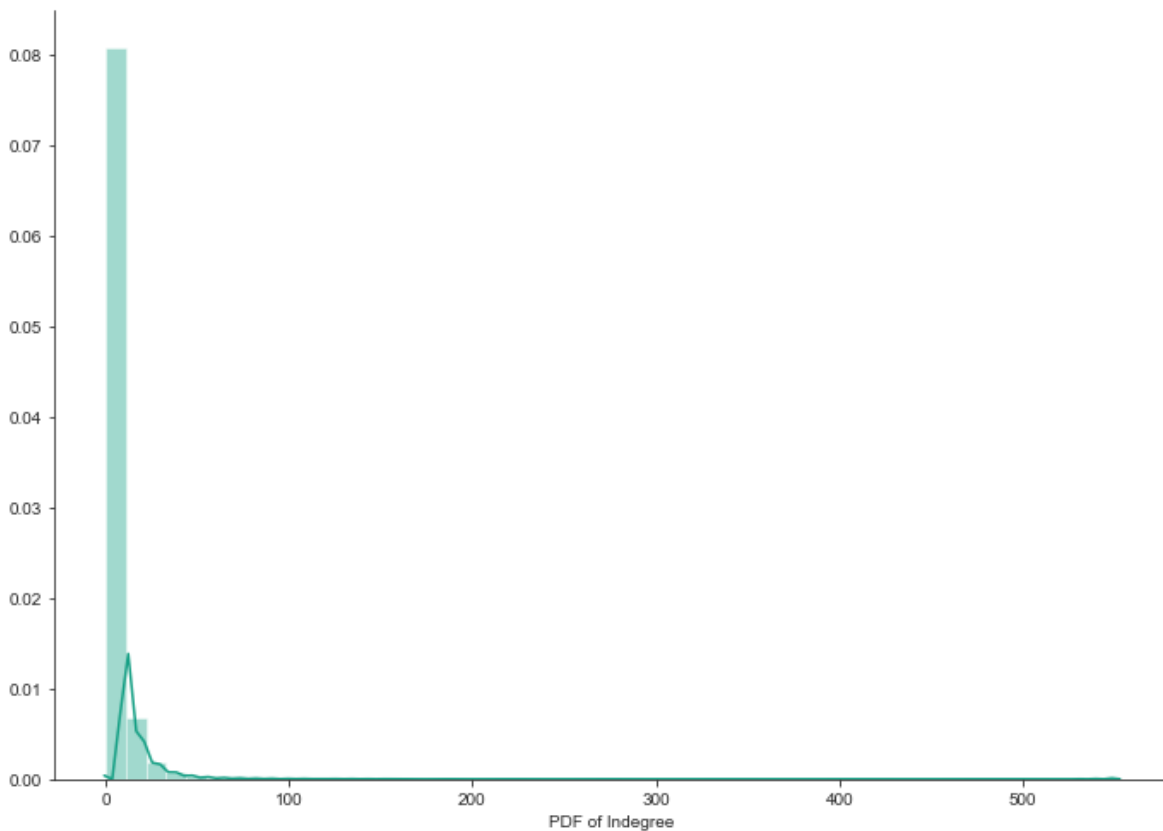
```

In [11]:

```

%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()

```

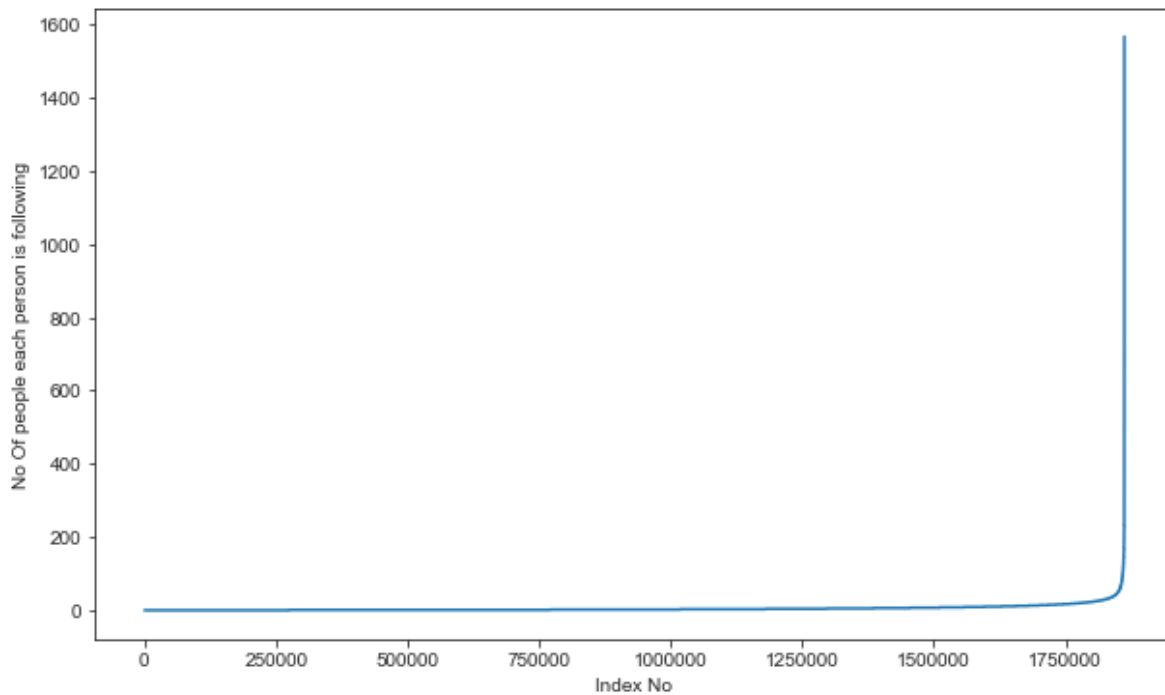


## 1.2 No of people each person is following



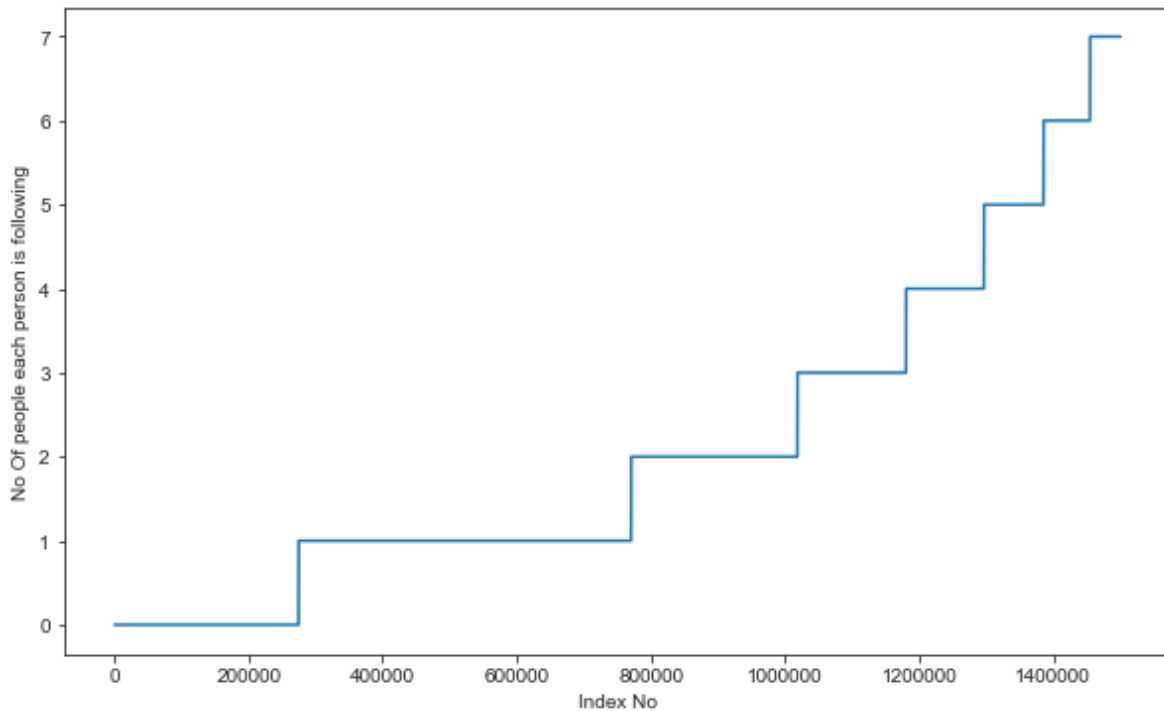
In [12]:

```
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()
```



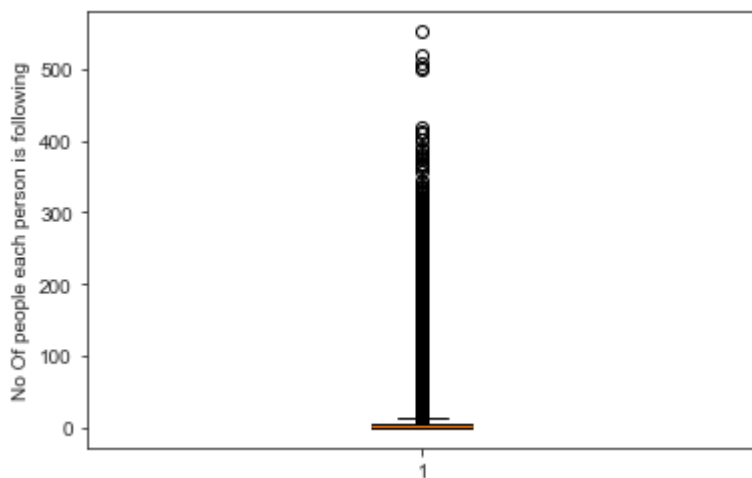
In [13]:

```
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 [14]:

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



In [15]:

```
### 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
```

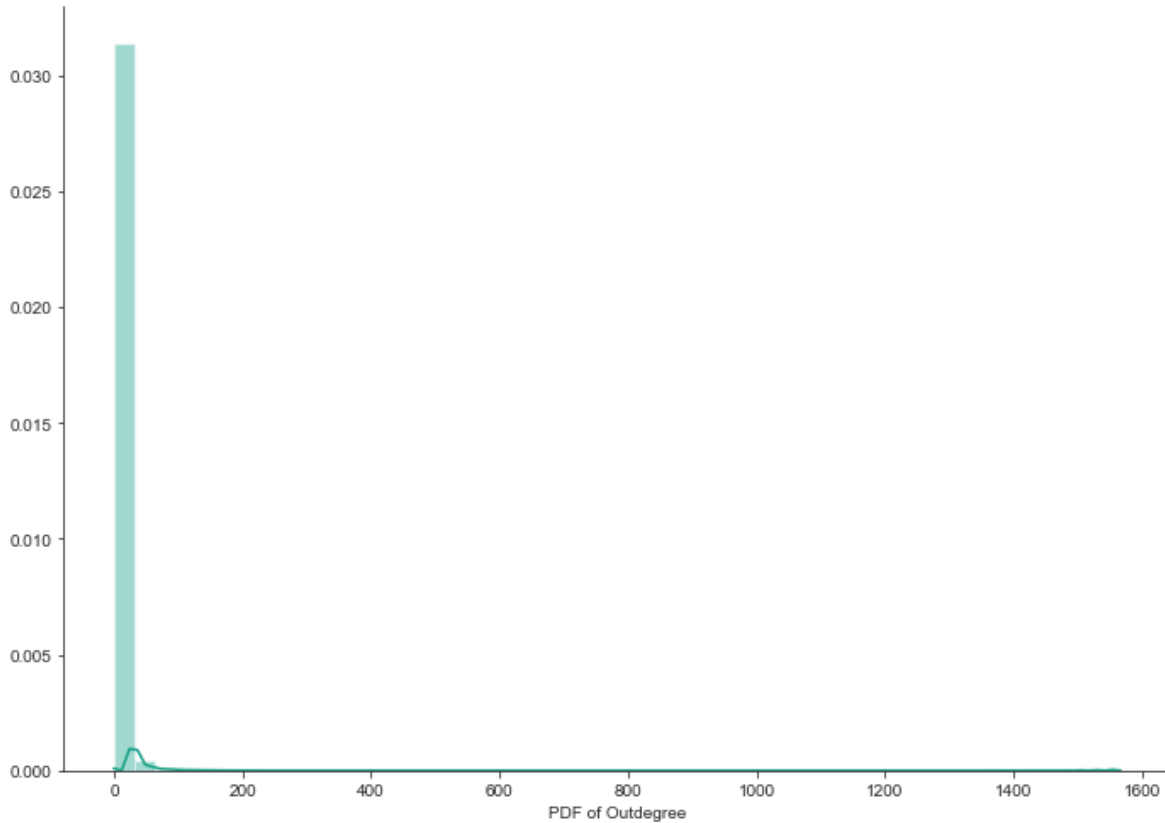
In [16]:

```
### 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 [17]:

```
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 [18]:

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

No of persons those are not following anyone are 274512 and % is 14.74111544  
2858524

In [19]:

```
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 [20]:

```

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 a

```

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

## 1.3 both followers + following

In [21]:

```

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()))

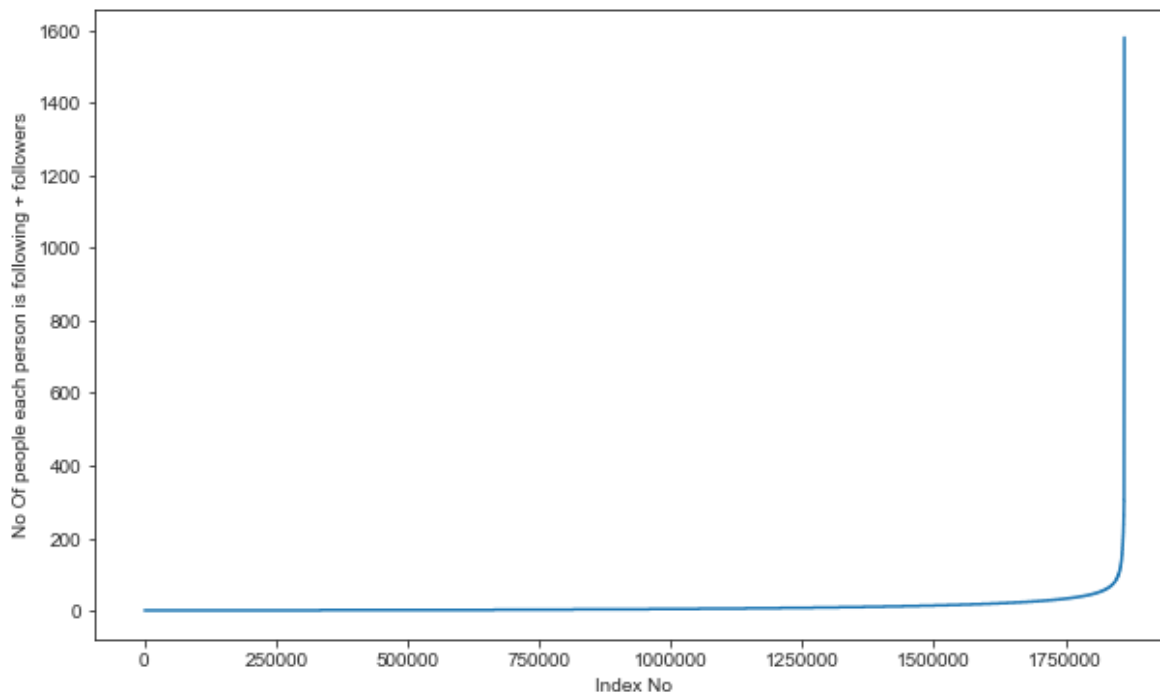
```

In [22]:

```

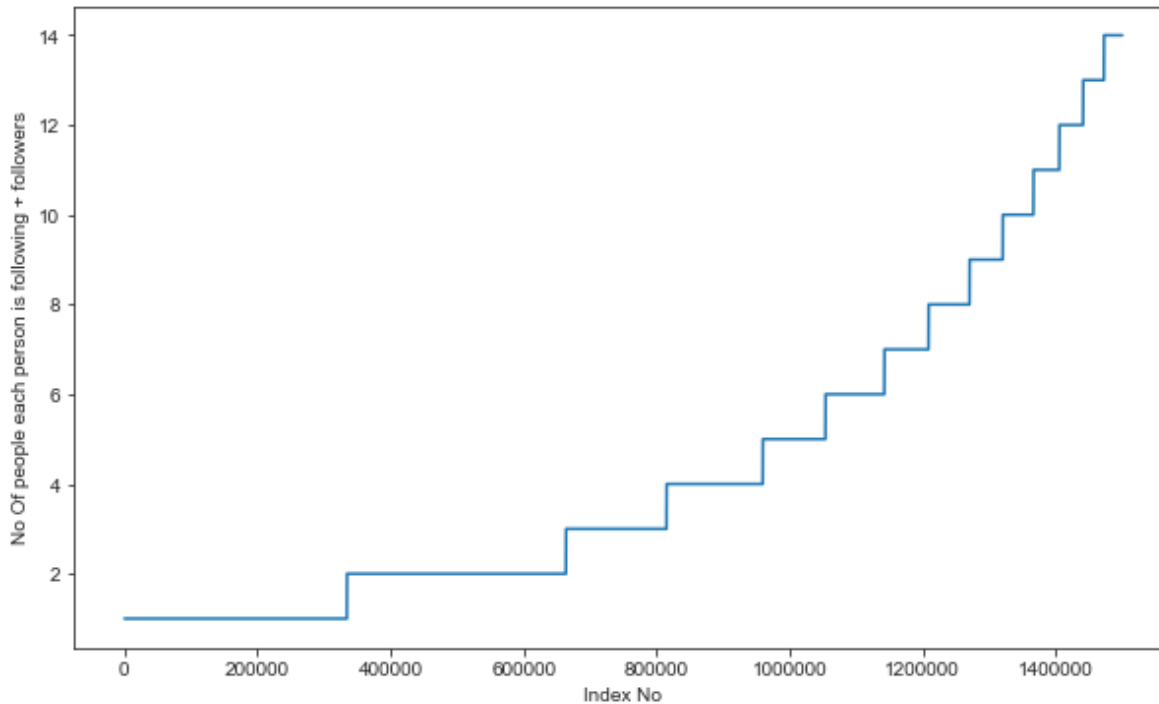
in_out_degree_sort = sorted(in_out_degree)
plt.figure(figsize=(10,6))
plt.plot(in_out_degree_sort)
plt.xlabel('Index No')
plt.ylabel('No Of people each person is following + followers')
plt.show()

```



In [23]:

```
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 [24]:

```
### 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 [25]:

```
### 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 [26]:

```
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 followers +
```

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

In [27]:

```
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 followers +
```

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

In [28]:

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

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

In [29]:

```
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 [30]:

```
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))
        except:
            continue
    else:
        continue
    pickle.dump(missing_edges, open('missing_edges_final.p', 'wb'))
else:
    missing_edges = pickle.load(open('missing_edges_final.p', 'rb'))
```

In [31]:

```
len(missing_edges)
```

Out[31]:

9437519

## 2.2 Training and Test data split:

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



In [32]:

```

from sklearn.model_selection import train_test_split
if (not os.path.isfile('data/after_eda/train_pos_after_eda.csv')) and (not os.path.isfile('
    #reading total data df
    df_pos = pd.read_csv('data/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 only
    #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)),
    X_train_neg, X_test_neg, y_train_neg, y_test_neg = train_test_split(df_neg, np.zeros(len(df_neg)),

    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_test_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('data/after_eda/train_pos_after_eda.csv', header=False, index=False)
    X_test_pos.to_csv('data/after_eda/test_pos_after_eda.csv', header=False, index=False)
    X_train_neg.to_csv('data/after_eda/train_neg_after_eda.csv', header=False, index=False)
    X_test_neg.to_csv('data/after_eda/test_neg_after_eda.csv', header=False, index=False)
else:
    #Graph from Traing data only
    del missing_edges

```

In [33]:

```

if (os.path.isfile('data/after_eda/train_pos_after_eda.csv')) and (os.path.isfile('data/aft
train_graph=nx.read_edgelist('data/after_eda/train_pos_after_eda.csv',delimiter=',',cre
test_graph=nx.read_edgelist('data/after_eda/test_pos_after_eda.csv',delimiter=',',creat
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 {} %'.f

```

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 7.1  
200735962845405 %

we have a cold start problem here

In [34]:

```
import os

X_train_pos = pd.read_csv('data/after_eda/train_pos_after_eda.csv', names=['source_node', 'de
X_test_pos = pd.read_csv('data/after_eda/test_pos_after_eda.csv', names=['source_node', 'de
X_train_neg = pd.read_csv('data/after_eda/train_neg_after_eda.csv', names=['source_node', 'de
X_test_neg = pd.read_csv('data/after_eda/test_neg_after_eda.csv', names=['source_node', 'de

X_train = X_train_pos.append(X_train_neg,ignore_index=True)
X_test = X_test_pos.append(X_test_neg,ignore_index=True)
y_train = pd.read_csv('data/train_y.csv')
y_test = pd.read_csv('data/test_y.csv')

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])
```

```
=====
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 [35]:

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

```
Data points in train data (15100030, 2)
Data points in test data (3775008, 2)
Shape of traget variable in train (15100029, 1)
Shape of traget variable in test (3775007, 1)
```

## 2.1 Jaccard Distance:

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

In [36]:

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

In [37]:

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

0.0

In [38]:

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

0.0

In [39]:

```
#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.predecesso
                                                                (len(set(train_graph.predecessors(a)).union(set(train_gra
    return sim
    except:
        return 0
```

In [40]:

```
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 [41]:

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

In [42]:

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

0.0

In [43]:

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

0

In [44]:

```
def cosine_for_followers(a,b):
    try:
        if len(set(train_graph.predecessors(a))) == 0 | len(set(train_graph.predecessors(b)
            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
                                                                return sim
    except:
        return 0
```

In [45]:

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

0.02886751345948129

0

### 3. Ranking Measures

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

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



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

## 3.1 Page Ranking

In [46]:

```
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 [47]:

```
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 [48]:

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

```
5.615699699389075e-07
```

## 4. Other Graph Features

### 4.1 Shortest path:

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

In [49]:

```
#Shortest path:Getting Shortest path between twoo nodes, if nodes have direct path i.e dire  
#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 [50]:

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

Out[50]:

10

In [51]:

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

Out[51]:

-1

## 4.2 Checking for same community

In [52]:

```
#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 [53]:

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

Out[53]:

0

In [54]:

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

Out[54]:

0

## 4.3 Adamic/Adar Index:



In [55]:

```
#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 [56]:

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

Out[56]:

0

## 4.4 Is person was following back:

In [57]:

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

In [58]:

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

Out[58]:

1

In [59]:

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

Out[59]:

0

## 4.5 Katz Centrality:

### 4.5 Katz Centrality:

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

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

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

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

The parameter

$$\beta$$

controls the initial centrality and

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

In [60]:

```
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 [61]:

```
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 [62]:

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

```
0.0007483800935562018
```

## Hits Score

### 4.6 Hits Score

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

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

In [63]:

```

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 [64]:

```

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

```

## Task 1: 4.7 Preferential Attachment

In [65]:

```

def preferential_attachment(a,b):
    try:
        return len(set(train_graph.predecessors(a))) * len(set(train_graph.predecessors(b)))
    except:
        return 0

```

In [66]:

```
preferential_attachment(1,189226)
```

Out[66]:

9

In [67]:

```
preferential_attachment(669354,1635354)
```

Out[67]:

0

## Task 2: 4.7 SVD Dot

In [68]:

```

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 [73]:

```

if os.path.isfile('data/after_eda/train_after_eda.csv'):
    filename = "data/after_eda/test_after_eda.csv"
    # you uncomment this line, if you dont know the length of the file name
    # here we have hardcoded the number of lines as 3775008
    # n_test = sum(1 for line in open(filename)) #number of records in file (excludes header)
    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 [74]:

```

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

```

Number of rows in the train data file: 15100028

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

Number of rows in the test data file: 3775006

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

In [95]:

```

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

```

Our train matrix size (100002, 3)

Out[95]:

	source_node	destination_node	indicator_link
0	273084	1505602	1
1	1443670	271786	1

In [96]:

```

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

```

Our test matrix size (50002, 3)

Out[96]:

	source_node	destination_node	indicator_link
0	848424	784690	1
1	1586725	1119166	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]:

```
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_follow
```

In [78]:

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

In [97]:

```

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['de
    df_final_test['jaccard_followers'] = df_final_test.apply(lambda row:
                                                                jaccard_for_followers(row['source_node'],row['de

    #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['de
    df_final_test['jaccard_followees'] = df_final_test.apply(lambda row:
                                                                jaccard_for_followees(row['source_node'],row['de

    #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['de
    df_final_test['cosine_followers'] = df_final_test.apply(lambda row:
                                                                cosine_for_followers(row['source_node'],row['de

    #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['de
    df_final_test['cosine_followees'] = df_final_test.apply(lambda row:
                                                                cosine_for_followees(row['source_node'],row['de

    df_final_train['num_followers_s'], df_final_train['num_followers_s'], \
    df_final_train['num_followees_s'], df_final_train['num_followees_d'], \
    df_final_train['inter_followers'], df_final_train['inter_followees']= compute_features_

    df_final_test['num_followers_s'], df_final_test['num_followers_s'], \
    df_final_test['num_followees_s'], df_final_test['num_followees_d'], \
    df_final_test['inter_followers'], df_final_test['inter_followees']= compute_features_st

    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 [98]:

```

from pandas import HDFStore
from pandas import read_hdf
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_
    #mapping adar index on test
    df_final_test['adar_index'] = df_final_test.apply(lambda row: calc_adar_in(row['source_

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

    #-----
    #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[
    ##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['s

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

    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

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





In [99]:

```

if not os.path.isfile('storage_sample_stage3.h5'):
    #mapping to pandas train
    df_final_train['weight_in'] = df_final_train.destination_node.apply(lambda x: Weight_in.g
    df_final_train['weight_out'] = df_final_train.source_node.apply(lambda x: Weight_out.ge

    #mapping to pandas test
    df_final_test['weight_in'] = df_final_test.destination_node.apply(lambda x: Weight_in.g
    df_final_test['weight_out'] = df_final_test.source_node.apply(lambda x: Weight_out.get(

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

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

    #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_p
    df_final_train['page_rank_d'] = df_final_train.destination_node.apply(lambda x: pr.get(x, m

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

    #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, mean_k
    df_final_train['katz_d'] = df_final_train.destination_node.apply(lambda x: katz.get(x, m

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

    #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(

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

    #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].ge
    df_final_train['authorities_d'] = df_final_train.destination_node.apply(lambda x: hits[

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

    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

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]:

```

#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 [84]:

```

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

```

In [85]:

```

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

```

In [86]:

```

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 [100]:

```

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', 'svd_u_s_6', 'svd_u_s_7', 'svd_u_s_8', 'svd_u_s_9', 'svd_u_s_10'],
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', 'svd_u_d_6', 'svd_u_d_7', 'svd_u_d_8', 'svd_u_d_9', 'svd_u_d_10'],
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', 'svd_v_s_6', 'svd_v_s_7', 'svd_v_s_8', 'svd_v_s_9', 'svd_v_s_10'],
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', 'svd_v_d_6', 'svd_v_d_7', 'svd_v_d_8', 'svd_v_d_9', 'svd_v_d_10'],
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', 'svd_u_s_7', 'svd_u_s_8', 'svd_u_s_9', 'svd_u_s_10'],
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', 'svd_u_d_7', 'svd_u_d_8', 'svd_u_d_9', 'svd_u_d_10'],
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', 'svd_v_s_7', 'svd_v_s_8', 'svd_v_s_9', 'svd_v_s_10'],
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', 'svd_v_d_7', 'svd_v_d_8', 'svd_v_d_9', 'svd_v_d_10'],
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()

```

## ML Models

### Adding new features

- Preferential Attachment
- SVD Dot

In [101]:

```

if not os.path.isfile('storage_sample_stage5.h5'):

    #mapping preferential attachment on train
    df_final_train['pref_at'] = df_final_train.apply(lambda row: preferential_attachment(row), axis=1)
    #mapping preferential attachment on test
    df_final_test['pref_at'] = df_final_test.apply(lambda row: preferential_attachment(row), axis=1)

    #-----

    #adding svd_dot
    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)

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

```

In [102]:

```

y_train = df_final_train['indicator_link']
y_test = df_final_test['indicator_link']

```

In [103]:

```
df_final_train.columns
```

Out[103]:

```

Index(['source_node', 'destination_node', 'indicator_link',
      'jaccard_followers', 'jaccard_followees', 'cosine_followers',
      'cosine_followees', 'num_followers_s', 'num_followees_s',
      'num_followees_d', 'inter_followers', 'inter_followees', 'adar_index',
      'follows_back', 'same_comp', 'shortest_path', 'weight_in', 'weight_out',
      'weight_f1', 'weight_f2', 'weight_f3', 'weight_f4', 'page_rank_s',
      'page_rank_d', 'katz_s', 'katz_d', 'hubs_s', 'hubs_d', 'authorities_s',
      'authorities_d', 'svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4',
      'svd_u_s_5', 'svd_u_s_6', 'svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
      'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6', 'svd_v_s_1', 'svd_v_s_2',
      'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6', 'svd_v_d_1',
      'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6',
      'pref_at', 'svd_dot_u', 'svd_dot_v'],
      dtype='object')

```

In [104]:

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

## Random Forest model with hyperparameter tuning

In [105]:

```

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
estimators = [10,50,100,250,450]
train_scores = []
test_scores = []
for i in estimators:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=5, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=52, min_samples_split=120,
                                min_weight_fraction_leaf=0.0, n_estimators=i, n_jobs=-1, random_state=25, verbose=
    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')

```

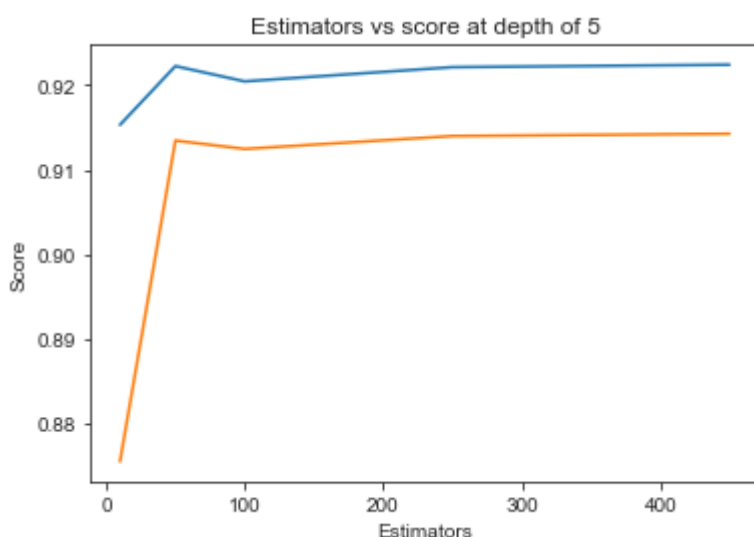
```

Estimators = 10 Train Score 0.9152933524950435 test Score 0.875474000784553
1
Estimators = 50 Train Score 0.9222421263122812 test Score 0.913431203466117
1
Estimators = 100 Train Score 0.9204354847111528 test Score 0.91244608100999
47
Estimators = 250 Train Score 0.9221182502136174 test Score 0.91397419436317
33
Estimators = 450 Train Score 0.9224079500307614 test Score 0.91423757371524
85

```

Out[105]:

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



In [106]:

```

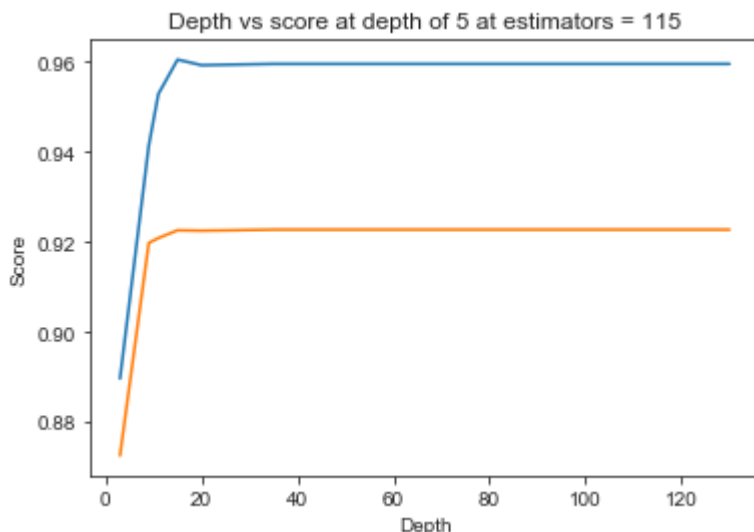
depths = [3,9,11,15,20,35,50,70,130]
train_scores = []
test_scores = []
for i in depths:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=i, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=52, min_samples_split=120,
                                min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=-1, random_state=25, verbose=0)
    clf.fit(df_final_train, y_train)
    train_sc = f1_score(y_train, clf.predict(df_final_train))
    test_sc = f1_score(y_test, clf.predict(df_final_test))
    test_scores.append(test_sc)
    train_scores.append(train_sc)
    print('depth = ', i, 'Train Score', train_sc, 'test Score', test_sc)
plt.plot(depths, train_scores, label='Train Score')
plt.plot(depths, test_scores, label='Test Score')
plt.xlabel('Depth')
plt.ylabel('Score')
plt.title('Depth vs score at depth of 5 at estimators = 115')
plt.show()

```

```

depth = 3 Train Score 0.8896223341076553 test Score 0.872528036983135
depth = 9 Train Score 0.9413566513168106 test Score 0.9197426346088723
depth = 11 Train Score 0.9528000902758543 test Score 0.9208178832503443
depth = 15 Train Score 0.9604212498981505 test Score 0.9225502419644555
depth = 20 Train Score 0.9591786749397886 test Score 0.9224553684744903
depth = 35 Train Score 0.9594686961223385 test Score 0.9226896799712334
depth = 50 Train Score 0.9594686961223385 test Score 0.9226896799712334
depth = 70 Train Score 0.9594686961223385 test Score 0.9226896799712334
depth = 130 Train Score 0.9594686961223385 test Score 0.9226896799712334

```





In [107]:

```

from sklearn.ensemble import RandomForestClassifier
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_train_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 [0.95888362 0.95834514 0.9533171  0.95761826 0.96074153]
mean train scores [0.95985373 0.95938841 0.95406853 0.95874168 0.96208994]

```

In [108]:

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

```

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

```

## Random Forest with best hyperparameter

In [109]:

```

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 [110]:

```

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

```

## Evaluating model performance

In [111]:

```

from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score
from sklearn.metrics import f1_score
from sklearn.metrics import recall_score

print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_test_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_test_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_test_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_test_pred)))

```

Accuracy on test set: 92.776%  
Precision on test set: 0.978  
Recall on test set: 0.875  
F1-Score on test set: 0.924

## Confusion matrix

In [112]:

```

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

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

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

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

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

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

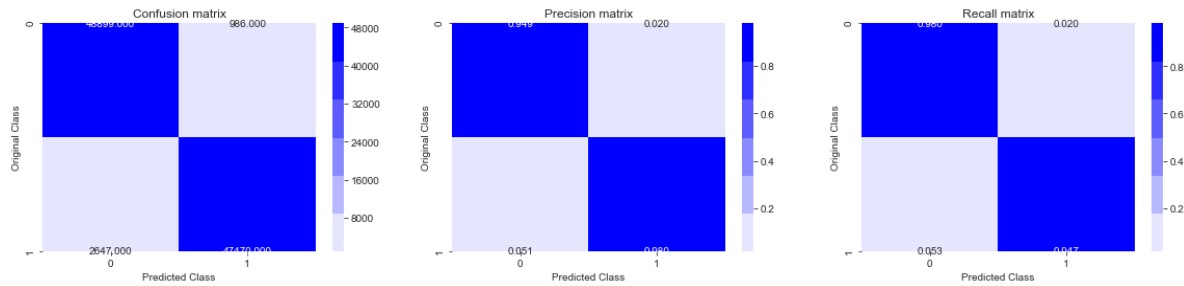
    plt.show()

```

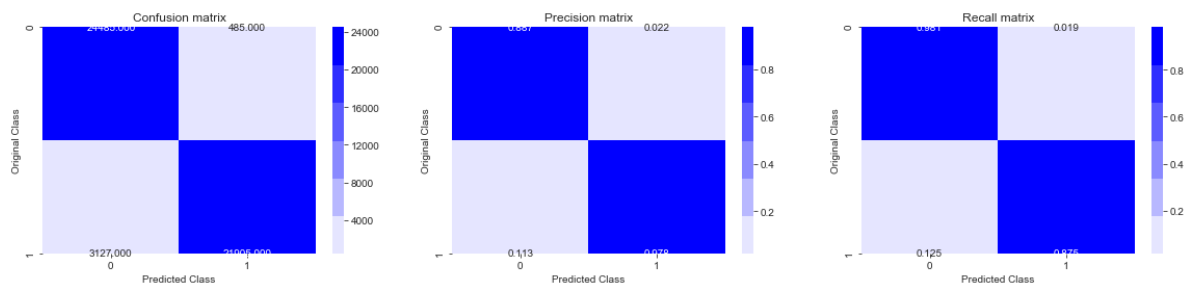
In [113]:

```
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



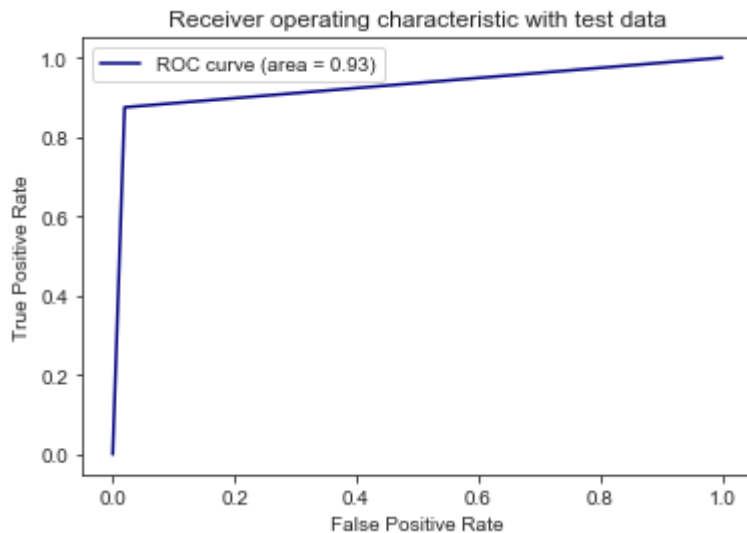
### Test confusion\_matrix



## AUC-ROC Curve

In [114]:

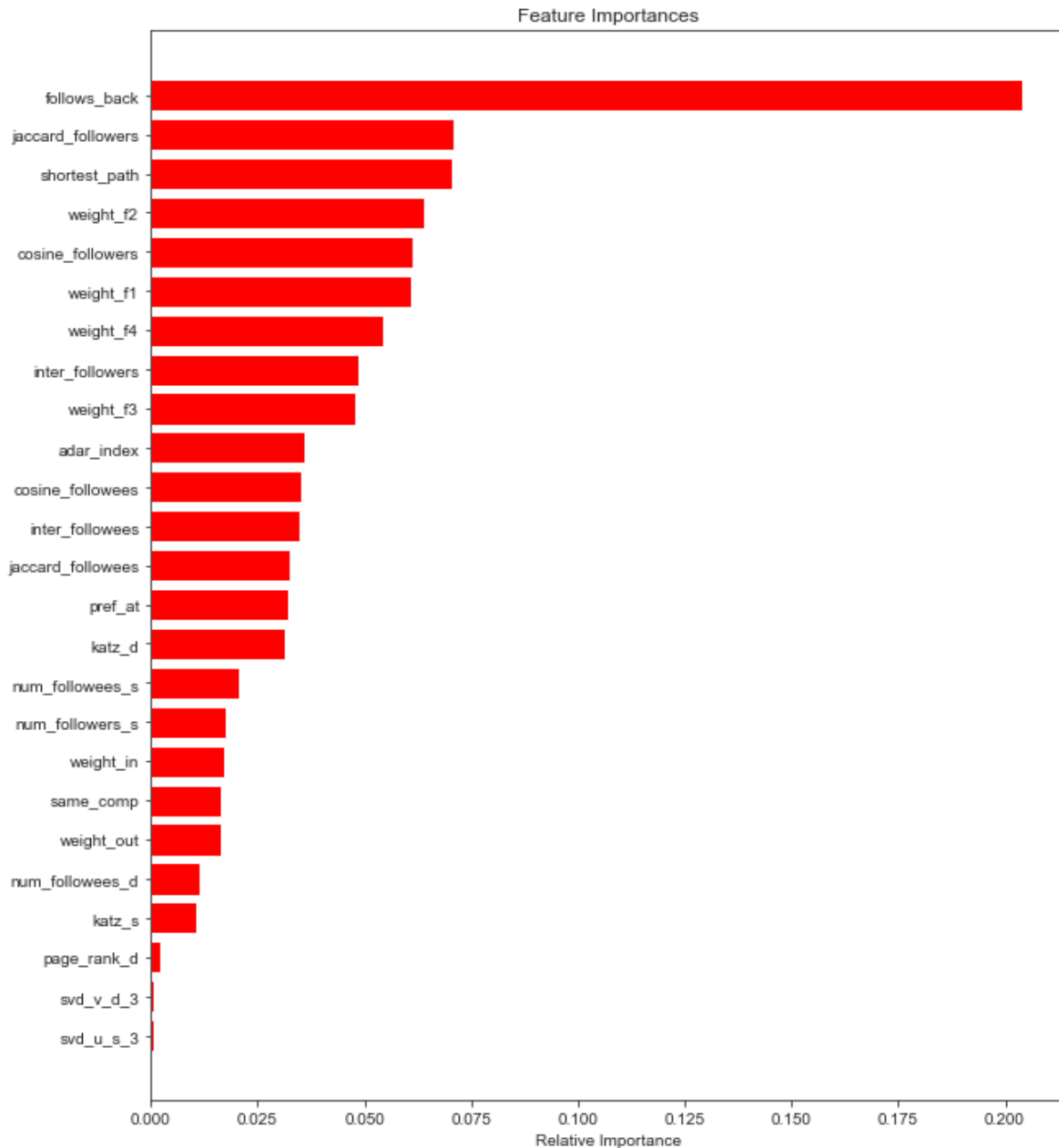
```
from sklearn.metrics import roc_curve, auc
fpr, tpr, ths = roc_curve(y_test, y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy', label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic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



## Important features

In [115]:

```
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()
```



## Task: XGBoost Model with hyperparameter tuning

### Hyperparameter tuning

In [121]:

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

clf = xgb.XGBClassifier()
param_grid = {"n_estimators": sp_randint(100,150),
              "max_depth": sp_randint(2,10)}

model = RandomizedSearchCV(clf, param_distributions=param_grid, return_train_score = True)

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

```
mean train scores [0.96675066 0.97280054 0.9947801  0.97571548 0.9797004  0.
99268515
 0.98751025 0.96902061 0.96649566 0.98765525]
mean test scores [0.96608068 0.97245055 0.97924042 0.97427051 0.97632047 0.9
7916042
 0.97814044 0.96850063 0.96585068 0.97786044]
```

In [122]:

```
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=9,
              min_child_weight=1, missing=None, n_estimators=137, 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)
```

## XGBoost model with best parameters

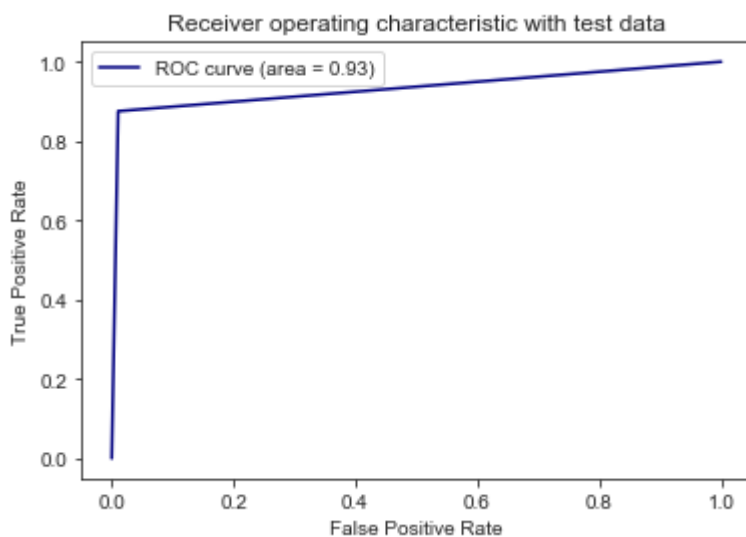
In [123]:

```
#XGBboot classifier
clf=xgb.XGBClassifier(max_depth=9,n_estimators=137)
#fitting model
clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)
```

## AUC ROC Curve

In [124]:

```
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()
```



## Evaluating model performance

In [125]:

```

from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score
from sklearn.metrics import f1_score
from sklearn.metrics import recall_score

print("Accuracy on test set: %0.3f%%"%(accuracy_score(y_test, y_test_pred)*100))
print("Precision on test set: %0.3f"%(precision_score(y_test, y_test_pred)))
print("Recall on test set: %0.3f"%(recall_score(y_test, y_test_pred)))
print("F1-Score on test set: %0.3f"%(f1_score(y_test, y_test_pred)))

```

Accuracy on test set: 93.234%  
Precision on test set: 0.988  
Recall on test set: 0.875  
F1-Score on test set: 0.928

## Confusion matrix

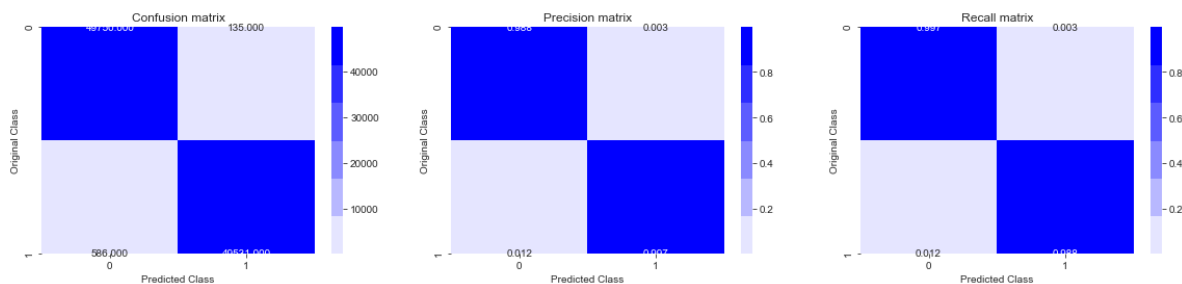
In [126]:

```

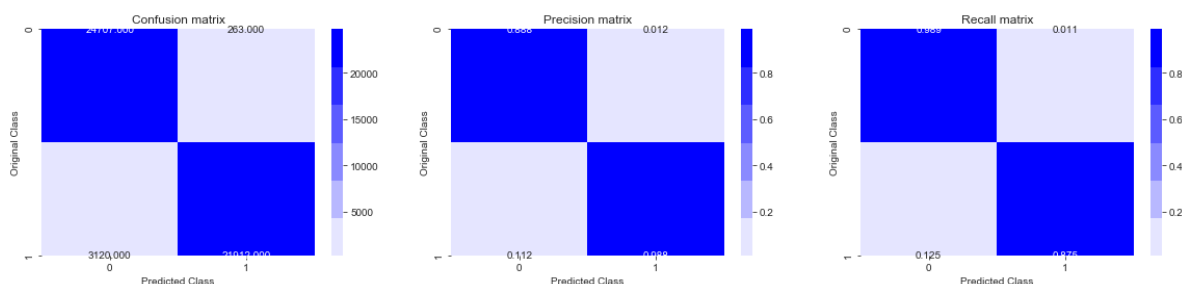
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)

```

Train confusion\_matrix



Test confusion\_matrix

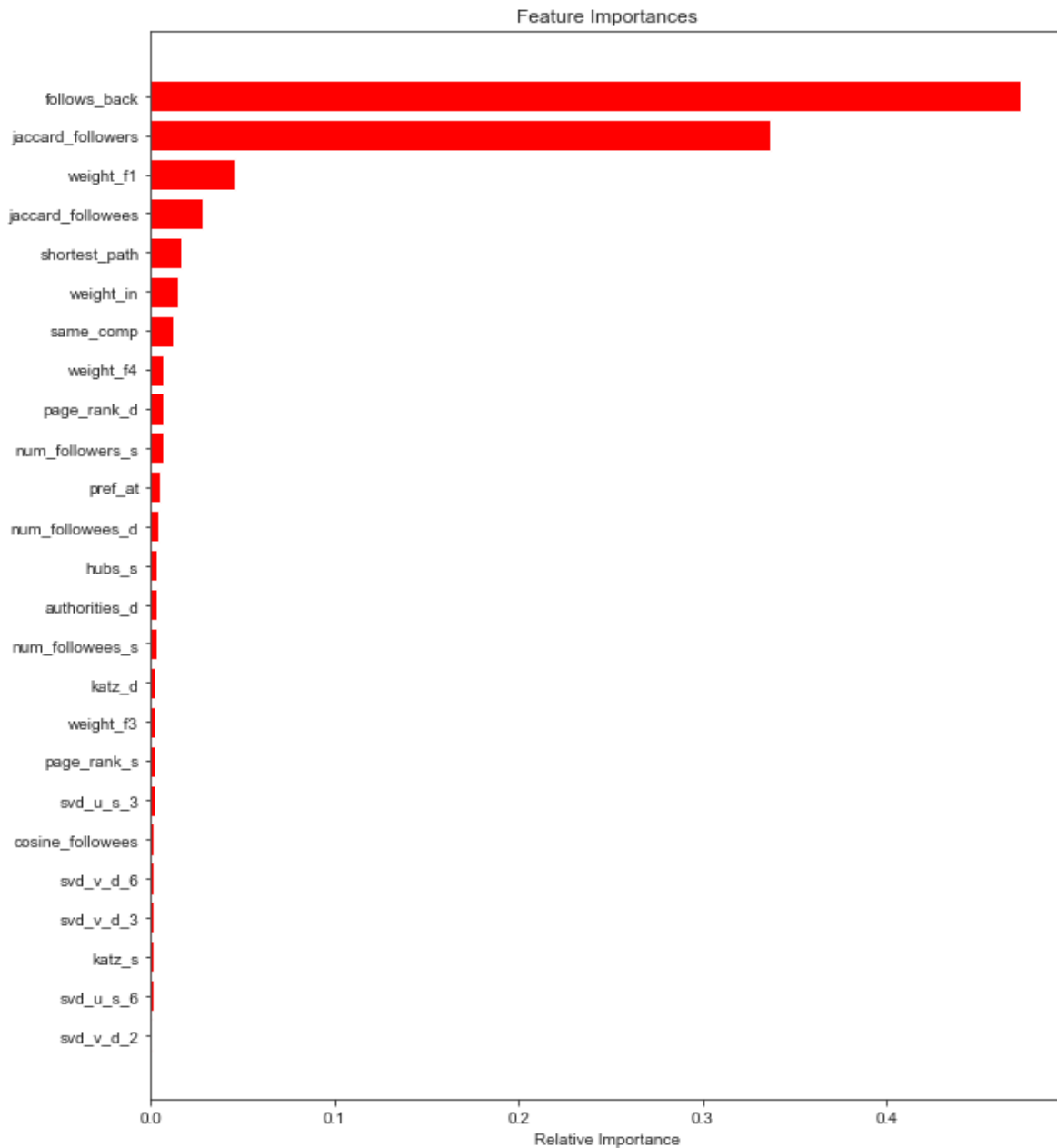


## Important features



In [127]:

```
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()
```



# Conclusion

In [128]:

```
# Please compare all your models using Prettytable library # http://zetcode.com/python/prettytable
from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["Model", "Best param", "F1-Score", "Accuracy on test set"]
x.add_row(["Random Forest", "Max Depth:14 , n_estimators : 121", 0.92, "92.77%"])
x.add_row(["XGBoost", "Max Depth:7 , Estimators : 117", 0.92, "93.23%"])
print(x)
```

```
+-----+-----+-----+-----+
-----+
|      Model      |      Best param      | F1-Score | Accuracy on
test set |
+-----+-----+-----+-----+
-----+
| Random Forest | Max Depth:14 , n_estimators : 121 |    0.92    |      92.7
7%      |
|    XGBoost    | Max Depth:7 , Estimators : 117 |    0.92    |      93.2
3%      |
+-----+-----+-----+-----+
-----+
```

## Procedure:

1. Directed graph have been created out of the given dataset.
2. Generated samples of bad links for each link from the directed graph.
3. Added features like jaccard distance, cosine similarity etc and two new features namely preferential attachment and SVD dot.
4. Build Random forest model and XGBoost with the best set of hyperparameters obtained using hyperparameter tuning for both the models.
5. Summarized the models performance using pretty table.