# Social network Graph Link Prediction - Facebook Challenge

In [1]:
```python
#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 aritmetic operations on arrays
# matplotlib: used to plot graphs
import matplotlib
import matplotlib.pylab 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
```

## 1. Reading Data

In [2]:
```python
if os.path.isfile('data/after_eda/train_pos_after_eda.csv'):
    train_graph=nx.read_edgelist('data/after_eda/train_pos_after_eda.csv',de
    print(nx.info(train_graph))
else:
    print("please run the FB_EDA.ipynb or download the files from drive")
```

```
Name:
Type: DiGraph
Number of nodes: 1780722
Number of edges: 7550015
Average in degree:   4.2399
Average out degree:   4.2399
```

In [ ]:
```
```

In [ ]:　　1

# 2. Similarity measures

## 2.1 Jaccard Distance:

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

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

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

In [4]:
```
1   #one test case
2   print(jaccard_for_followees(273084,1505602))
```

0.0

In [5]:
```
1   #node 1635354 not in graph
2   print(jaccard_for_followees(273084,1505602))
```

0.0

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

In [7]:
```
1   print(jaccard_for_followers(273084,470294))
```

0

```
In [8]:   1  #node 1635354 not in graph
          2  print(jaccard_for_followees(669354,1635354))
```

0

## 2.2 Cosine distance

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

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

```
In [10]:  1  print(cosine_for_followees(273084,1505602))
```

0.0

```
In [11]:  1  print(cosine_for_followees(273084,1635354))
```

0

```
In [12]:  1  def cosine_for_followers(a,b):
          2      try:
          3
          4          if len(set(train_graph.predecessors(a))) == 0  | len(set(train_graph
          5              return 0
          6          sim = (len(set(train_graph.predecessors(a)).intersection(set(train_g
          7                                  (math.sqrt(len(set(train_graph.predeces
          8          return sim
          9      except:
         10          return 0
```

```
In [13]:  1  print(cosine_for_followers(2,470294))
```

0.02886751345948129

```
In [14]:  1  print(cosine_for_followers(669354,1635354))
```

0

# Preferential attachment

```python
In [15]:   1  def preferential_attachment_followers(a,b):
           2      try:
           3          if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.pr
           4              return 0
           5          p_a = len(set(train_graph.successors(a)).intersection(set(train_grap
           6      except:
           7          return 0
           8      return p_a
           9
```

```python
In [16]:   1  def preferential_attachment_followees(a,b):
           2      try:
           3          if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.pr
           4              return 0
           5          p_a = len(set(train_graph.predecessors(a)).intersection(set(train_gr
           6      except:
           7          return 0
           8      return p_a
           9
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

```python
In [ ]:   1
```

# 3. Ranking Measures

https://networkx.github.io/documentation/networkx-
1.10/reference/generated/networkx.algorithms.link_analysis.pagerank_alg.pagerank.html
(https://networkx.github.io/documentation/networkx-
1.10/reference/generated/networkx.algorithms.link_analysis.pagerank_alg.pagerank.html)

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



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

## 3.1 Page Ranking

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

```
In [17]:  1  if not os.path.isfile('data/fea_sample/page_rank.p'):
          2      pr = nx.pagerank(train_graph, alpha=0.85)
          3      pickle.dump(pr,open('data/fea_sample/page_rank.p','wb'))
          4  else:
          5      pr = pickle.load(open('data/fea_sample/page_rank.p','rb'))
```

```
In [18]:  1  print('min',pr[min(pr, key=pr.get)])
          2  print('max',pr[max(pr, key=pr.get)])
          3  print('mean',float(sum(pr.values())) / len(pr))
```

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

```
In [19]:  1  #for imputing to nodes which are not there in Train data
          2  mean_pr = float(sum(pr.values())) / len(pr)
          3  print(mean_pr)
```

```
5.615699699389075e-07
```

# 4. Other Graph Features

## 4.1 Shortest path:

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

```
In [20]:    1  #if has direct edge then deleting that edge and calculating shortest path
            2  def compute_shortest_path_length(a,b):
            3      p=-1
            4      try:
            5          if train_graph.has_edge(a,b):
            6              train_graph.remove_edge(a,b)
            7              p= nx.shortest_path_length(train_graph,source=a,target=b)
            8              train_graph.add_edge(a,b)
            9          else:
           10              p= nx.shortest_path_length(train_graph,source=a,target=b)
           11          return p
           12      except:
           13          return -1
```

```
In [21]:    1  #testing
            2  compute_shortest_path_length(77697, 826021)
```

Out[21]:  10

```
In [22]:    1  #testing
            2  compute_shortest_path_length(669354,1635354)
```

Out[22]:  -1

## 4.2 Checking for same community

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

```
In [24]:   1  belongs_to_same_wcc(861, 1659750)
```

Out[24]:  0

```
In [25]:   1  belongs_to_same_wcc(669354,1635354)
```

Out[25]:  0

## 4.3 Adamic/Adar Index:

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

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

```
In [26]:    1  #adar index
            2  def calc_adar_in(a,b):
            3      sum=0
            4      try:
            5          n=list(set(train_graph.successors(a)).intersection(set(train_graph.s
            6          if len(n)!=0:
            7              for i in n:
            8                  sum=sum+(1/np.log10(len(list(train_graph.predecessors(i)))))
            9              return sum
           10          else:
           11              return 0
           12      except:
           13          return 0
```

```
In [27]:    1  calc_adar_in(1,189226)
```

Out[27]: 0

```
In [28]:    1  calc_adar_in(669354,1635354)
```

Out[28]: 0

## 4.4 Is the person following back:

```
In [29]:    1  def follows_back(a,b):
            2      if train_graph.has_edge(b,a):
            3          return 1
            4      else:
            5          return 0
```

```
In [30]:    1  follows_back(1,189226)
```

Out[30]: 1

```
In [31]:    1  follows_back(669354,1635354)
```

Out[31]: 0

## 4.5 Katz Centrality:

https://en.wikipedia.org/wiki/Katz_centrality (https://en.wikipedia.org/wiki/Katz_centrality)

https://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

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

```python
In [32]: 1  if not os.path.isfile('data/fea_sample/katz.p'):
         2      katz = nx.katz.katz_centrality(train_graph,alpha=0.005,beta=1)
         3      pickle.dump(katz,open('data/fea_sample/katz.p','wb'))
         4  else:
         5      katz = pickle.load(open('data/fea_sample/katz.p','rb'))
```

```python
In [33]: 1  print('min',katz[min(katz, key=katz.get)])
         2  print('max',katz[max(katz, key=katz.get)])
         3  print('mean',float(sum(katz.values())) / len(katz))
```

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

```python
In [34]: 1  mean_katz = float(sum(katz.values())) / len(katz)
         2  print(mean_katz)
```

```
0.0007483800935562018
```

## 4.6 Hits Score

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

https://en.wikipedia.org/wiki/HITS_algorithm (https://en.wikipedia.org/wiki/HITS_algorithm)

```python
In [35]: 1  if not os.path.isfile('data/fea_sample/hits.p'):
         2      hits = nx.hits(train_graph, max_iter=100, tol=1e-08, nstart=None, normal
         3      pickle.dump(hits,open('data/fea_sample/hits.p','wb'))
         4  else:
         5      hits = pickle.load(open('data/fea_sample/hits.p','rb'))
```

```python
In [36]: 1  print('min',hits[0][min(hits[0], key=hits[0].get)])
         2  print('max',hits[0][max(hits[0], key=hits[0].get)])
         3  print('mean',float(sum(hits[0].values())) / len(hits[0]))
```

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

# 5. Featurization

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

```
In [37]:  1  import random
          2  if os.path.isfile('data/after_eda/train_after_eda.csv'):
          3      filename = "data/after_eda/train_after_eda.csv"
          4      # you uncomment this line, if you dont know the lentgh of the file name
          5      # here we have hardcoded the number of lines as 15100030
          6      # n_train = sum(1 for line in open(filename)) #number of records in file
          7      n_train =  15100028
          8      s = 100000 #desired sample size
          9      skip_train = sorted(random.sample(range(1,n_train+1),n_train-s))
         10      #https://stackoverflow.com/a/22259008/4084039
```

```
In [38]:  1  if os.path.isfile('data/after_eda/train_after_eda.csv'):
          2      filename = "data/after_eda/test_after_eda.csv"
          3      # you uncomment this line, if you dont know the lentgh of the file name
          4      # here we have hardcoded the number of lines as 3775008
          5      # n_test = sum(1 for line in open(filename)) #number of records in file
          6      n_test = 3775006
          7      s = 50000 #desired sample size
          8      skip_test = sorted(random.sample(range(1,n_test+1),n_test-s))
          9      #https://stackoverflow.com/a/22259008/4084039
```

```
In [39]:  1  print("Number of rows in the train data file:", n_train)
          2  print("Number of rows we are going to elimiate in train data are",len(skip_t
          3  print("Number of rows in the test data file:", n_test)
          4  print("Number of rows we are going to elimiate in test data are",len(skip_te
```

```
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 [40]:  1  df_final_train = pd.read_csv('data/after_eda/train_after_eda.csv', skiprows=
          2  df_final_train['indicator_link'] = pd.read_csv('data/train_y.csv', skiprows=
          3  print("Our train matrix size ",df_final_train.shape)
          4  df_final_train.head(2)
```

```
Our train matrix size  (100002, 3)
```

Out[40]:

| | source_node | destination_node | indicator_link |
|---|---|---|---|
| **0** | 273084 | 1505602 | 1 |
| **1** | 1814537 | 613441 | 1 |

In [41]:
```python
df_final_test = pd.read_csv('data/after_eda/test_after_eda.csv', skiprows=sk
df_final_test['indicator_link'] = pd.read_csv('data/test_y.csv', skiprows=sk
print("Our test matrix size ",df_final_test.shape)
df_final_test.head(2)
```

Our test matrix size  (50002, 3)

Out[41]:

|   | source_node | destination_node | indicator_link |
|---|---|---|---|
| 0 | 848424 | 784690 | 1 |
| 1 | 548473 | 1721521 | 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 [46]:
```python
if not os.path.isfile('data/fea_sample/storage_sample_stage1.h5'):
    #mapping jaccrd followers to train and test data
    df_final_train['jaccard_followers'] = df_final_train.apply(lambda row:
                                    jaccard_for_followers(row['sourc
    df_final_test['jaccard_followers'] = df_final_test.apply(lambda row:
                                    jaccard_for_followers(row['sourc

    #mapping jaccrd followees to train and test data
    df_final_train['jaccard_followees'] = df_final_train.apply(lambda row:
                                    jaccard_for_followees(row['sourc
    df_final_test['jaccard_followees'] = df_final_test.apply(lambda row:
                                    jaccard_for_followees(row['sourc


        #mapping jaccrd followers to train and test data
    df_final_train['cosine_followers'] = df_final_train.apply(lambda row:
                                    cosine_for_followers(row['source
    df_final_test['cosine_followers'] = df_final_test.apply(lambda row:
                                    cosine_for_followers(row['source

    #mapping jaccrd followees to train and test data
    df_final_train['cosine_followees'] = df_final_train.apply(lambda row:
                                    cosine_for_followees(row['source
    df_final_test['cosine_followees'] = df_final_test.apply(lambda row:
                                    cosine_for_followees(row['source
```

In [60]:
```python
asdad = df_final_train.apply(lambda row:jaccard_for_followers(row['source_no
```

In [47]:
```python
# #if anything not there in train graph then adding mean page rank
# df_final_train['preferential_attachment_s'] = df_final_train.source_node.a
# df_final_train['preferential_attachment_d'] = df_final_train.destination_n

# df_final_test['preferential_attachment_s'] = df_final_test.source_node.app
# df_final_test['preferential_attachment_d'] = df_final_test.destination_nod
```

```python
In [48]:
def compute_features_stage1(df_final):
    #calculating no of followers followees for source and destination
    #calculating intersection of followers and followees for source and dest
    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_
```

```python
In [49]:
if not os.path.isfile('data/fea_sample/storage_sample_stage1.h5'):
    df_final_train['num_followers_s'], df_final_train['num_followers_d'], \
    df_final_train['num_followees_s'], df_final_train['num_followees_d'], \
    df_final_train['inter_followers'], df_final_train['inter_followees']= co

    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']= comp

    hdf = HDFStore('data/fea_sample/storage_sample_stage1.h5')
    hdf.put('train_df',df_final_train, format='table', data_columns=True)
    hdf.put('test_df',df_final_test, format='table', data_columns=True)
    hdf.close()
else:
    df_final_train = read_hdf('data/fea_sample/storage_sample_stage1.h5', 't
    df_final_test = read_hdf('data/fea_sample/storage_sample_stage1.h5', 'te
```

## 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 [50]:   1  if not os.path.isfile('data/fea_sample/storage_sample_stage2.h5'):
           2      #mapping adar index on train
           3      df_final_train['adar_index'] = df_final_train.apply(lambda row: calc_ada
           4      #mapping adar index on test
           5      df_final_test['adar_index'] = df_final_test.apply(lambda row: calc_adar_
           6
           7      #-----------------------------------------------------------------
           8      #mapping followback or not on train
           9      df_final_train['follows_back'] = df_final_train.apply(lambda row: follow
          10
          11      #mapping followback or not on test
          12      df_final_test['follows_back'] = df_final_test.apply(lambda row: follows_
          13
          14      #-----------------------------------------------------------------
          15      #mapping same component of wcc or not on train
          16      df_final_train['same_comp'] = df_final_train.apply(lambda row: belongs_t
          17
          18      ##mapping same component of wcc or not on train
          19      df_final_test['same_comp'] = df_final_test.apply(lambda row: belongs_to_
          20
          21      #-----------------------------------------------------------------
          22      #mapping shortest path on train
          23      df_final_train['shortest_path'] = df_final_train.apply(lambda row: compu
          24      #mapping shortest path on test
          25      df_final_test['shortest_path'] = df_final_test.apply(lambda row: compute
          26
          27      hdf = HDFStore('data/fea_sample/storage_sample_stage2.h5')
          28      hdf.put('train_df',df_final_train, format='table', data_columns=True)
          29      hdf.put('test_df',df_final_test, format='table', data_columns=True)
          30      hdf.close()
          31  else:
          32      df_final_train = read_hdf('data/fea_sample/storage_sample_stage2.h5', 't
          33      df_final_test = read_hdf('data/fea_sample/storage_sample_stage2.h5', 'te
```

## 5.4 Adding new set of features

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

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

4. katz of source
5. katz of dest
6. hubs of source
7. hubs of dest
8. authorities_s of source
9. authorities_s of dest

**Weight Features**

In order to determine the similarity of nodes, an edge weight value was calculated between nodes. Edge weight decreases as the neighbor count goes up. Intuitively, consider one million people following a celebrity on a social network then chances are most of them never met each other or the celebrity. On the other hand, if a user has 30 contacts in his/her social network, the chances are higher that many of them know each other. `credit` - Graph-based Features for Supervised Link Prediction William Cukierski, Benjamin Hamner, Bo Yang

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

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

```
In [54]:
1  #weight for source and destination of each link
2  Weight_in = {}
3  Weight_out = {}
4  for i in  tqdm(train_graph.nodes()):
5      s1=set(train_graph.predecessors(i))
6      w_in = 1.0/(np.sqrt(1+len(s1)))
7      Weight_in[i]=w_in
8
9      s2=set(train_graph.successors(i))
10     w_out = 1.0/(np.sqrt(1+len(s2)))
11     Weight_out[i]=w_out
12
13 #for imputing with mean
14 mean_weight_in = np.mean(list(Weight_in.values()))
15 mean_weight_out = np.mean(list(Weight_out.values()))
```

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

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


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

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

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

# 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 [42]:
```python
def svd(x, S):
    try:
        z = sadj_dict[x]
        return S[z]
    except:
        return [0,0,0,0,0,0]
```

In [43]:
```python
#for svd features to get feature vector creating a dict node val and inedx i
sadj_col = sorted(train_graph.nodes())
sadj_dict = { val:idx for idx,val in enumerate(sadj_col)}
```

In [44]:
```python
Adj = nx.adjacency_matrix(train_graph,nodelist=sorted(train_graph.nodes())).
```

In [45]:
```python
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,)
```

```python
In [0]:    1  if not os.path.isfile('data/fea_sample/storage_sample_stage4.h5'):
           2      #=================================================================
           3
           4      df_final_train[['svd_u_s_1', 'svd_u_s_2','svd_u_s_3', 'svd_u_s_4', 'svd_
           5      df_final_train.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)
           6
           7      df_final_train[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd
           8      df_final_train.destination_node.apply(lambda x: svd(x, U)).apply(pd.Seri
           9      #=================================================================
          10
          11      df_final_train[['svd_v_s_1','svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_
          12      df_final_train.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
          13
          14      df_final_train[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd
          15      df_final_train.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Se
          16      #=================================================================
          17
          18      df_final_test[['svd_u_s_1', 'svd_u_s_2','svd_u_s_3', 'svd_u_s_4', 'svd_u
          19      df_final_test.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)
          20
          21      df_final_test[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_
          22      df_final_test.destination_node.apply(lambda x: svd(x, U)).apply(pd.Serie
          23
          24      #=================================================================
          25
          26      df_final_test[['svd_v_s_1','svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v
          27      df_final_test.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
          28
          29      df_final_test[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_
          30      df_final_test.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Ser
          31      #=================================================================
          32
          33      hdf = HDFStore('data/fea_sample/storage_sample_stage4.h5')
          34      hdf.put('train_df',df_final_train, format='table', data_columns=True)
          35      hdf.put('test_df',df_final_test, format='table', data_columns=True)
          36      hdf.close()
```

```python
In [ ]:    1
```

```python
In [ ]:    1
```

```python
In [ ]:    1
```

```python
In [ ]:    1
```

```python
In [104]:  1  from pandas import read_hdf
           2  df_final_train = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'train
           3  df_final_test = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'test_d
```

In [105]:
```
1  df_final_train.describe()
```

Out[105]:

| | source_node | destination_node | indicator_link | jaccard_followers | jaccard_followees | cosine_ |
|---|---|---|---|---|---|---|
| count | 1.000020e+05 | 1.000020e+05 | 100002.000000 | 100002.0 | 100002.000000 | 1000 |
| mean | 9.305942e+05 | 9.280010e+05 | 0.500490 | 0.0 | 0.040088 | |
| std | 5.376040e+05 | 5.383887e+05 | 0.500002 | 0.0 | 0.104141 | |
| min | 3.000000e+00 | 7.000000e+00 | 0.000000 | 0.0 | 0.000000 | |
| 25% | 4.654765e+05 | 4.598682e+05 | 0.000000 | 0.0 | 0.000000 | |
| 50% | 9.309050e+05 | 9.290790e+05 | 1.000000 | 0.0 | 0.000000 | |
| 75% | 1.395303e+06 | 1.392905e+06 | 1.000000 | 0.0 | 0.000000 | |
| max | 1.862198e+06 | 1.862218e+06 | 1.000000 | 0.0 | 0.833333 | |

8 rows × 54 columns

In [96]:
```
1  # prepared and stored the data from machine learning models
2  # pelase check the FB_Models.ipynbb
```

In [97]:
```
1  svd_dot_train_U_U = []
2  for i,a in tqdm(df_final_train.iterrows(), total = df_final_train.shape[0]):
3      s = sadj_dict.get(a.source_node, 'NA')
4      d = sadj_dict.get(a.destination_node, 'NA')
5      if (s!='NA') and (d!='NA'):
6          aaa = np.dot(U[s,:], U[d,:])
7          svd_dot_train_U_U.append(aaa)
8      else:
9          svd_dot_train_U_U.append(0)
```

```
100%|████████████████████████████████████████████████████████████| 1
00002/100002 [00:07<00:00, 14036.04it/s]
```

In [98]:
```
1  svd_dot_train_U_V = []
2  for i,a in tqdm(df_final_train.iterrows(), total = df_final_train.shape[0]):
3      s = sadj_dict.get(a.source_node, 'NA')
4      d = sadj_dict.get(a.destination_node, 'NA')
5      if (s!='NA') and (d!='NA'):
6          aaa = np.dot(U[s,:], V[:,d])
7          svd_dot_train_U_V.append(aaa)
8      else:
9          svd_dot_train_U_V.append(0)
```

```
100%|████████████████████████████████████████████████████████████| 1
00002/100002 [00:07<00:00, 13601.72it/s]
```

In [99]:
```python
svd_dot_train_V_V = []
for i,a in tqdm(df_final_train.iterrows(), total = df_final_train.shape[0]):
    s = sadj_dict.get(a.source_node, 'NA')
    d = sadj_dict.get(a.destination_node, 'NA')
    if (s!='NA') and (d!='NA'):
        aaa = np.dot(V[:,s], V[:,d])
        svd_dot_train_V_V.append(aaa)
    else:
        svd_dot_train_V_V.append(0)
```

```
100%|██████████████████████████████████████████████████████████| 1
00002/100002 [00:07<00:00, 13043.87it/s]
```

In [100]:
```python
svd_dot_test_U_U = []
for i,a in tqdm(df_final_test.iterrows(), total = df_final_test.shape[0]):
    s = sadj_dict.get(a.source_node, 'NA')
    d = sadj_dict.get(a.destination_node, 'NA')
    if (s!='NA') and (d!='NA'):
        aaa = np.dot(U[s,:], U[d,:])
        svd_dot_test_U_U.append(aaa)
    else:
        svd_dot_test_U_U.append(0)
```

```
100%|██████████████████████████████████████████████████████████|
50002/50002 [00:03<00:00, 13228.33it/s]
```

In [101]:
```python
svd_dot_test_U_V = []
for i,a in tqdm(df_final_test.iterrows(), total = df_final_test.shape[0]):
    s = sadj_dict.get(a.source_node, 'NA')
    d = sadj_dict.get(a.destination_node, 'NA')
    if (s!='NA') and (d!='NA'):
        aaa = np.dot(U[s,:], V[:,d])
        svd_dot_test_U_V.append(aaa)
    else:
        svd_dot_test_U_V.append(0)
```

```
100%|██████████████████████████████████████████████████████████|
50002/50002 [00:03<00:00, 14008.32it/s]
```

In [102]:
```python
svd_dot_test_V_V = []
for i,a in tqdm(df_final_test.iterrows(), total = df_final_test.shape[0]):
    s = sadj_dict.get(a.source_node, 'NA')
    d = sadj_dict.get(a.destination_node, 'NA')
    if (s!='NA') and (d!='NA'):
        aaa = np.dot(V[:,s], V[:,d])
        svd_dot_test_V_V.append(aaa)
    else:
        svd_dot_test_V_V.append(0)
```

```
100%|██████████████████████████████████████████████████████████|
50002/50002 [00:03<00:00, 14118.65it/s]
```

In [106]:
```python
if not os.path.isfile('data/fea_sample/storage_sample_stage5.h5'):
    #=====================================================================

    df_final_train['svd_dot_U_U'] = svd_dot_train_U_U
    df_final_train['svd_dot_U_V'] = svd_dot_train_U_V
    df_final_train['svd_dot_V_V'] = svd_dot_train_V_V
    df_final_test['svd_dot_U_U'] = svd_dot_test_U_U
    df_final_test['svd_dot_U_V'] = svd_dot_test_U_V
    df_final_test['svd_dot_V_V'] = svd_dot_test_V_V
    df_final_train['preferential_attachment_followers'] = df_final_train.app
        lambda row:preferential_attachment_followers(row['source_node'], row
    df_final_train['preferential_attachment_followees'] = df_final_train.app
        lambda row:preferential_attachment_followees(row['source_node'], row
    df_final_test['preferential_attachment_followers'] = df_final_test.apply
        lambda row:preferential_attachment_followers(row['source_node'], row
    df_final_test['preferential_attachment_followees'] = df_final_test.apply
        lambda row:preferential_attachment_followees(row['source_node'], row
    hdf = HDFStore('data/fea_sample/storage_sample_stage5.h5')
    hdf.put('train_df',df_final_train, format='table', data_columns=True)
    hdf.put('test_df',df_final_test, format='table', data_columns=True)
    hdf.close()
```

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]: