

# prabhudayala@gmail.com\_FB\_featurization

November 24, 2019

## Social network Graph Link Prediction - Facebook Challenge

```
[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  
from pandas import HDFStore, DataFrame  
from pandas import read_hdf  
from scipy.sparse.linalg import svds, eigs  
import gc  
from tqdm import tqdm
```

# 1 1. Reading Data

```
[2]: 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',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
      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

## 2 2. Similarity measures

### 2.1 2.1 Jaccard Distance:

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

$$j = \frac{|X \cap Y|}{|X \cup Y|} \quad (1)$$

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

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

0.0

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

0.0

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

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

0

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

0

## 2.2 Cosine distance

$$\text{CosineDistance} = \frac{|X \cap Y|}{|X| \cdot |Y|} \quad (2)$$

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

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

0.0

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

0

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

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

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

0.02886751345948129

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

0

## 2.3 3. Ranking Measures

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

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.

## 2.4 3.1 Page Ranking

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

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

```
[16]: 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
```

```
[17]: #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
```

## 3 4. Other Graph Features

### 3.1 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.

```
[18]: #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
```

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

```
[19]: 10
```

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

```
[20]: -1
```

### 3.2 4.2 Checking for same community

```
[21]: #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
```

```
[22]: belongs_to_same_wcc(861, 1659750)
```

```
[22]: 0
```

```
[23]: belongs_to_same_wcc(669354,1635354)
```

```
[23]: 0
```

### 3.3 4.3 Adamic/Adar Index:

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

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

```
[24]: #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
```

```
[25]: calc_adar_in(1,189226)
```

```
[25]: 0
```

```
[26]: calc_adar_in(669354,1635354)
```

```
[26]: 0
```

### 3.4 4.4 Is person was following back:

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

```
[28]: follows_back(1,189226)
```

```
[28]: 1
```

```
[29]: follows_back(669354,1635354)
```

```
[29]: 0
```

### 3.5 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/> 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

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

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

[31]: 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
```

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

```
0.0007483800935562018
```

### 3.6 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)

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

[34]: 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
```



## 4 5. Featurization

### 4.1 5. 1 Reading a sample of Data from both train and test

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

[36]: 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 lentgh of the file name
    # here we have hardcoded the number of lines as 3775008
    # n_test = sum(1 for line in open(filename)) #number of records in file
    →(excludes 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

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

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

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

Our train matrix size (100002, 3)

```
[38]: source_node destination_node indicator_link
0      273084      1505602             1
1      960433      1371463             1
```

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

Our test matrix size (50002, 3)

```
[39]: source_node destination_node indicator_link
0      848424      784690             1
1     1645357     1526467             1
```

## 4.2 5.2 Adding a set of features

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

jaccard\_followers  
jaccard\_followees  
cosine\_followers  
cosine\_followees  
num\_followers\_s  
num\_followees\_s  
num\_followers\_d  
num\_followees\_d  
inter\_followers  
inter\_followees

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

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

```

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

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

```

```

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

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

        inter_followers.append(len(s1.intersection(d1)))

```

```

        inter_followees.append(len(s2.intersection(d2)))

    return num_followers_s, num_followers_d, num_followees_s, num_followees_d, \
    →inter_followers, inter_followees

```

```

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

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

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

```

### 4.3 5.3 Adding new set of features

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

- adar index
- is following back
- belongs to same weakly connect components
- shortest path between source and destination

```

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

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

```

```

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

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

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

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

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

```

## 4.4 5.4 Adding new set of features

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

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

Page Ranking of source

Page Ranking of dest

katz of source  
 katz of dest  
 hubs of source  
 hubs of dest  
 authorities\_s of source  
 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

$$W = \frac{1}{\sqrt{1 + |X|}} \quad (3)$$

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

```
[44]: #weight for source and destination of each link
Weight_in = {}
Weight_out = {}
for i in tqdm(train_graph.nodes()):
    s1=set(train_graph.predecessors(i))
    w_in = 1.0/(np.sqrt(1+len(s1)))
    Weight_in[i]=w_in

    s2=set(train_graph.successors(i))
    w_out = 1.0/(np.sqrt(1+len(s2)))
    Weight_out[i]=w_out

#for imputing with mean
mean_weight_in = np.mean(list(Weight_in.values()))
mean_weight_out = np.mean(list(Weight_out.values()))
```

100%||  
 1780722/1780722 [00:10<00:00, 168728.52it/s]

```
[45]: if not os.path.isfile('data/fea_sample/storage_sample_stage3.h5'):
    #mapping to pandas train
    df_final_train['weight_in'] = df_final_train.destination_node.apply(lambda x:
    ↳Weight_in.get(x,mean_weight_in))
    df_final_train['weight_out'] = df_final_train.source_node.apply(lambda x:
    ↳Weight_out.get(x,mean_weight_out))

    #mapping to pandas test
    df_final_test['weight_in'] = df_final_test.destination_node.apply(lambda x:
    ↳Weight_in.get(x,mean_weight_in))
```

```

df_final_test['weight_out'] = df_final_test.source_node.apply(lambda x:
↳Weight_out.get(x,mean_weight_out))

#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)

```

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

```

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

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

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

```

```

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

□
→#=====

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

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

□
→#=====

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

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

□
→#=====

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

```



## 4.5 5.5 Adding new set of features

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

SVD features for both source and destination

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

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

[49]: Adj = nx.adjacency_matrix(train_graph,nodelist=sorted(train_graph.nodes())).
      →asfptype()

[50]: 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,)

```
[51]: if not os.path.isfile('data/fea_sample/storage_sample_stage4.h5'):
      □
      →#=====

      df_final_train[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4',
      →'svd_u_s_5', 'svd_u_s_6']] = \
          df_final_train.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)

      df_final_train[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4',
      →'svd_u_d_5', 'svd_u_d_6']] = \
          df_final_train.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)
      □
      →#=====

      df_final_train[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4',
      →'svd_v_s_5', 'svd_v_s_6']] = \
          df_final_train.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)
```

```

df_final_train[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4',
→ 'svd_v_d_5', 'svd_v_d_6']] = \
df_final_train.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.
→Series)

□
→#=====

df_final_test[['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4',
→ 'svd_u_s_5', 'svd_u_s_6']] = \
df_final_test.source_node.apply(lambda x: svd(x, U)).apply(pd.Series)

df_final_test[['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4',
→ 'svd_u_d_5', 'svd_u_d_6']] = \
df_final_test.destination_node.apply(lambda x: svd(x, U)).apply(pd.Series)

□
→#=====

df_final_test[['svd_v_s_1', 'svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4',
→ 'svd_v_s_5', 'svd_v_s_6'],] = \
df_final_test.source_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)

df_final_test[['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4',
→ 'svd_v_d_5', 'svd_v_d_6']] = \
df_final_test.destination_node.apply(lambda x: svd(x, V.T)).apply(pd.Series)

□
→#=====

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

```

[52]: *# prepared and stored the data from machine learning models*  
*# pelase check the FB\_Models.ipynb*

## 4.6 5.6 Adding new set of features

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

SVD DOT features of source and destination

```

[53]: def svd_dot(x, y, S, P):
    try:
        a = sadj_dict[x]
        b = sadj_dict[y]
        return [np.dot(S[a],S[b]),np.dot(P.T[a],P.T[b])]
    except:

```

```
return [0,0]
```

```
[54]: #for svd features to get feature vector creating a dict node val and index in_
      ↪svd vector
```

```
sadj_col = sorted(train_graph.nodes())
sadj_dict = { val:idx for idx,val in enumerate(sadj_col)}
```

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

```
[56]: 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,)
```

```
[57]: if not os.path.isfile('data/fea_sample/storage_sample_stage5.h5'):
      df_final_train[['svd_dot_u','svd_dot_v']] = df_final_train.apply(lambda x:
      ↪svd_dot(x.get('source_node'), x.get('destination_node'),U,V),axis=1).
      ↪apply(pd.Series)
      df_final_test[['svd_dot_u','svd_dot_v']] = df_final_train.apply(lambda x:
      ↪svd_dot(x.get('source_node'), x.get('destination_node'),U,V),axis=1).
      ↪apply(pd.Series)
      hdf = HDFStore('data/fea_sample/storage_sample_stage5.h5')
      hdf.put('train_df',df_final_train, format='table', data_columns=True)
      hdf.put('test_df',df_final_test, format='table', data_columns=True)
      hdf.close()
```

## 4.7 5.7 Adding new set of features

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

Preferential Attachment features of source and destination

```
[58]: #for followees
def preferential_attachment(a,b):
    try:
        w = len(set(train_graph.predecessors(a)))
        x = len(set(train_graph.predecessors(b)))
        y = len(set(train_graph.successors(a)))
        z = len(set(train_graph.successors(b)))
        return [w*x, y*z]
    except:
        return [0, 0]
```

```

[59]: if not os.path.isfile('data/fea_sample/storage_sample_stage6.h5'):
        df_final_train[['preferential_attachment_p', 'preferential_attachment_s']] =
        →df_final_train.apply(lambda row:
                                □
        →preferential_attachment(row['source_node'], row['destination_node']), axis=1).
        →apply(pd.Series)
        df_final_test[['preferential_attachment_p', 'preferential_attachment_s']] =
        →df_final_train.apply(lambda row:
                                □
        →preferential_attachment(row['source_node'], row['destination_node']), axis=1).
        →apply(pd.Series)
        hdf = HDFStore('data/fea_sample/storage_sample_stage6.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()

```