

FB_featurization

September 6, 2019

Social network Graph Link Prediction - Facebook Challenge

```
[265]: #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

```
[266]: 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/>

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

```
[267]: #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
```

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

0.0

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

0.0

```
[270]: #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
```

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

0

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

0

2.2 Cosine distance

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

```
[273]: #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
```

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

0.0

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

0

```
[276]: 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
```

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

0.02886751345948129

```
[278]: 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

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>

```
[279]: 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'))
```

```
[280]: 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
```

```
[281]: #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.

```
[282]: #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
```

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

```
[283]: 10
```

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

```
[284]: -1
```

3.2 4.2 Checking for same community

```
[285]: #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
```

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

```
[286]: 0
```

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

```
[287]: 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.

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

```
[288]: #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
```

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

```
[289]: 0
```

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

```
[290]: 0
```

3.4 4.4 Is person was following back:

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

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

```
[292]: 1
```

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

```
[293]: 0
```

3.5 4.5 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

λ

The parameter

$$\beta$$

controls the initial centrality and

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

```
[294]: 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'))

[295]: 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
```

```
[296]: 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

```
[297]: 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'))

[298]: 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

```
[299]: 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
```

```
[300]: 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
```

```
[301]: 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

```
[302]: 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)

```
[302]: source_node destination_node indicator_link
0      273084      1505602             1
1      842723       524527             1
```

```
[303]: 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)

```
[303]: source_node destination_node indicator_link
0      848424       784690             1
1      539116      1687019             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

```
[304]: 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)

```

```

[305]: 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

```

```

[306]: 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

```

[307]: 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

```
[308]: #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 [00:18<00:00, 98596.66it/s]

```
[309]: 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:
    →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)

```

[310]: `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 Preferential Attachment Feature

Compute the preferential attachment score of all node pairs for train and test data

```
[311]: #Preferential Attachment for followers in Train data
df_final_train['preferential_attachment_followers'] =
    df_final_train['num_followers_s'] * df_final_train['num_followers_d']
df_final_train.head()
```

```
[311]:
```

	source_node	destination_node	indicator_link	jaccard_followers	\
0	273084	1505602	1	0	
1	1492633	1370536	1	0	
2	992126	1128784	1	0	
3	1027527	285795	1	0	
4	663497	827488	1	0	

	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	\
0	0.000000	0.000000	0.000000	11	
1	0.000000	0.000000	0.000000	16	
2	0.000000	0.000000	0.000000	3	
3	0.000000	0.235702	0.000000	8	
4	0.090909	0.096225	0.169031	3	

	num_followers_d	num_followees_s	...	weight_f4	page_rank_s	\
0	6	15	...	0.877964	2.045290e-06	
1	3	6	...	1.255929	2.418525e-06	
2	1	2	...	2.414214	7.204663e-07	
3	9	5	...	1.132724	5.674660e-07	
4	6	5	...	1.194461	3.106595e-07	

	page_rank_d	katz_s	katz_d	hubs_s	hubs_d	\
0	3.459963e-07	0.000773	0.000756	1.943132e-13	1.941103e-13	
1	2.817230e-07	0.000791	0.000744	2.561722e-17	1.200162e-17	
2	4.042322e-07	0.000743	0.000735	9.362040e-14	7.161209e-17	
3	6.759004e-07	0.000762	0.000766	4.847751e-220	3.974254e-220	
4	4.604112e-07	0.000743	0.000755	8.934659e-17	9.554532e-14	

	authorities_s	authorities_d	preferential_attachment_followers
0	9.226339e-16	2.231877e-15	66
1	7.148065e-16	8.256346e-16	48
2	9.977100e-15	6.583609e-16	3
3	1.392266e-218	1.374740e-218	72
4	3.229567e-17	1.750384e-15	18

[5 rows x 32 columns]

```
[312]: #Preferential Attachment for followees in Train data
df_final_train['preferential_attachment_followees'] =
    df_final_train['num_followees_s'] * df_final_train['num_followees_d']
```

```
df_final_train.head()
```

```
[312]:
```

	source_node	destination_node	indicator_link	jaccard_followers	\
0	273084	1505602	1	0	
1	1492633	1370536	1	0	
2	992126	1128784	1	0	
3	1027527	285795	1	0	
4	663497	827488	1	0	

	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	\
0	0.000000	0.000000	0.000000	11	
1	0.000000	0.000000	0.000000	16	
2	0.000000	0.000000	0.000000	3	
3	0.000000	0.235702	0.000000	8	
4	0.090909	0.096225	0.169031	3	

	num_followers_d	num_followees_s	...	page_rank_s	page_rank_d	\
0	6	15	...	2.045290e-06	3.459963e-07	
1	3	6	...	2.418525e-06	2.817230e-07	
2	1	2	...	7.204663e-07	4.042322e-07	
3	9	5	...	5.674660e-07	6.759004e-07	
4	6	5	...	3.106595e-07	4.604112e-07	

	katz_s	katz_d	hubs_s	hubs_d	authorities_s	\
0	0.000773	0.000756	1.943132e-13	1.941103e-13	9.226339e-16	
1	0.000791	0.000744	2.561722e-17	1.200162e-17	7.148065e-16	
2	0.000743	0.000735	9.362040e-14	7.161209e-17	9.977100e-15	
3	0.000762	0.000766	4.847751e-220	3.974254e-220	1.392266e-218	
4	0.000743	0.000755	8.934659e-17	9.554532e-14	3.229567e-17	

	authorities_d	preferential_attachment_followers	\
0	2.231877e-15	66	
1	8.256346e-16	48	
2	6.583609e-16	3	
3	1.374740e-218	72	
4	1.750384e-15	18	

	preferential_attachment_followees
0	120
1	6
2	2
3	25
4	35

[5 rows x 33 columns]

```
[313]:
```

```
#Preferential Attachment for Train data
df_final_test['preferential_attachment_followers'] =
    df_final_test['num_followers_s'] * df_final_test['num_followers_d']
df_final_test.head()
```

```
[313]:
```

	source_node	destination_node	indicator_link	jaccard_followers	\
0	848424	784690	1	0	
1	1591960	1350641	1	0	
2	1663641	58705	1	0	
3	667970	1000475	1	0	
4	1272737	988441	1	0	

	jaccard_followees	cosine_followers	cosine_followees	num_followers_s	\
0	0.000000	0.029161	0.000000	6	
1	0.000000	0.017495	0.000000	3	
2	0.117647	0.026218	0.218218	11	
3	0.147059	0.083333	0.304290	9	
4	0.000000	0.000000	0.000000	16	

	num_followers_d	num_followees_s	...	weight_f4	page_rank_s	\
0	14	6	...	1.014128	6.557971e-07	
1	33	2	...	1.326199	4.334026e-07	
2	23	14	...	0.720522	8.375801e-07	
3	20	9	...	0.850673	4.693406e-07	
4	0	20	...	1.025406	9.567645e-07	

	page_rank_d	katz_s	katz_d	hubs_s	hubs_d	\
0	1.559547e-06	0.000754	0.000786	3.243237e-16	1.745627e-16	
1	3.985817e-06	0.000743	0.000859	1.945488e-16	1.905493e-14	
2	3.089958e-06	0.000776	0.000821	5.421010e-13	1.205587e-12	
3	1.803976e-06	0.000767	0.000808	5.105743e-17	1.120746e-15	
4	5.615700e-07	0.000796	0.000748	8.427032e-17	0.000000e+00	

	authorities_s	authorities_d	preferential_attachment_followers
0	2.969838e-15	9.269213e-14	84
1	4.404545e-15	8.468971e-15	99
2	5.114413e-13	5.554424e-12	253
3	2.813467e-17	1.656798e-15	180
4	7.052574e-17	0.000000e+00	0

[5 rows x 32 columns]

```
[314]: #Preferential Attachment for Train data
df_final_test['preferential_attachment_followees'] =
    df_final_test['num_followees_s'] * df_final_test['num_followees_d']
df_final_test.head()
```

```

[314]: source_node destination_node indicator_link jaccard_followers \
0      848424      784690      1      0
1      1591960      1350641      1      0
2      1663641      58705      1      0
3      667970      1000475      1      0
4      1272737      988441      1      0

      jaccard_followees cosine_followers cosine_followees num_followers_s \
0      0.000000      0.029161      0.000000      6
1      0.000000      0.017495      0.000000      3
2      0.117647      0.026218      0.218218      11
3      0.147059      0.083333      0.304290      9
4      0.000000      0.000000      0.000000      16

      num_followers_d num_followees_s ... page_rank_s page_rank_d \
0      14      6 ... 6.557971e-07 1.559547e-06
1      33      2 ... 4.334026e-07 3.985817e-06
2      23      14 ... 8.375801e-07 3.089958e-06
3      20      9 ... 4.693406e-07 1.803976e-06
4      0      20 ... 9.567645e-07 5.615700e-07

      katz_s katz_d hubs_s hubs_d authorities_s \
0 0.000754 0.000786 3.243237e-16 1.745627e-16 2.969838e-15
1 0.000743 0.000859 1.945488e-16 1.905493e-14 4.404545e-15
2 0.000776 0.000821 5.421010e-13 1.205587e-12 5.114413e-13
3 0.000767 0.000808 5.105743e-17 1.120746e-15 2.813467e-17
4 0.000796 0.000748 8.427032e-17 0.000000e+00 7.052574e-17

      authorities_d preferential_attachment_followers \
0 9.269213e-14      84
1 8.468971e-15      99
2 5.554424e-12      253
3 1.656798e-15      180
4 0.000000e+00      0

      preferential_attachment_followees
0      54
1      96
2      336
3      270
4      0

```

[5 rows x 33 columns]

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 features for both source and destination

```
[315]: def svd(x, S):  
        try:  
            z = sadj_dict[x]  
            return S[z]  
        except:  
            return [0,0,0,0,0,0]  
  
[316]: #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)}  
  
[317]: Adj = nx.adjacency_matrix(train_graph,nodelist=sorted(train_graph.nodes())).  
        →asfptype()  
  
[318]: 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,)
```

```
[319]: # 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)

#=====

```

```

[320]: svd_dot_u_train = df_final_train['svd_u_s_1'].
→dot(df_final_train['svd_u_d_1'])+df_final_train['svd_u_s_2'].
→dot(df_final_train['svd_u_d_2'])+\
df_final_train['svd_u_s_3'].
→dot(df_final_train['svd_u_d_3'])+df_final_train['svd_u_s_4'].
→dot(df_final_train['svd_u_d_4'])+\
df_final_train['svd_u_s_5'].
→dot(df_final_train['svd_u_d_5'])+df_final_train['svd_u_s_6'].
→dot(df_final_train['svd_u_d_6'])

df_final_train['svd_dot_u'] = svd_dot_u_train

```

```

[321]: svd_dot_v_train = df_final_train['svd_v_s_1'].
→dot(df_final_train['svd_v_d_1'])+df_final_train['svd_v_s_2'].
→dot(df_final_train['svd_v_d_2'])+\
df_final_train['svd_v_s_3'].
→dot(df_final_train['svd_v_d_3'])+df_final_train['svd_v_s_4'].
→dot(df_final_train['svd_v_d_4'])+\
df_final_train['svd_v_s_5'].
→dot(df_final_train['svd_v_d_5'])+df_final_train['svd_v_s_6'].
→dot(df_final_train['svd_v_d_6'])

df_final_train['svd_dot_v'] = svd_dot_v_train

```

```

[322]: svd_dot_u_test = df_final_test['svd_u_s_1'].
→dot(df_final_test['svd_u_d_1'])+df_final_test['svd_u_s_2'].
→dot(df_final_test['svd_u_d_2'])+\

```

```

df_final_test['svd_u_s_3'].
    ↳dot(df_final_test['svd_u_d_3'])+df_final_test['svd_u_s_4'].
    ↳dot(df_final_test['svd_u_d_4'])+\
df_final_test['svd_u_s_5'].
    ↳dot(df_final_test['svd_u_d_5'])+df_final_test['svd_u_s_6'].
    ↳dot(df_final_test['svd_u_d_6'])

df_final_test['svd_dot_u'] = svd_dot_u_test

```

```

[323]: svd_dot_v_test = df_final_test['svd_v_s_1'].
    ↳dot(df_final_test['svd_v_d_1'])+df_final_test['svd_v_s_2'].
    ↳dot(df_final_test['svd_v_d_2'])+\
df_final_test['svd_v_s_3'].
    ↳dot(df_final_test['svd_v_d_3'])+df_final_test['svd_v_s_4'].
    ↳dot(df_final_test['svd_v_d_4'])+\
df_final_test['svd_v_s_5'].
    ↳dot(df_final_test['svd_v_d_5'])+df_final_test['svd_v_s_6'].
    ↳dot(df_final_test['svd_v_d_6'])

df_final_test['svd_dot_v'] = svd_dot_v_test

```

```

[324]: 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()

```

```

[ ]:

```