

Social network Graph Link Prediction - Facebook Challenge

Problem statement:

Given a directed social graph, have to predict missing links to recommend users (Link Prediction in graph)

Data Overview

Taken data from facebook's recruiting challenge on kaggle <https://www.kaggle.com/c/FacebookRecruiting> data contains two columns source and destination eac edge in graph - Data columns (total 2 columns):

- source_node int64
- destination_node int64

Mapping the problem into supervised learning problem:

- Generated training samples of good and bad links from given directed graph and for each link got some features like the followed back, page rank, katz score, adar index, some svd features of adj matrix, some weight features etc. based on these features to predict link.

Saved successfully!



[ne/kleinber/link-pred.pdf](#)

- <https://www3.nd.edu/~dai/publications/lichtenwalter2010new.pdf>
- https://kaggle2.blob.core.windows.net/forum-message-attachments/2594/supervised_link_prediction
- <https://www.youtube.com/watch?v=2M77Hgy17cg>

Business objectives and constraints:

- No low-latency requirement.
- Probability of prediction is useful to recommend highest probability links

▼ Performance metric for supervised learning:

- Both precision and recall is important so F1 score is good choice
- Confusion matrix

```
from google.colab import drive
drive.mount('/content/drive/')
```

🔗 Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=9473189

Enter your authorization code:
.....

Mounted at /content/drive/

```
cd /content/drive/My Drive/training/Facebook/
```

📁 /content/drive/My Drive/training/Facebook

```
import warnings
warnings.filterwarnings("ignore")

import csv
import pandas as pd
import datetime
import time

import numpy as np

import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib import rcParams
from sklearn.cluster import MiniBatchKMeans, KMeans
import math
import pickle
import os
```

Saved successfully!

```
import warnings
import networkx as nx
import pdb
import pickle

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

📁

Displaying a sub graph

```

if not os.path.isfile('train_woheader_sample.csv'):
    pd.read_csv('data/train.csv', nrows=100).to_csv('train_woheader_sample.csv', header=False,

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

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


```

▼ 1. Exploratory Data Analysis

```

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

```

 The number of unique persons 1862220

▼ 1.1 No of followers for each person

Saved successfully! 

```

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

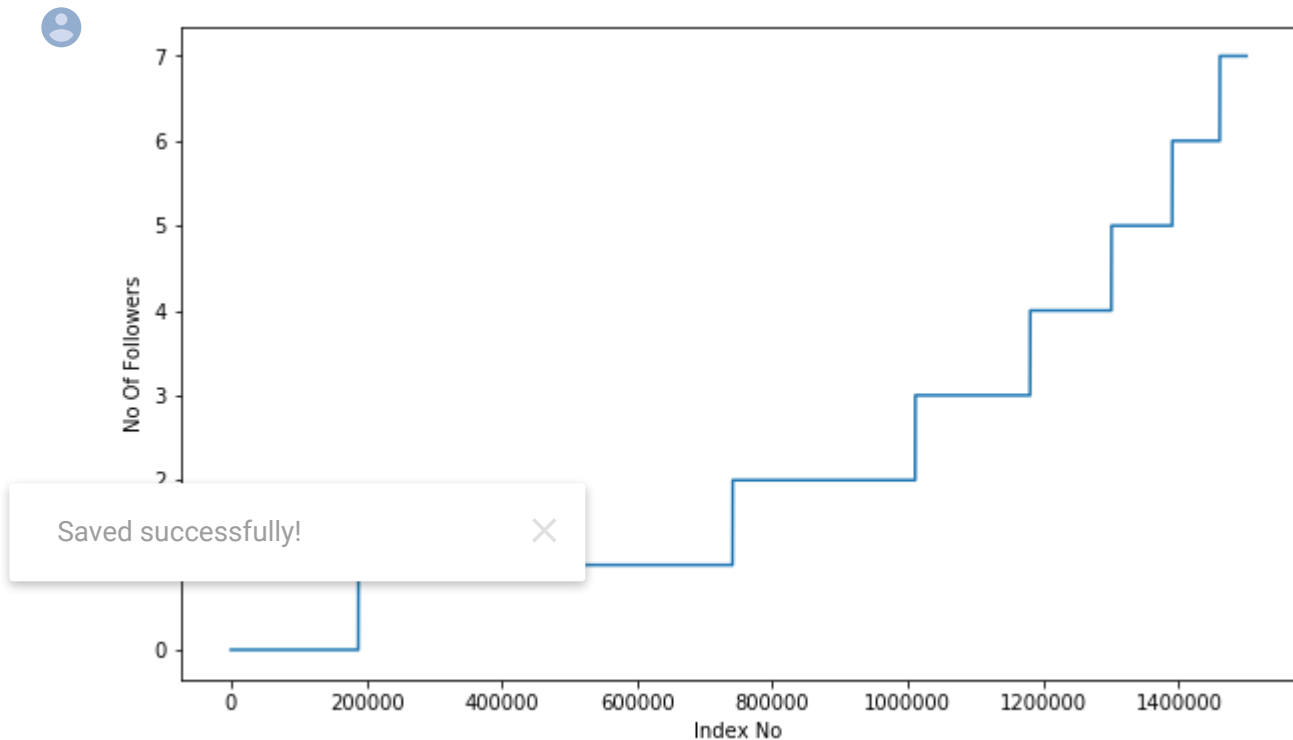
```



```

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

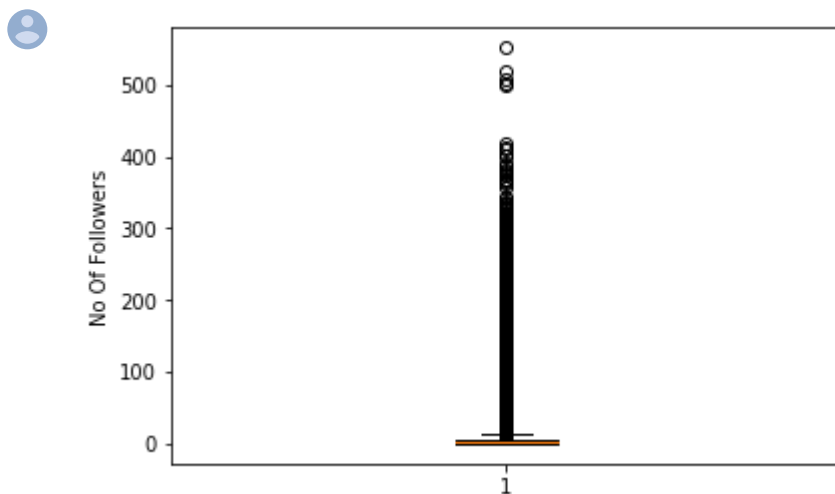
```



```

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

```



```

### 90-100 percentile
for i in range(0,11):

```

```
print(90+i,'percentile value is',np.percentile(indegree_dist,90+i))
```



```
90 percentile value is 12.0
91 percentile value is 13.0
92 percentile value is 14.0
93 percentile value is 15.0
94 percentile value is 17.0
95 percentile value is 19.0
96 percentile value is 21.0
97 percentile value is 24.0
98 percentile value is 29.0
99 percentile value is 40.0
100 percentile value is 552.0
```

99% of data having followers of 40 only.

```
### 99-100 percentile
```

```
for i in range(10,110,10):
```

```
    print(99+(i/100),'percentile value is',np.percentile(indegree_dist,99+(i/100)))
```



```
99.1 percentile value is 42.0
99.2 percentile value is 44.0
99.3 percentile value is 47.0
```

Saved successfully!



```
99.6 percentile value is 61.0
99.7 percentile value is 70.0
99.8 percentile value is 84.0
99.9 percentile value is 112.0
100.0 percentile value is 552.0
```

```
%matplotlib inline
```

```
sns.set_style('ticks')
```

```
fig, ax = plt.subplots()
```

```
fig.set_size_inches(11.7, 8.27)
```

```
sns.distplot(indegree_dist, color='#16A085')
```

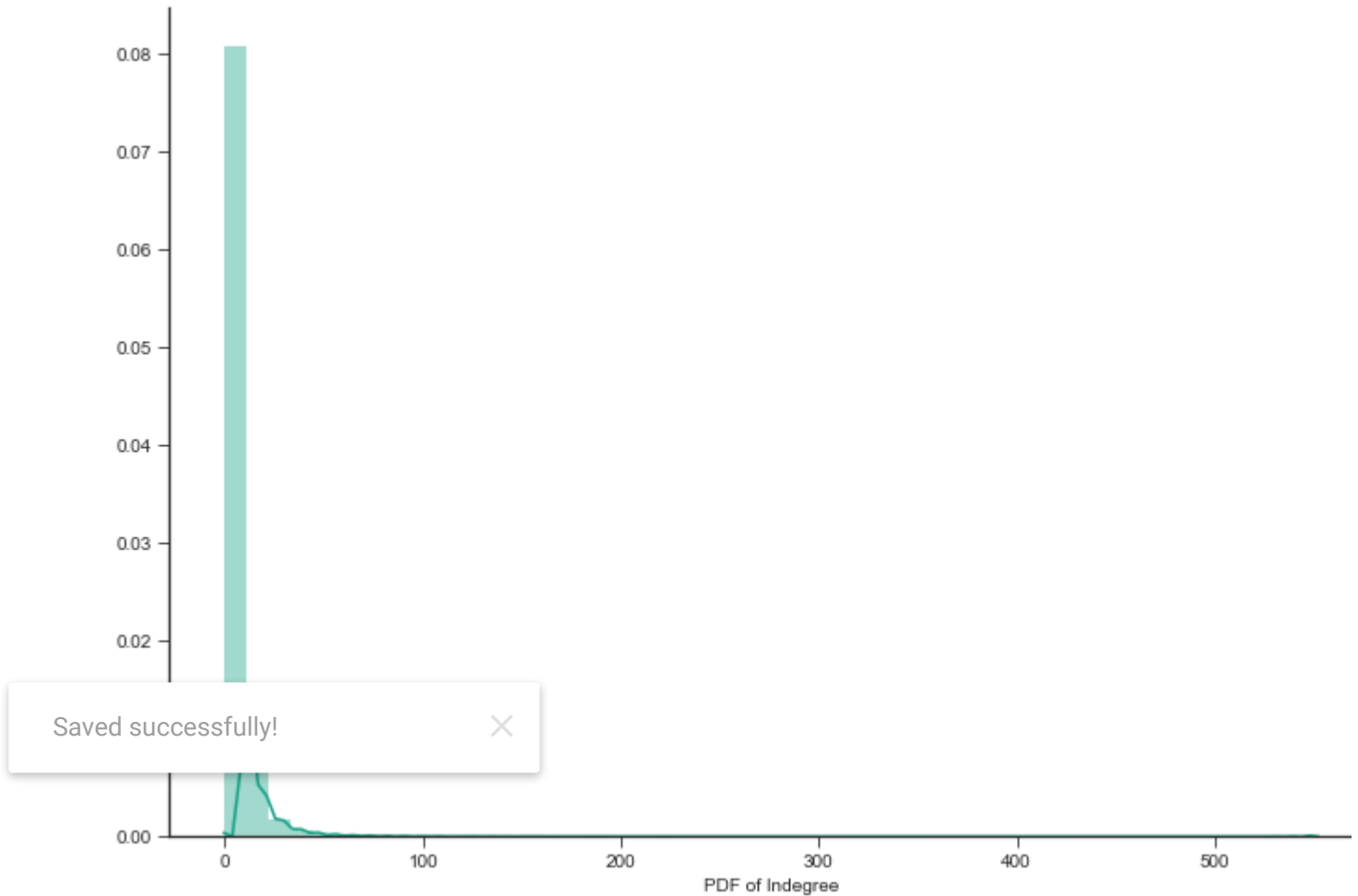
```
plt.xlabel('PDF of Indegree')
```

```
sns.despine()
```

```
#plt.show()
```



D:\installed\Anaconda3\lib\site-packages\matplotlib\axes_axes.py:6571: UserWarning: The warnings.warn("The 'normed' kwarg is deprecated, and has been "



► 1.2 No of people each person is following

↳ 9 cells hidden

► 1.3 both followers + following

↳ 9 cells hidden

▼ 2. Posing a problem as classification problem

► 2.1 Generating some edges which are not present in graph for supervised le

Generated Bad links from graph which are not in graph and whose shortest path is greater than 2.

↳ 2 cells hidden

▶ 2.2 Training and Test data split:

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

↳ 7 cells hidden

▼ 1. Reading Data

```
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.Graph)
    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

Saved successfully!



▼ 2. Similarity measures

▼ 2.1 Jaccard Distance:

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

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

```
#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
```

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

 0.0

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

 0.0

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

print(jaccard_for_followers(273084,470294))
```

Saved successfully! 

```
#node 1635354 not in graph
print(jaccard_for_followees(669354,1635354))
```

 0

▼ 2.2 Cosine distance

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

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

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





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

 0

```
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)))) * math.sqrt(len(set(train_graph.predecessors(b)))))
        return sim
    except:
        return 0
```

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

 0.02886751345948129

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

Saved successfully!

▼ 3. Ranking Measures

https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx.algorithms.link_analysis.pagerank_numpy.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 A has 5 incoming links; the probability of following one of those links is 20%, so Page A has a PageRank of 0.2. Page B has no incoming links, so its PageRank is 0.0. Page C has two incoming links; the probability of following one of those links is 25%, so Page C has a PageRank of 0.25. Page D has one incoming link, so its PageRank is 0.1. Page E has three incoming links; the probability of following one of those links is 20%, so Page E has a PageRank of 0.3. (The 15% likelihood of jumping to an arbitrary page corresponds to a damping factor of 0.85.) Without damping, all web surfers would eventually end up on the same page, so all other pages would have PageRank zero. In the presence of damping, Page A effectively links to all pages in the network, so all pages have a PageRank greater than zero. The damping factor is a value between 0 and 1, representing the probability of following a link versus jumping to an arbitrary page.

▼ 3.1 Page Ranking

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

```
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'))
```

```
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
```

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

```
5.615699699389075e-07
```

4. Other Graph Features

4.1 Shortest path:

Saved successfully!



If nodes have direct path i.e directly connected then we are removing the path.

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

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

```
10
```

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

```
-1
```

▼ 4.2 Checking for same community

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

Saved successfully!

```
else:
    for i in wcc:
        if a in i:
            index= i
            break
    if(b in index):
        return 1
    else:
        return 0
```

belongs_to_same_wcc(861, 1659750)

 0

belongs_to_same_wcc(669354,1635354)

 0

▼ 4.3 Adamic/Adar Index:

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

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

```
#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
```

```
calc_adar_in(1,189226)
```

 0

```
calc_adar_in(669354,1635354)
```

 0

Saved successfully!



back:

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

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

 1

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

 0

▼ 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 and 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

 β

controls the initial centrality and

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

```
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'))

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
```

Saved successfully!

```
mean_katz = float(sum(katz.values())) / len(katz)
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 node value based on outgoing links.

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

```
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'))

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
```

5. Featurization

5.1 Reading a sample of Data from both train and test

```
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

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
    # 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

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))

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

Saved successfully!




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



Our train matrix size (100002, 3)

	source_node	destination_node	indicator_link
0	273084	1505602	1
1	832016	1543415	1

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

 Our test matrix size (50002, 3)

	source_node	destination_node	indicator_link
0	848424	784690	1
1	483294	1255532	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

Saved successfully!

8. num_followees_d
9. inter_followers
10. inter_followees

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

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

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

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

```

df_final_train['cosine_followees'] = df_final_train.apply(lambda row:
    cosine_for_followees(row['source_node'],row['dest
df_final_test['cosine_followees'] = df_final_test.apply(lambda row:
    cosine_for_followees(row['source_node'],row['dest

```

```

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_follower

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_st

    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_stag

    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)

```

Saved successfully!




```

hdf.close()
else:
    df_final_train = pd.read_hdf('data/fea_sample/storage_sample_stage1.h5', 'train_df', mode='r')
    df_final_test = pd.read_hdf('data/fea_sample/storage_sample_stage1.h5', 'test_df', mode='r')

```

▼ 5.3 Adding new set of features

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

1. adar index
2. is following back
3. belongs to same weakly connect components
4. shortest path between source and destination

```

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_no'], row['target_no']), axis=1)
    #mapping adar index on test
    df_final_test['adar_index'] = df_final_test.apply(lambda row: calc_adar_in(row['source_no'], row['target_no']), axis=1)

    #mapping follows back or not on train
    df_final_train['follows_back'] = df_final_train.apply(lambda row: follows_back(row['source_no'], row['target_no']), axis=1)

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

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

    #-----
    #mapping shortest path on train
    df_final_train['shortest_path'] = df_final_train.apply(lambda row: compute_shortest_path_length(row['source_no'], row['target_no']), axis=1)
    #mapping shortest path on test
    df_final_test['shortest_path'] = df_final_test.apply(lambda row: compute_shortest_path_length(row['source_no'], row['target_no']), 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 = pd.read_hdf('data/fea_sample/storage_sample_stage2.h5', 'train_df', mode='r')
    df_final_test = pd.read_hdf('data/fea_sample/storage_sample_stage2.h5', 'test_df', mode='r')

```

Saved successfully!



ain

df_final_train.apply(lambda row: follows_back(row['source_no'], row['target_no']), axis=1)

▼ 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

Saved successfully!



▼ Weight Features

In order to determine the similarity of nodes, an edge weight value was calculated between nodes. Edge weight decr count goes up. Intuitively, consider one million people following a celebrity on a social network then chances are mc other or the celebrity. On the other hand, if a user has 30 contacts in his/her social network, the chances are higher t each other. credit - Graph-based Features for Supervised Link Prediction William Cukierski, Benjamin Hamner, Bo Y

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

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

```
#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()))

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))
    df_final_train['weight_out'] = df_final_train.source_node.apply(lambda x: Weight_out.get(x))

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

    #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'] = df_final_test.weight_in + 1*df_final_test.weight_out
    df_final_test['weight_f4'] = df_final_test.weight_in + 2*df_final_test.weight_out

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))

```

Saved successfully!



```

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(
df_final_train['authorities_d'] = df_final_train.destination_node.apply(lambda x: hits[1].get(

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

```

5.5 Adding new set of features

Saved successfully!



for both train and test data points

1. SVD features for both source and destination

```

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

#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)}

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

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,)
```

▼ Adding new feature Preferential attachment

```
def followee_attachment(u1,u2):
    try:
        u_1 = len(set(train_graph.successors(u1)))
        u_2 = len(set(train_graph.successors(u2)))
        return(u_1*u_2)
    except:
        return 0

def follower_attachment(user1,user2):
    try:
        u_1 = len(set(train_graph.predecessors(u1)))
        u_2 = len(set(train_graph.predecessors(u2)))
```

Saved successfully!

```
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()
```

```
# prepared and stored the data from machine learning models
# pelase check the FB_Models.ipynb
```

```
df_final_train.columns
```

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

Saved successfully!

▼ Dot product between svd features

```
s1, s2, s3, s4, s5, s6 = df_final_train['svd_v_s_1'], df_final_train['svd_v_s_2'], df_final_train['svd_v_s_3'], df_final_train['svd_v_s_4'], df_final_train['svd_v_s_5'], df_final_train['svd_v_s_6']
s7, s8, s9, s10, s11, s12 = df_final_train['svd_u_s_1'], df_final_train['svd_u_s_2'], df_final_train['svd_u_s_3'], df_final_train['svd_u_s_4'], df_final_train['svd_u_s_5'], df_final_train['svd_u_s_6']
```

```
d1, d2, d3, d4, d5, d6 = df_final_train['svd_v_d_1'], df_final_train['svd_v_d_2'], df_final_train['svd_v_d_3'], df_final_train['svd_v_d_4'], df_final_train['svd_v_d_5'], df_final_train['svd_v_d_6']
d7, d8, d9, d10, d11, d12 = df_final_train['svd_u_d_1'], df_final_train['svd_u_d_2'], df_final_train['svd_u_d_3'], df_final_train['svd_u_d_4'], df_final_train['svd_u_d_5'], df_final_train['svd_u_d_6']
```

```
dot_svd_fea = []
for i in range(len(np.array(s1))):
    a1, a2=[], []
    a1.append(np.array(s1[i]))
    a1.append(np.array(s2[i]))
    a1.append(np.array(s3[i]))
    a1.append(np.array(s4[i]))
    a1.append(np.array(s5[i]))
    a1.append(np.array(s6[i]))
    a1.append(np.array(s7[i]))
    a1.append(np.array(s8[i]))
```

```

a1.append(np.array(s8[i]))
a1.append(np.array(s9[i]))
a1.append(np.array(s10[i]))
a1.append(np.array(s11[i]))
a1.append(np.array(s12[i]))
a2.append(np.array(d1[i]))
a2.append(np.array(d2[i]))
a2.append(np.array(d3[i]))
a2.append(np.array(d4[i]))
a2.append(np.array(d5[i]))
a2.append(np.array(d6[i]))
a2.append(np.array(d7[i]))
a2.append(np.array(d8[i]))
a2.append(np.array(d9[i]))
a2.append(np.array(d10[i]))
a2.append(np.array(d11[i]))
a2.append(np.array(d12[i]))
dot_svd_fea.append(np.dot(a1,a2))
df_final_train['dot_svd'] = dot_svd_fea

```

```

s1, s2, s3, s4, s5, s6 = df_final_test['svd_v_s_1'], df_final_test['svd_v_s_2'], df_final_test['svd_v_s_3'], df_final_test['svd_v_s_4'], df_final_test['svd_v_s_5'], df_final_test['svd_v_s_6']
s7, s8, s9, s10, s11, s12 = df_final_test['svd_u_s_1'], df_final_test['svd_u_s_2'], df_final_test['svd_u_s_3'], df_final_test['svd_u_s_4'], df_final_test['svd_u_s_5'], df_final_test['svd_u_s_6']

```

Saved successfully!



```

df_final_test['svd_v_d_1'], df_final_test['svd_v_d_2'], df_final_test['svd_v_d_3'], df_final_test['svd_v_d_4'], df_final_test['svd_v_d_5'], df_final_test['svd_v_d_6'], df_final_test['svd_u_d_1'], df_final_test['svd_u_d_2'], df_final_test['svd_u_d_3'], df_final_test['svd_u_d_4'], df_final_test['svd_u_d_5'], df_final_test['svd_u_d_6']

```

```

dot_svd_fea_t = []
for i in range(len(np.array(s1))):
    a1, a2=[], []
    a1.append(np.array(s1[i]))
    a1.append(np.array(s2[i]))
    a1.append(np.array(s3[i]))
    a1.append(np.array(s4[i]))
    a1.append(np.array(s5[i]))
    a1.append(np.array(s6[i]))
    a1.append(np.array(s7[i]))
    a1.append(np.array(s8[i]))
    a1.append(np.array(s9[i]))
    a1.append(np.array(s10[i]))
    a1.append(np.array(s11[i]))
    a1.append(np.array(s12[i]))
    a2.append(np.array(d1[i]))
    a2.append(np.array(d2[i]))
    a2.append(np.array(d3[i]))
    a2.append(np.array(d4[i]))
    a2.append(np.array(d5[i]))
    a2.append(np.array(d6[i]))
    a2.append(np.array(d7[i]))
    a2.append(np.array(d8[i]))
    a2.append(np.array(d9[i]))
    a2.append(np.array(d10[i]))

```

```

a2.append(np.array(d11[i]))
a2.append(np.array(d12[i]))
dot_svd_fea_t.append(np.dot(a1,a2))
df_final_test['dot_svd'] = dot_svd_fea_t

```

```

df_final_train['followee_attachment'] = df_final_train.apply(lambda x: followee_attachment(x[
df_final_test['followee_attachment'] = df_final_test.apply(lambda x: followee_attachment(x['s

```

```

df_final_train['follower_attachment'] = df_final_train.apply(lambda x: follower_attachment(x[
df_final_test['follower_attachment'] = df_final_test.apply(lambda x: follower_attachment(x['s

```

```

#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 sys #System information on arrays

```

Saved successfully!

```

import matplotlib.pyplot as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots
from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
import math
import pickle
import os
# to install xgboost: pip3 install xgboost
import xgboost as xgb

```

```

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

```

```

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

```



```
df_final_train.columns
```

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

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

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

Saved successfully!

```
from sklearn.model_selection import RandomizedSearchCV
clf = xgb.XGBClassifier()
param_dist = {"n_estimators": [100, 150, 200, 250, 300],
              "max_depth": [2, 4, 6],
              "learning_rate": [0.1, 0.2, 0.3],
              "min_child_weight": [2, 4, 6] }
model = RandomizedSearchCV(clf, param_dist, n_iter=5, cv=5, scoring='f1', random_state=42)
```

```
model.fit(df_final_train, y_train)
print(model.best_estimator_)
```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=1, gamma=0,
              learning_rate=0.2, max_delta_step=0, max_depth=4,
              min_child_weight=4, missing=None, n_estimators=200, n_jobs=1,
              nthread=None, objective='binary:logistic', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=None, subsample=1, verbosity=1)
```

```
clf=xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bytree=1, gamma=0, learning_rate=0.2, max_delta_step=0,
                      max_depth=4, min_child_weight=4, missing=None, n_estimators=200,
                      n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
                      reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
```

```

silent=True, subsample=1, verbosity=1)
clf.fit(df_final_train, y_train)

```

```

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
               colsample_bynode=1, colsample_bytree=1, gamma=0,
               learning_rate=0.2, max_delta_step=0, max_depth=4,
               min_child_weight=4, missing=None, n_estimators=200, n_jobs=1,
               nthread=None, objective='binary:logistic', random_state=0,
               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
               silent=True, subsample=1, verbosity=1)

```

```

y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)

```

```

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

```

```

Train f1 score 0.986434342216327
Test f1 score 0.9147326800756403

```

```

from sklearn.metrics import confusion_matrix
confusion_matrix(y_test,y_test_pred):
confusion_matrix(y_test,y_test_pred)

```

Saved successfully!

```

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

```

```

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

```

```

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

```

```

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

```

```

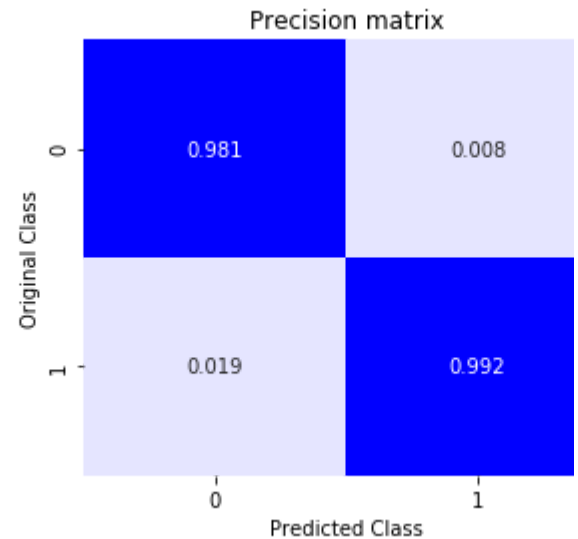
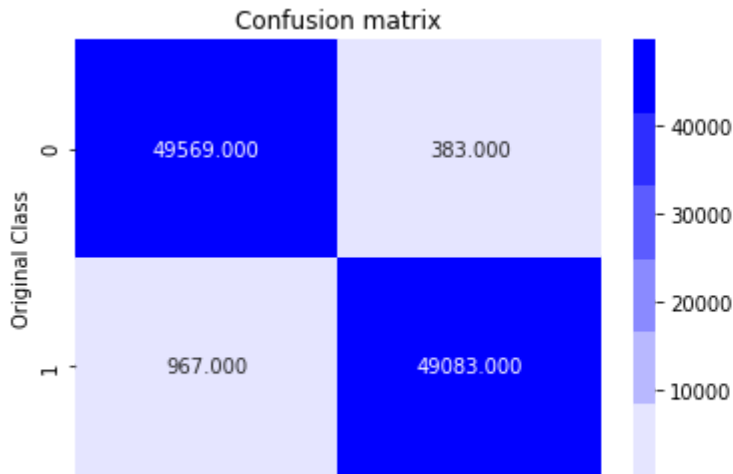
plt.subplot(1, 3, 3)
# representing B in heatmap format
sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.title("Recall matrix")

```

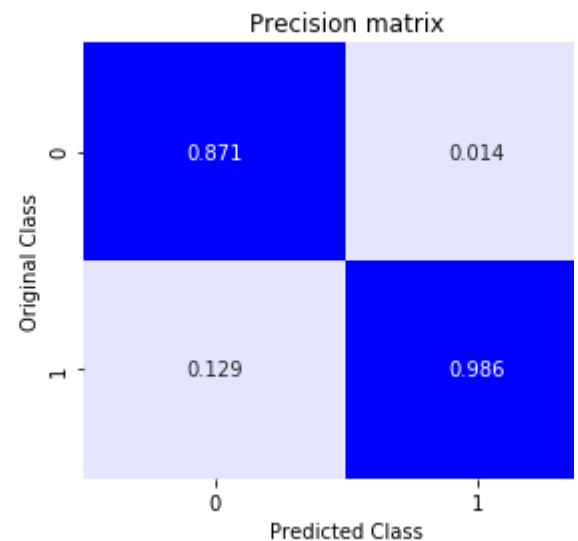
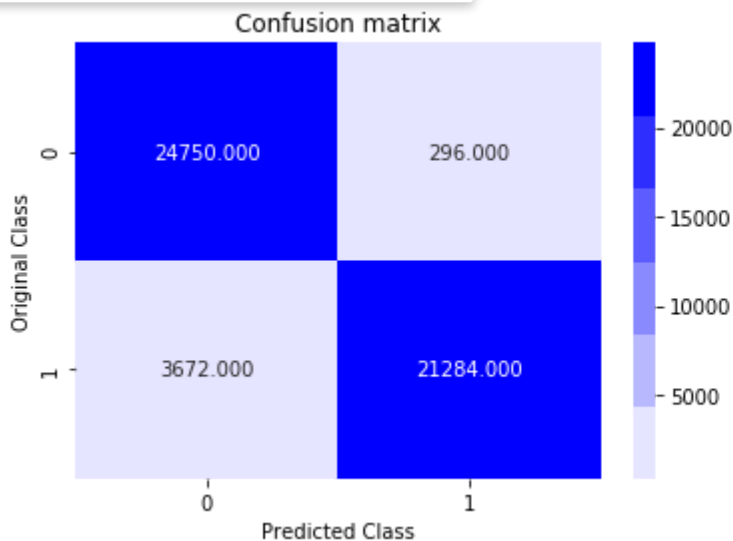
```
plt.show()
```

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

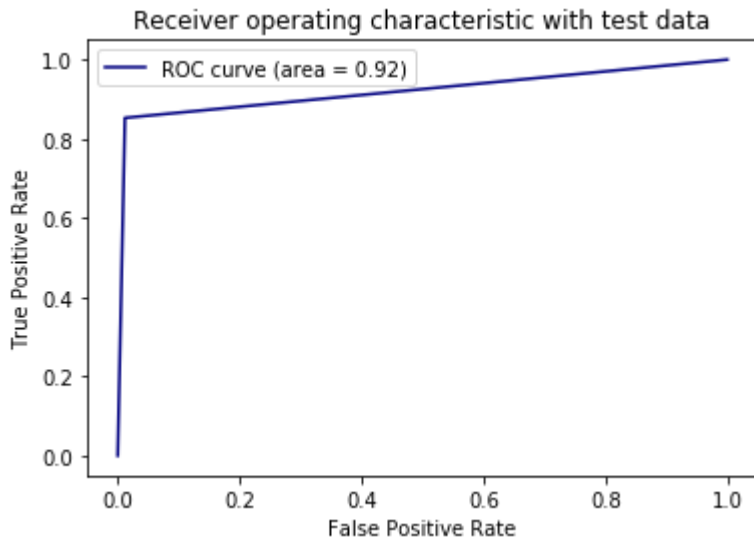
☞ Train confusion_matrix



Saved successfully!



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



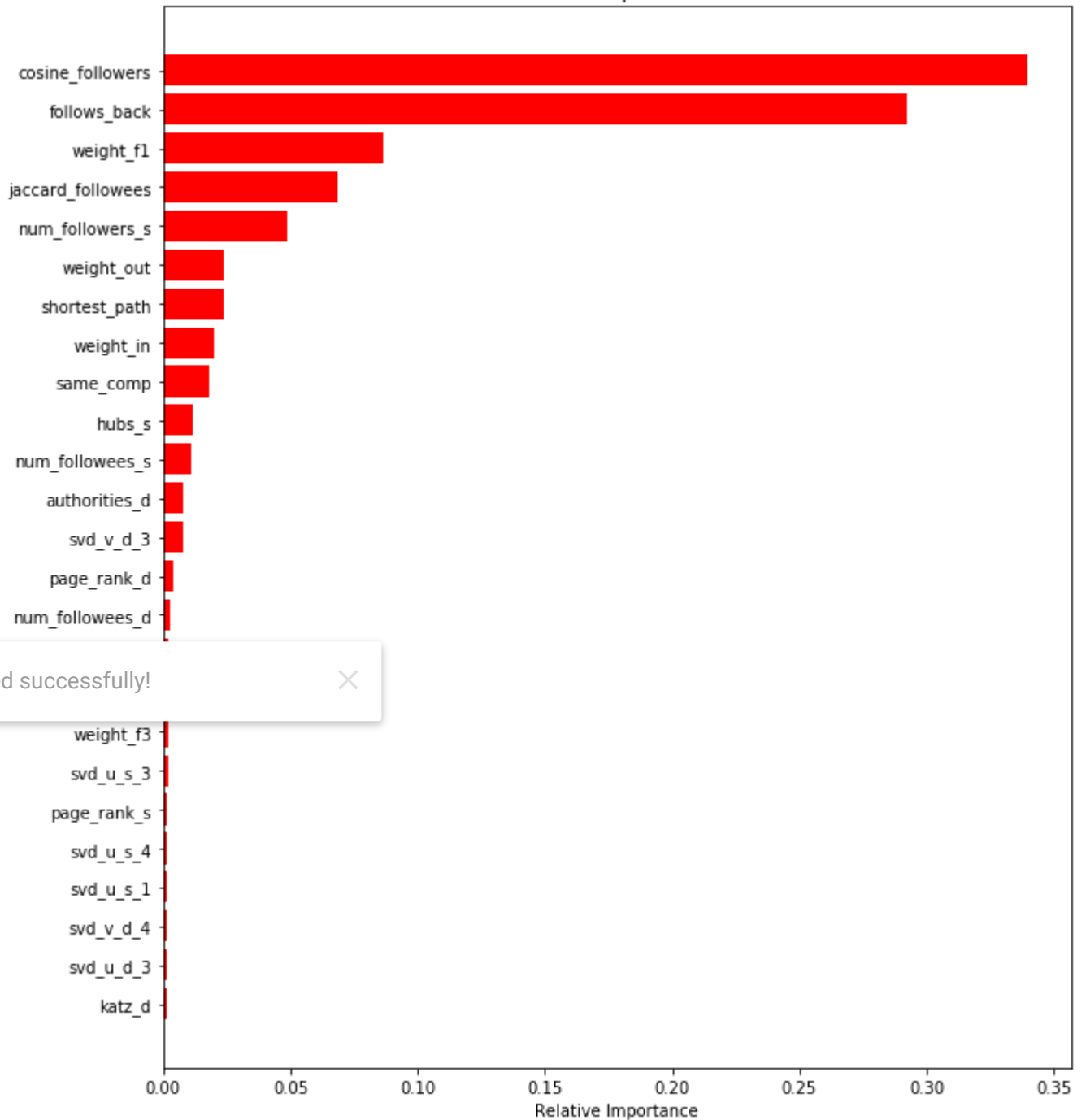
```
features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-25:]
plt.figure(figsize=(10,12))
plt.title('Feature Importances')
plt.xlabel('Relative Importance')
plt.show()
```

Saved successfully!

ances[indices], color='r', align='center')
tures[i] for i in indices])



Feature Importances



▼ Using Random forest

```

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

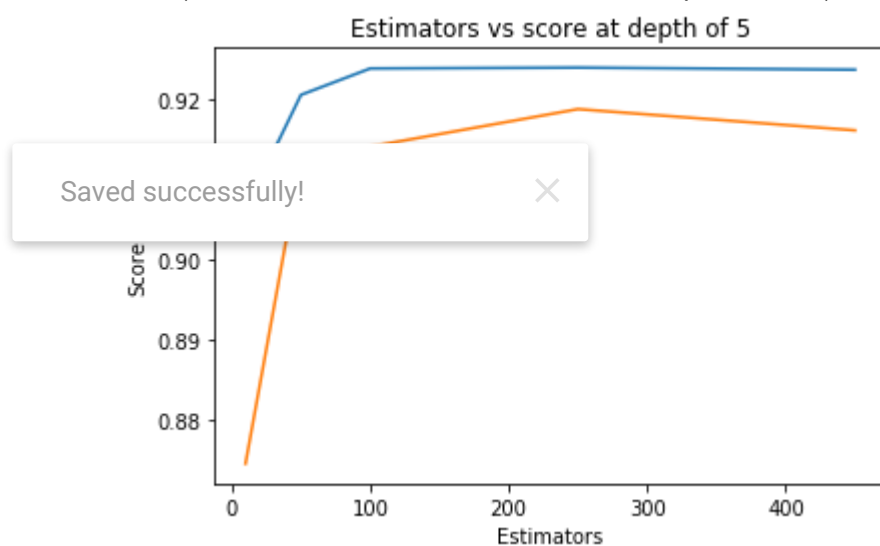
```

```

min_weight_fraction_leaf=0.0, n_estimators=i, n_jobs=-1, random_state=25, verbose=0
clf.fit(df_final_train, y_train)
train_sc = f1_score(y_train, clf.predict(df_final_train))
test_sc = f1_score(y_test, clf.predict(df_final_test))
test_scores.append(test_sc)
train_scores.append(train_sc)
print('Estimators = ', i, 'Train Score', train_sc, 'test Score', test_sc)
plt.plot(estimators, train_scores, label='Train Score')
plt.plot(estimators, test_scores, label='Test Score')
plt.xlabel('Estimators')
plt.ylabel('Score')
plt.title('Estimators vs score at depth of 5')

```

Estimators = 10 Train Score 0.9063252121775113 test Score 0.8745605278006858
 Estimators = 50 Train Score 0.9205725512208812 test Score 0.9125653355634538
 Estimators = 100 Train Score 0.9238690848446947 test Score 0.9141199714153599
 Estimators = 250 Train Score 0.9239789348046863 test Score 0.9188007232664732
 Estimators = 450 Train Score 0.9237190618658074 test Score 0.9161507685828595
 Text(0.5, 1, 'Estimators vs score at depth of 5')



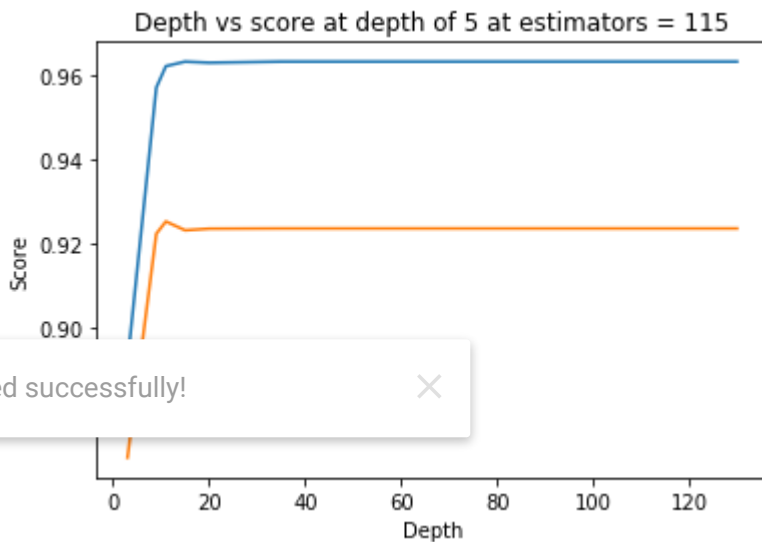
```

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

```

```
plt.xlabel('Depth')
plt.ylabel('Score')
plt.title('Depth vs score at depth of 5 at estimators = 115')
plt.show()
```

depth = 3 Train Score 0.8916120853581238 test Score 0.8687934859875491
 depth = 9 Train Score 0.9572226298198419 test Score 0.9222953031452904
 depth = 11 Train Score 0.9623451340902863 test Score 0.9252318758281279
 depth = 15 Train Score 0.9634267621927706 test Score 0.9231288356496615
 depth = 20 Train Score 0.9631629153051491 test Score 0.9235051024711141
 depth = 35 Train Score 0.9634333127085721 test Score 0.9235601652753184
 depth = 50 Train Score 0.9634333127085721 test Score 0.9235601652753184
 depth = 70 Train Score 0.9634333127085721 test Score 0.9235601652753184
 depth = 130 Train Score 0.9634333127085721 test Score 0.9235601652753184



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

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


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

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

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

```
mean test scores [0.96225043 0.96215493 0.96057081 0.96194015 0.96330005]
mean train scores [0.96294922 0.96266735 0.96115674 0.96263457 0.96430539]
```

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

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


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

```
clf.fit(df_final_train,y_train)
```

```
(train)
(est)
```

Saved successfully!

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

 Train f1 score 0.9652533106548414
Test f1 score 0.9241678239279553

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

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

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

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

    plt.subplot(1, 3, 2)
```



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

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

```
plt.show()
```

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

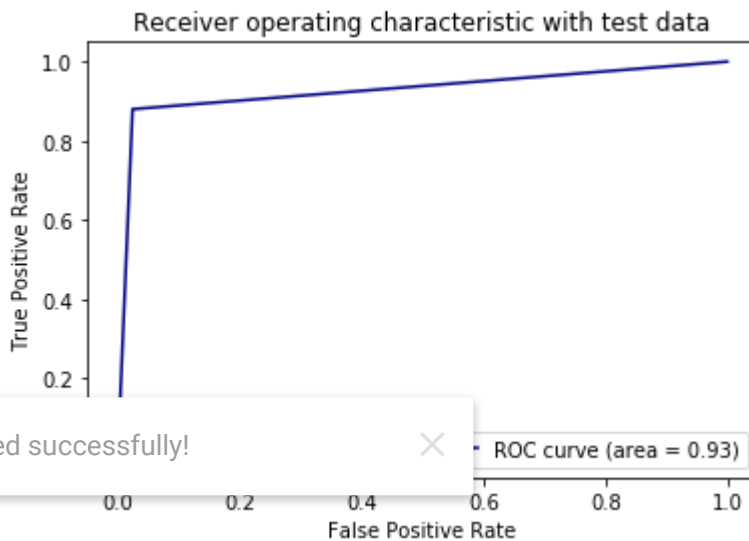


Saved successfully!



Train confusion matrix

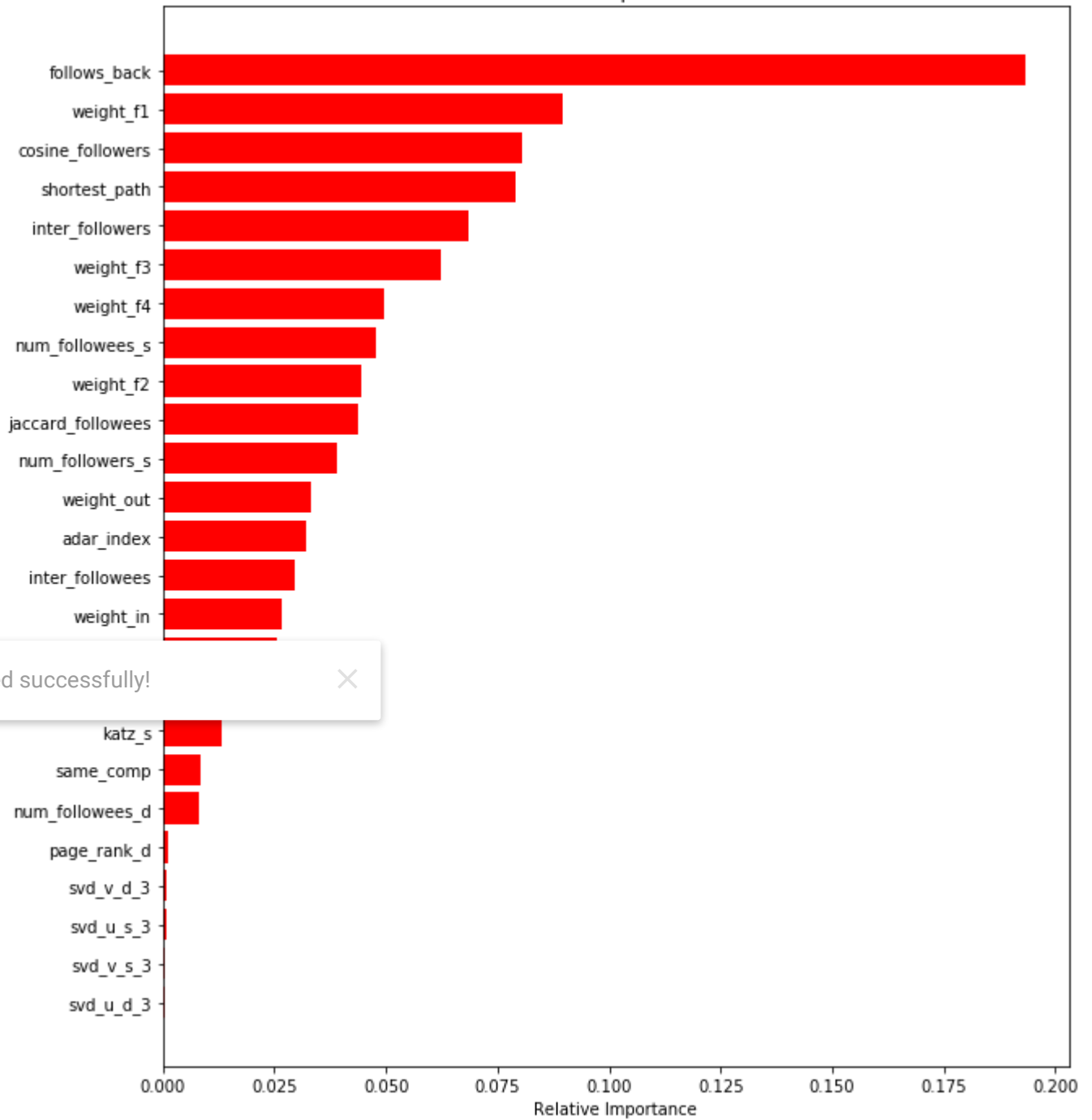
```
from sklearn.metrics import roc_curve, auc
fpr,tpr,ths = roc_curve(y_test,y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic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



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



Feature Importances



```

from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Model", "Train f1-score", "Test f1-score"]
x.add_row(['Random forest', '0.964', '0.921'])
x.add_row(['XGB00ST', '0.98', '0.91'])
print(x)

```



Model	Train f1-Score	Test f1-Score
Random Forest	0.964	0.921
XGB00ST	0.98	0.91

Saved successfully!

