

```
[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#plt
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
from pandas import HDFStore,DataFrame
from pandas import read_hdf
from scipy.sparse.linalg import svds, eig
import gc
from sklearn.metrics import log_loss

from tqdm import tqdm
from sklearn.calibration import CalibratedClassifierCV
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import f1_score
from sklearn.model_selection import train_test_split

In [2]: #Reading File train
if os.path.isfile('train_pos_after_edu.csv'):
    train_graphmx.read_edgelist('train_pos_after_edu.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: DiGraph
Number of nodes: 1780722
Number of edges: 7550815
Average in degree: 4.2399
Average out degree: 4.2399
```

Creating a new feature Preferential Attachment

```
In [3]: def PA_follower(a,b):
    try:
        if len(set(train_graph.successors(a))) == 0 | len(set(train_graph.successors(b))) == 0:
            return 0
        return len(set(train_graph.successors(a)))*len(set(train_graph.successors(b)))
    except:
        #print("Something went wrong in PA_follower please check it once")
        return 0

In [4]: def PA_followe(a,b):
    try:
        if len(set(train_graph.predecessors(a))) == 0 | len(set(train_graph.predecessors(b))) == 0:
            return 0
        return len(set(train_graph.predecessors(a)))*len(set(train_graph.predecessors(b)))
    except:
        return 0

In [5]: #Again reading the data frame to get source and destination nodes
df_final_train = read_hdf('storage_sample_stage4.h5', 'train_df',mode='r')
df_final_test = read_hdf('storage_sample_stage4.h5', 'test_df',mode='r')

In [6]: df_final_train['PA_followers'] = df_final_train.apply(lambda row:
    PA_follower(row['source_node'],row['destination_node']),axis=1)
df_final_test['PA_followers'] = df_final_test.apply(lambda row:
    PA_follower(row['source_node'],row['destination_node']),axis=1)

In [7]: df_final_train['PA_followe'] = df_final_train.apply(lambda row:
    PA_followe(row['source_node'],row['destination_node']),axis=1)
df_final_test['PA_followe'] = df_final_test.apply(lambda row:
    PA_followe(row['source_node'],row['destination_node']),axis=1)
```

Creating another New feature SVD_DOT

```
In [8]: def svd_dot_fun(a,b):
    temp=0
    for i,j in zip(a,b):
        temp=temp+j
    return temp

In [9]: #Creating another Feature svd_dot
U=[ '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']
V=[ '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']

In [10]: df_final_test['svd_dot']=df_final_test.apply(lambda row:svd_dot_fun(row[U],row[V]),axis=1)
df_final_train['svd_dot']=df_final_train.apply(lambda row: svd_dot_fun(row[U],row[V]),axis=1)

In [11]: #writing the df_final_train, df_final_test into csv files with new features svd_dot and Preferential Attachment
df_final_train.to_csv('df_final_train.csv')
df_final_test.to_csv('df_final_test.csv')

In [12]: #why suing these two files, we can avoid running above code
df_final_train=pd.read_csv('df_final_train.csv')
df_final_test=pd.read_csv('df_final_test.csv')
```

Creating a train,test,CV split

```
In [13]: X_train,X_cv=train_test_split(df_final_train,test_size=0.20)
y_train=X_train.indicator_link
y_cv=X_cv.indicator_link
y_test=df_final_test.indicator_link
X_test=df_final_test

In [14]: X_train.drop(['source_node','destination_node','indicator_link'],inplace=True,axis=1)
X_cv.drop(['source_node','destination_node','indicator_link'],inplace=True,axis=1)
X_test.drop(['source_node','destination_node','indicator_link'],inplace=True,axis=1)

In [15]: 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()
```

Building a model using with XGBoost

```
In [16]: alpha=[10,50,100,500,1000,2000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=1,nthread=1)
    x_cfl.fit(X_train,y_train)
    sig_cfl = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_cfl.fit(X_train,y_train)
    predict_y = sig_cfl.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_, eps=1e-15))

for i in range(len(cv_log_error_array)):
    print ('log_loss for c = ',alpha[i],',is',cv_log_error_array[i])

best_alpha = np.argmin(cv_log_error_array)

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))

plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=1)
x_cfl.fit(X_train,y_train)
sig_cfl = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_cfl.fit(X_train,y_train)

[22:11:04] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:05] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:05] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:06] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:07] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:08] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:08] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:10] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:14] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:16] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:17] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:19] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:22] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:24] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:27] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:29] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:31] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:33] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:43] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:11:58] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:04] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:10] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:16] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:17] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:27] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:12:49] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:13:09] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:13:18] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:13:40] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:14:51] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

log_loss for c = 10 is 0.00013687540893760807
log_loss for c = 50 is 0.0001372631785593055
log_loss for c = 100 is 0.0001372031987497085
log_loss for c = 500 is 0.0001372031987497085
log_loss for c = 1000 is 0.00013720320452998077
log_loss for c = 2000 is 0.0001372032153564643

Cross Validation Error for each alpha
Error measure
0.00013720
0.00013715
0.00013710
0.00013705
0.00013695
0.00013690
0
250 500 750 1000 1250 1500 1750 2000
Alpha
```

```
[22:15:09] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:15:10] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:15:11] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:15:12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:15:12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:15:13] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[22:15:13] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release.1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

Out[16]: CalibratedClassifierCV(base_estimator=XGBClassifier(base_score=0.5,
    booster='gbtree',
    colsample_bylevel=1,
    colsample_bynode=1,
    colsample_bytree=1,
    enable_categorical=False,
    gamma=0,
    gpu_id=-1,
    importance_type=None,
    interaction_constraints='',
    learning_rate=0.300000012,
    max_delta_step=0,
    max_depth=6,
    min_child_weight=1,
    missing=None,
    monotone_constraints=(),
    n_estimators=10, n_jobs=12,
    nthread=-1,
    num_parallel_tree=1,
    predictor='auto',
    random_state=0, reg_alpha=0,
    reg_lambda=1,
    sample_weight=1,
    subsample=1,
    tree_method='exact',
    validate_parameters=1,
    verbosity=None))

In [17]: predict_y = sig_cfl.predict_proba(X_train)
print ("For values of best alpha = ", alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
predict_y = sig_cfl.predict_proba(X_cv)
print ("For values of best alpha = ", alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
predict_y = sig_cfl.predict_proba(X_test)
print ("For values of best alpha = ", alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
plot_confusion_matrix(y_test, sig_cfl.predict(X_test))

For values of best alpha = 10 The train log loss is: 0.00013983210767471455
For values of best alpha = 10 The cross validation log loss is: 0.00013687540893760807
For values of best alpha = 10 The test log loss is: 0.00013924075018745561

Confusion matrix
0 49952.000 0.000
1 0.000 50050.000
0 1
1
0 1
1
Predicted Class
```

Summary

- 1.In this case study we add two features Preferential Attachment and svd_dot
- 2.Build a XGBoost model with best hyperparameter of alpha 10, and got a test loss of 0.0001 which is the best value for this model.

- 1. Confusion matrix also shows a great results

Steps followed to slove the FaceBook prediction caseStudy

- 1.We defined the machine learning problem, i.e to predict the whether a relation might exists in the future between two persons or not
- 2.After seeing the dataset we analysed that we have only, possitive class data i.e we have only graph data where a link is present . From that we can say that we have only possitive class, so we added some random data as class 0, where no link is present between them.
- 3.Now to handle the graph data we will use a library called networkx which will handle the graph data, this module will play an important role in finding the various metrics about the directed graph.
- 4.Now we will do some feature engineering, to get new features such as Jaccard& cosine similarities, PageRank, Shortest path,Adar index etc.
- 5.We made a train test split randomly as we dont have any timestamp data.
- 6.Now based on the above data we built various models, such as linear Regression, Randomforest,XGBoost etc.and calculated various metric related to that models and found that the above XGBoost model will perform well.