# FB_Models

September 6, 2019

Social network Graph Link Prediction - Facebook Challenge

```python
[32]: #Importing Libraries
      # please do go through this python notebook:
      import warnings
      warnings.filterwarnings("ignore")

      import csv
      import pandas as pd#pandas to create small dataframes
      import datetime #Convert to unix time
      import time #Convert to unix time
      # if numpy is not installed already : pip3 install numpy
      import numpy as np#Do aritmetic operations on arrays
      # matplotlib: used to plot graphs
      import matplotlib
      import matplotlib.pylab as plt
      import seaborn as sns#Plots
      from matplotlib import rcParams#Size of plots
      from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
      import math
      import pickle
      import os
      # to install xgboost: pip3 install xgboost
      import xgboost as xgb

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

```
[33]: #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')
```

```
[34]: df_final_train.columns
```

```
[34]: Index(['source_node', 'destination_node', 'indicator_link',
             'jaccard_followers', 'jaccard_followees', 'cosine_followers',
             'cosine_followees', 'num_followers_s', 'num_followers_d',
             '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',
             'preferential_attachment_followers',
             'preferential_attachment_followees', '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', 'svd_dot_u', 'svd_dot_v'],
            dtype='object')
```

```
[35]: y_train = df_final_train.indicator_link
      y_test = df_final_test.indicator_link
```

```
[36]: 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)
```

### 0.0.1  1.1 Random Forest

```
[37]: estimators = [10,50,100,250,450]
      train_scores = []
      test_scores = []
      for i in tqdm(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,warm_start=False)
          clf.fit(df_final_train,y_train)
```