

# Social network Graph Link Prediction - Facebook Challenge

```
In [1]: 1 #Importing Libraries
2 # please do go through this python notebook:
3 import warnings
4 warnings.filterwarnings("ignore")
5
6 import csv
7 import pandas as pd#pandas to create small dataframes
8 import datetime #Convert to unix time
9 import time #Convert to unix time
10 # if numpy is not installed already : pip3 install numpy
11 import numpy as np#Do arithmetic operations on arrays
12 # matplotlib: used to plot graphs
13 import matplotlib
14 import matplotlib.pyplot as plt
15 import seaborn as sns#Plots
16 from matplotlib import rcParams#Size of plots
17 from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
18 import math
19 import pickle
20 import os
21 # to install xgboost: pip3 install xgboost
22 import xgboost as xgb
23
24 import warnings
25 import networkx as nx
26 import pdb
27 import pickle
28 from pandas import HDFStore, DataFrame
29 from pandas import read_hdf
30 from scipy.sparse.linalg import svds, eigs
31 import gc
32 from tqdm import tqdm
33 from sklearn.ensemble import RandomForestClassifier
34 from sklearn.metrics import f1_score
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\ensemble\weight\_boosting.py:  
29: DeprecationWarning: numpy.core.umath\_tests is an internal NumPy module and  
should not be imported. It will be removed in a future NumPy release.  
from numpy.core.umath\_tests import inner1d

```
In [0]: 1 #reading
2 from pandas import read_hdf
3 df_final_train = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'train')
4 df_final_test = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'test_d')
```

```
In [0]: 1 df_final_train.columns
```

```
Out[3]: 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'],  
             dtype='object')
```

```
In [0]: 1 y_train = df_final_train.indicator_link  
        2 y_test = df_final_test.indicator_link
```

```
In [0]: 1 df_final_train.drop(['source_node', 'destination_node', 'indicator_link'], axis=  
        2 df_final_test.drop(['source_node', 'destination_node', 'indicator_link'], axis=
```

```

In [0]: 1 estimators = [10,50,100,250,450]
2 train_scores = []
3 test_scores = []
4 for i in estimators:
5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion
6         max_depth=5, max_features='auto', max_leaf_nodes=None,
7         min_impurity_decrease=0.0, min_impurity_split=None,
8         min_samples_leaf=52, min_samples_split=120,
9         min_weight_fraction_leaf=0.0, n_estimators=i, n_jobs=-1,random_s
10     clf.fit(df_final_train,y_train)
11     train_sc = f1_score(y_train,clf.predict(df_final_train))
12     test_sc = f1_score(y_test,clf.predict(df_final_test))
13     test_scores.append(test_sc)
14     train_scores.append(train_sc)
15     print('Estimators = ',i,'Train Score',train_sc,'test Score',test_sc)
16 plt.plot(estimators,train_scores,label='Train Score')
17 plt.plot(estimators,test_scores,label='Test Score')
18 plt.xlabel('Estimators')
19 plt.ylabel('Score')
20 plt.title('Estimators vs score at depth of 5')

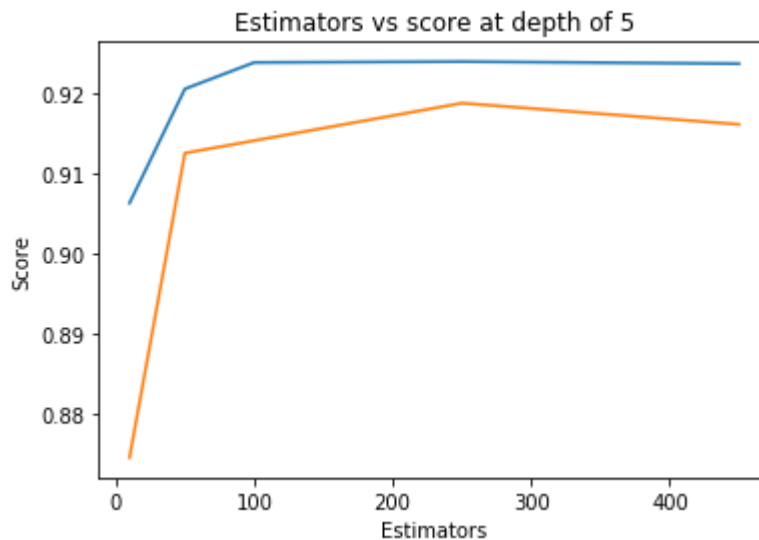
```

```

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

```

Out[6]: Text(0.5,1,'Estimators vs score at depth of 5')



```

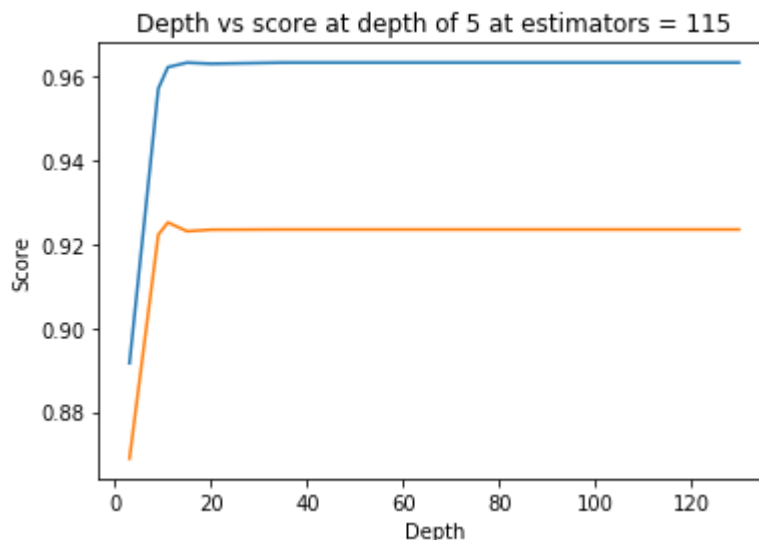
In [0]: 1 depths = [3,9,11,15,20,35,50,70,130]
2 train_scores = []
3 test_scores = []
4 for i in depths:
5     clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion
6         max_depth=i, max_features='auto', max_leaf_nodes=None,
7         min_impurity_decrease=0.0, min_impurity_split=None,
8         min_samples_leaf=52, min_samples_split=120,
9         min_weight_fraction_leaf=0.0, n_estimators=115, n_jobs=-1,random
10     clf.fit(df_final_train,y_train)
11     train_sc = f1_score(y_train,clf.predict(df_final_train))
12     test_sc = f1_score(y_test,clf.predict(df_final_test))
13     test_scores.append(test_sc)
14     train_scores.append(train_sc)
15     print('depth = ',i,'Train Score',train_sc,'test Score',test_sc)
16 plt.plot(depths,train_scores,label='Train Score')
17 plt.plot(depths,test_scores,label='Test Score')
18 plt.xlabel('Depth')
19 plt.ylabel('Score')
20 plt.title('Depth vs score at depth of 5 at estimators = 115')
21 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

```



```
In [0]: 1 from sklearn.metrics import f1_score
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.metrics import f1_score
4 from sklearn.model_selection import RandomizedSearchCV
5 from scipy.stats import randint as sp_randint
6 from scipy.stats import uniform
7
8 param_dist = {"n_estimators":sp_randint(105,125),
9               "max_depth": sp_randint(10,15),
10              "min_samples_split": sp_randint(110,190),
11              "min_samples_leaf": sp_randint(25,65)}
12
13 clf = RandomForestClassifier(random_state=25,n_jobs=-1)
14
15 rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
16                                n_iter=5,cv=10,scoring='f1',random_state=
17
18 rf_random.fit(df_final_train,y_train)
19 print('mean test scores',rf_random.cv_results_['mean_test_score'])
20 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]

```
In [0]: 1 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)
```

```
In [0]: 1 clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='g
2         max_depth=14, max_features='auto', max_leaf_nodes=None,
3         min_impurity_decrease=0.0, min_impurity_split=None,
4         min_samples_leaf=28, min_samples_split=111,
5         min_weight_fraction_leaf=0.0, n_estimators=121, n_jobs=-1,
6         oob_score=False, random_state=25, verbose=0, warm_start=False)
```

```
In [0]: 1 clf.fit(df_final_train,y_train)
2 y_train_pred = clf.predict(df_final_train)
3 y_test_pred = clf.predict(df_final_test)
```

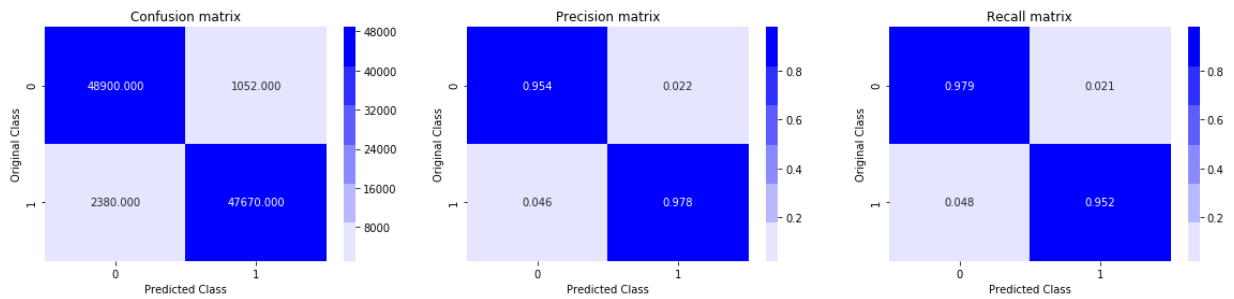
```
In [0]: 1 from sklearn.metrics import f1_score
2 print('Train f1 score',f1_score(y_train,y_train_pred))
3 print('Test f1 score',f1_score(y_test,y_test_pred))
```

Train f1 score 0.9652533106548414  
 Test f1 score 0.9241678239279553

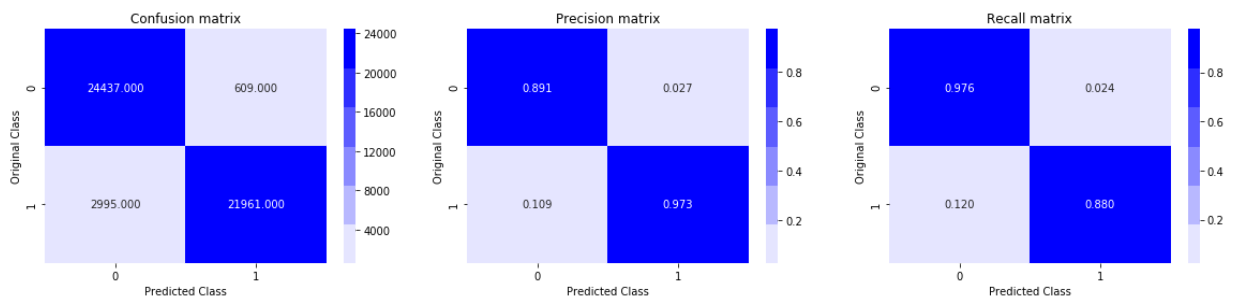
```
In [61]: 1 from sklearn.metrics import confusion_matrix
2 def plot_confusion_matrix(test_y, predict_y):
3     C = confusion_matrix(test_y, predict_y)
4
5     A = ((C.T)/(C.sum(axis=1))).T
6
7     B = (C/C.sum(axis=0))
8     plt.figure(figsize=(20,4))
9
10    labels = [0,1]
11    # representing A in heatmap format
12    cmap=sns.light_palette("blue")
13    plt.subplot(1, 3, 1)
14    sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yti
15    plt.xlabel('Predicted Class')
16    plt.ylabel('Original Class')
17    plt.title("Confusion matrix")
18
19    plt.subplot(1, 3, 2)
20    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yti
21    plt.xlabel('Predicted Class')
22    plt.ylabel('Original Class')
23    plt.title("Precision matrix")
24
25    plt.subplot(1, 3, 3)
26    # representing B in heatmap format
27    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yti
28    plt.xlabel('Predicted Class')
29    plt.ylabel('Original Class')
30    plt.title("Recall matrix")
31
32    plt.show()
```

```
In [0]: 1 print('Train confusion_matrix')
2 plot_confusion_matrix(y_train,y_train_pred)
3 print('Test confusion_matrix')
4 plot_confusion_matrix(y_test,y_test_pred)
```

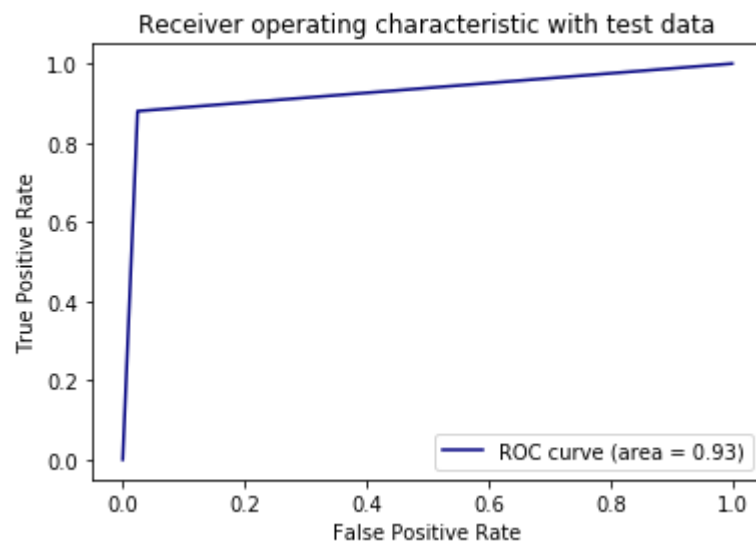
Train confusion\_matrix



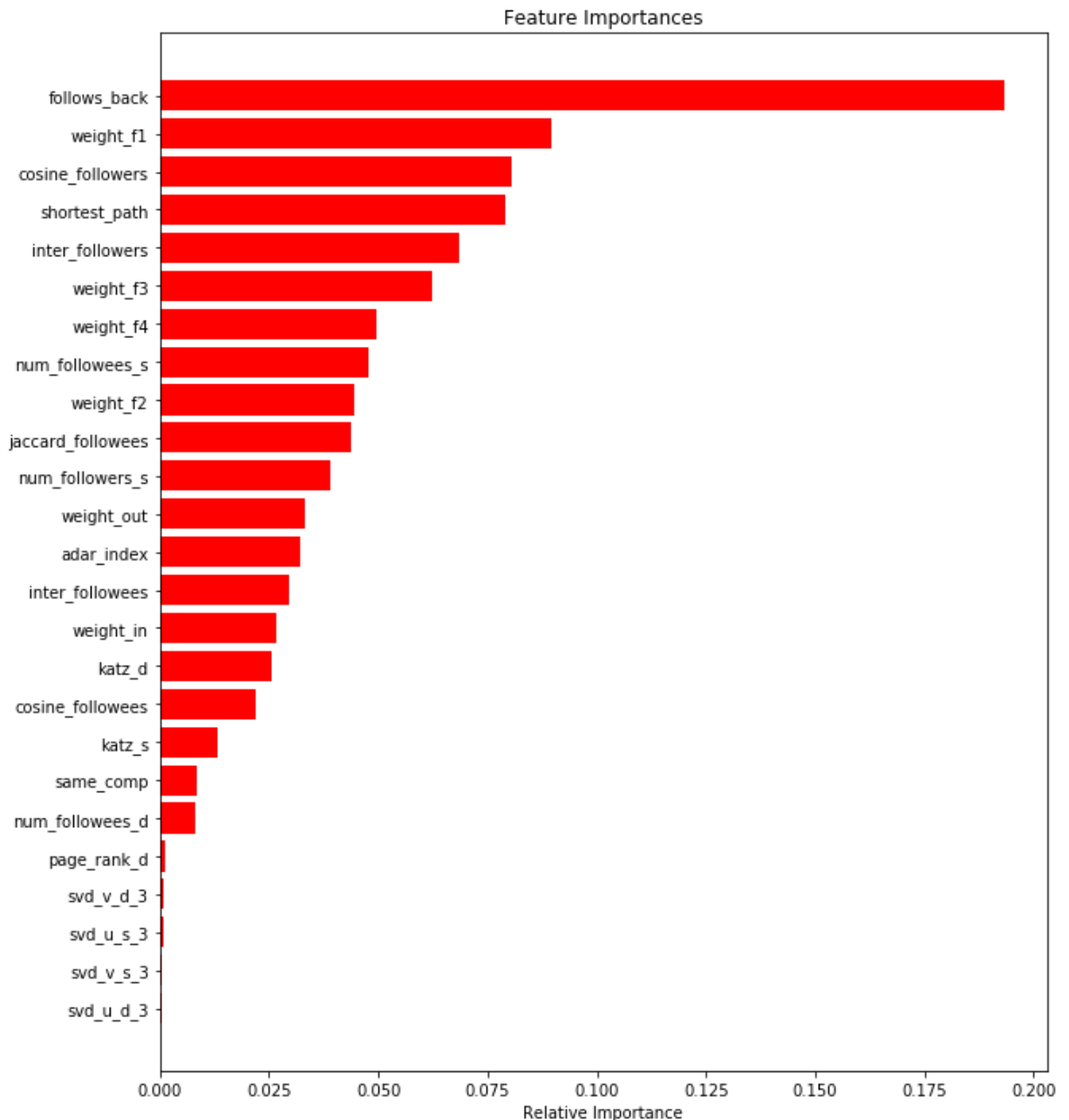
Test confusion\_matrix



```
In [0]: 1 from sklearn.metrics import roc_curve, auc
2 fpr, tpr, ths = roc_curve(y_test,y_test_pred)
3 auc_sc = auc(fpr, tpr)
4 plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('Receiver operating characteristic with test data')
8 plt.legend()
9 plt.show()
```



```
In [0]: 1 features = df_final_train.columns
2 importances = clf.feature_importances_
3 indices = (np.argsort(importances))[-25:]
4 plt.figure(figsize=(10,12))
5 plt.title('Feature Importances')
6 plt.barh(range(len(indices)), importances[indices], color='r', align='center')
7 plt.yticks(range(len(indices)), [features[i] for i in indices])
8 plt.xlabel('Relative Importance')
9 plt.show()
```



## Assignments:

1. Add another feature called Preferential Attachment with followers and followees data of vertex. you can check about Preferential Attachment in below link <http://be.amazd.com/link-prediction/> (<http://be.amazd.com/link-prediction/>)



2. Add feature called svd\_dot. you can calculate svd\_dot as Dot product between source node svd and destination node svd features. you can read about this in below pdf  
[https://storage.googleapis.com/kaggle-forum-message-attachments/2594/supervised\\_link\\_prediction.pdf](https://storage.googleapis.com/kaggle-forum-message-attachments/2594/supervised_link_prediction.pdf) ([https://storage.googleapis.com/kaggle-forum-message-attachments/2594/supervised\\_link\\_prediction.pdf](https://storage.googleapis.com/kaggle-forum-message-attachments/2594/supervised_link_prediction.pdf))
3. Tune hyperparameters for XG boost with all these features and check the error metric.

```
In [31]: 1 #reading
2 from pandas import read_hdf
3 X_train = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'train_df', mode='a')
4 X_test = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'test_df', mode='a')
```

```
In [32]: 1 Y_train = X_train.indicator_link
2 Y_test = X_test.indicator_link
3 X_train.drop(['source_node', 'destination_node', 'indicator_link'], axis=1, inplace=True)
4 X_test.drop(['source_node', 'destination_node', 'indicator_link'], axis=1, inplace=True)
5 predictors = X_train.columns
```

```
In [55]: 1 predictors = [str(x) for x in X_train.columns]
```

```
In [33]: 1 X_train.columns
```

```
Out[33]: Index(['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',
               'svd_dot_U_U', 'svd_dot_U_V', 'svd_dot_V_V',
               'preferential_attachment_followers',
               'preferential_attachment_followees'],
              dtype='object')
```

In [34]: 1 X\_test.columns

```
Out[34]: Index(['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',
'svd_dot_U_U', 'svd_dot_U_V', 'svd_dot_V_V',
'preferential_attachment_followers',
'preferential_attachment_followees'],
dtype='object')
```

```
In [22]: 1 import xgboost
2 from xgboost import XGBClassifier
3 from sklearn import metrics
4 from xgboost import plot_importance
5
6 def plot_features(booster, figsize):
7     fig, ax = plt.subplots(1,1,figsize=figsize)
8
9     return plot_importance(booster=booster, ax=ax, max_num_features=40)
10
11 def xgb_model(model,train_, test_):
12     xgb_param = model.get_xgb_params()
13     xgtrain = xgboost.DMatrix(train_, label = Y_train)
14     cv_result = xgboost.cv(xgb_param, xgtrain,
15                             num_boost_round=model.get_params()['n_estimators'],
16                             metrics= 'auc',
17                             early_stopping_rounds=50)
18     model.set_params(n_estimators =cv_result.shape[0])
19     model.fit(train_, Y_train, eval_metric='auc')
20     train_predict = model.predict(train_)
21     train_predict_proba = model.predict_proba(train_)[:, -1]
22     print("Train Accuracy : %.4g" % metrics.accuracy_score(Y_train, train_pr
23     print("AUC Score (Train): %f" % metrics.roc_auc_score(Y_train, train_pre
24     print('#'*50)
25     test_predict = model.predict(test_)
26     test_predict_proba = model.predict_proba(test_)[:, -1]
27     print("Test Accuracy : %.4g" % metrics.accuracy_score(Y_test, test_predi
28     print("AUC Score (test): %f" % metrics.roc_auc_score(Y_test, test_predic
29     model.get_booster().feature_names = predictors
30     plot_features(model, (16,27))
31
32
```

```

In [23]: 1 from sklearn.model_selection import GridSearchCV
2 param_test1 = {
3     'max_depth':range(6,14,2),
4     'min_child_weight':range(1,6,2)
5 }
6 gsearch1 = GridSearchCV(
7     estimator = XGBClassifier(learning_rate =0.1,
8                               n_estimators=140,
9                               gamma=0,
10                              subsample=0.8,
11                              colsample_bytree=0.8,
12                              objective= 'binary:logistic',
13                              nthread=4,
14                              scale_pos_weight=1,
15                              seed=27),
16     param_grid = param_test1,
17     scoring='roc_auc',
18     n_jobs=3,
19     iid=False,
20     cv=5)
21 gsearch1.fit(X_train,Y_train)
22 # gsearch1.grid_scores_, gsearch1.best_params_, gsearch1.best_score_gb1, (15

```

```

Out[23]: GridSearchCV(cv=5, error_score='raise',
    estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_bytree=1,
    colsample_bynode=1, colsample_bytree=0.8, gamma=0,
    learning_rate=0.1, max_delta_step=0, max_depth=5,
    min_child_weight=1, missing=None, n_estimators=140, n_jobs=1,
    nthread=4, objective='binary:logistic', random_state=0, reg_alpha=0,
    reg_lambda=1, scale_pos_weight=1, seed=27, silent=None,
    subsample=0.8, verbosity=1),
    fit_params=None, iid=False, n_jobs=3,
    param_grid={'max_depth': range(6, 14, 2), 'min_child_weight': range(1,
    6, 2)},
    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
    scoring='roc_auc', verbose=0)

```

```

In [24]: 1 gsearch1.best_params_

```

```

Out[24]: {'max_depth': 12, 'min_child_weight': 5}

```

```
In [59]: 1 xgb1 = XGBClassifier(learning_rate =0.1,
2                 n_estimators=140,
3                 max_depth=12,
4                 min_child_weight=5,
5                 gamma=0,
6                 subsample=0.8,
7                 colsample_bytree=0.8,
8                 objective= 'binary:logistic',
9                 nthread=4,
10                scale_pos_weight=1,
11                seed=27)
12 xgb_model(xgb1,X_train,X_test)
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\preprocessing\label.py:151:
DeprecationWarning: The truth value of an empty array is ambiguous. Returning F
alse, but in future this will result in an error. Use `array.size > 0` to check
that an array is not empty.
```

```
    if diff:
```

```
Train Accuracy : 0.9974
```

```
AUC Score (Train): 0.999981
```

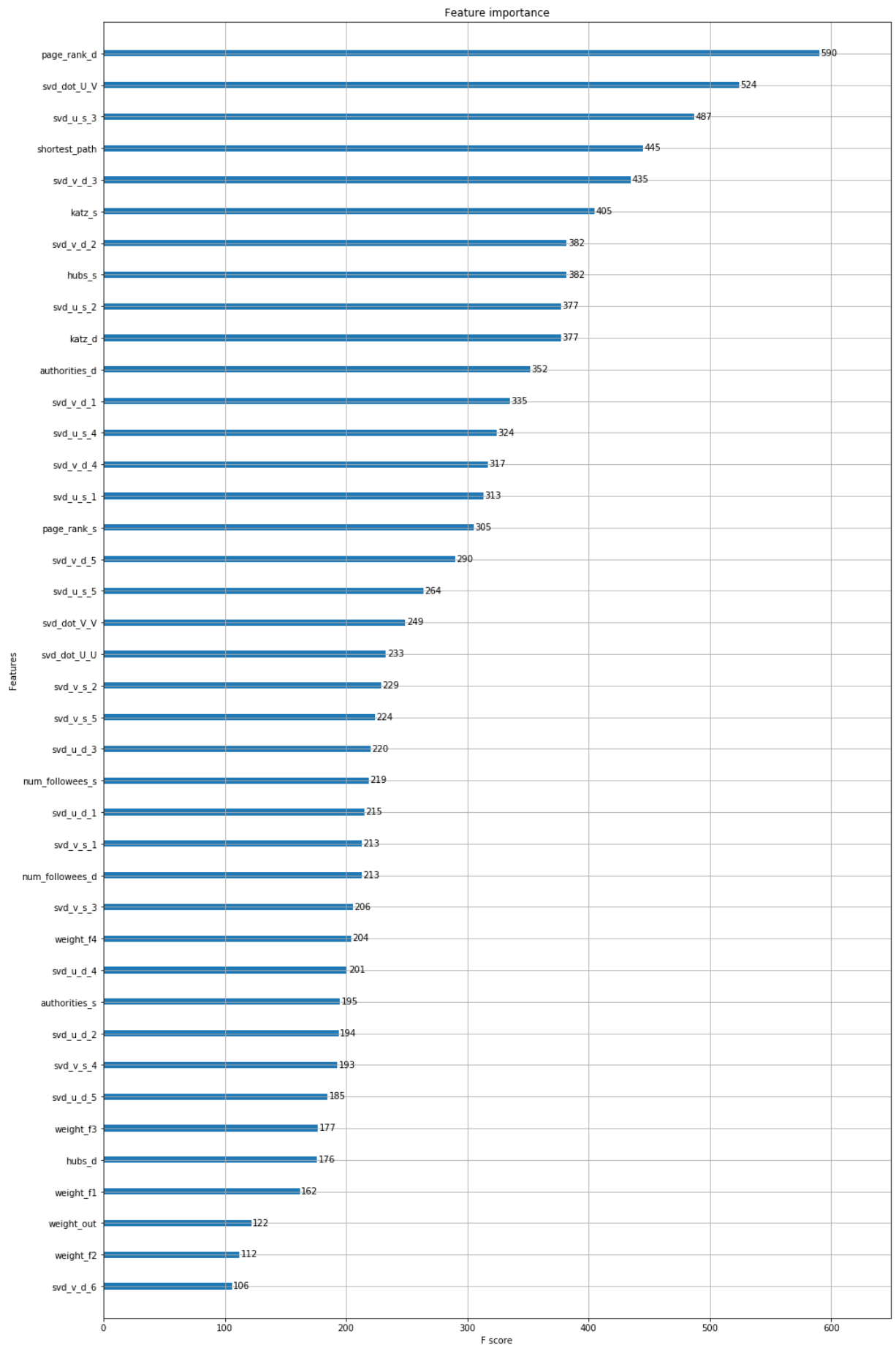
```
#####
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\preprocessing\label.py:151:
DeprecationWarning: The truth value of an empty array is ambiguous. Returning F
alse, but in future this will result in an error. Use `array.size > 0` to check
that an array is not empty.
```

```
    if diff:
```

```
Test Accuracy : 0.933
```

```
AUC Score (test): 0.948747
```



```
In [63]: 1 y_train_pred = xgb1.predict(X_train)
2 y_test_pred = xgb1.predict(X_test)
3 print('Train confusion_matrix')
4 plot_confusion_matrix(Y_train,y_train_pred)
5 print('Test confusion_matrix')
6 plot_confusion_matrix(Y_test,y_test_pred)
```

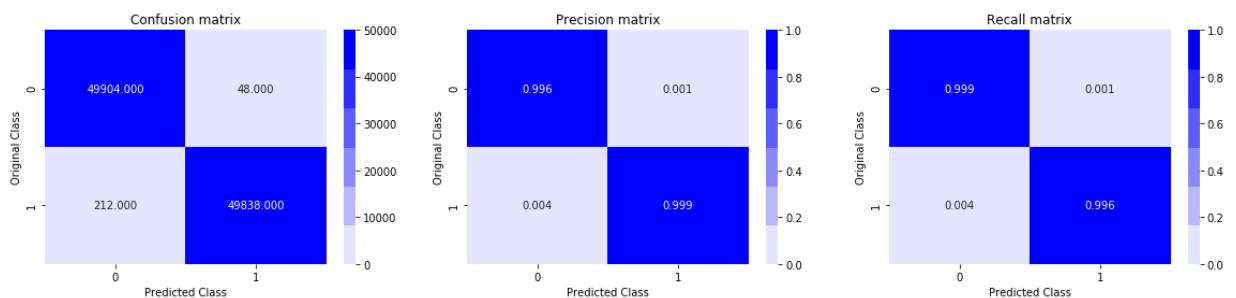
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\preprocessing\label.py:151: DeprecationWarning: The truth value of an empty array is ambiguous. Returning F  
 else, but in future this will result in an error. Use `array.size > 0` to check  
 that an array is not empty.

if diff:

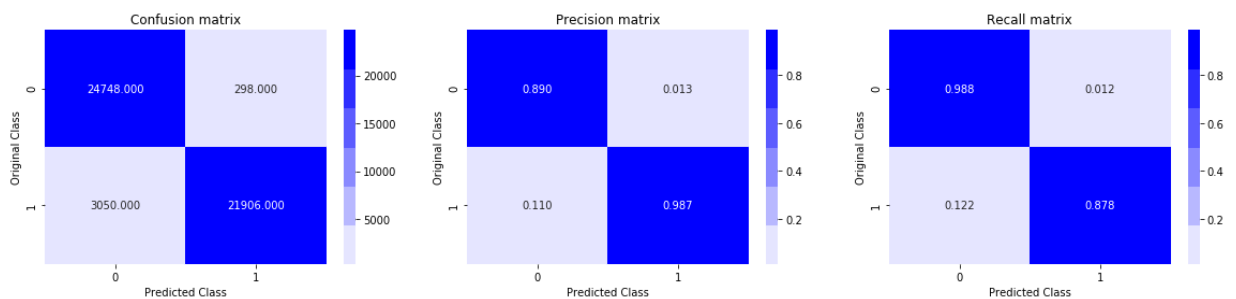
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\preprocessing\label.py:151: DeprecationWarning: The truth value of an empty array is ambiguous. Returning F  
 else, but in future this will result in an error. Use `array.size > 0` to check  
 that an array is not empty.

if diff:

Train confusion\_matrix



Test confusion\_matrix



In [ ]:

1

In [ ]:

1

In [ ]:

1

In [ ]:

1

In [ ]:

1

