

Social network Graph Link Prediction - Facebook Challenge

In [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#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
```

In [2]:

```
#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')
```

In [3]:

```
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)
```

In [0]:

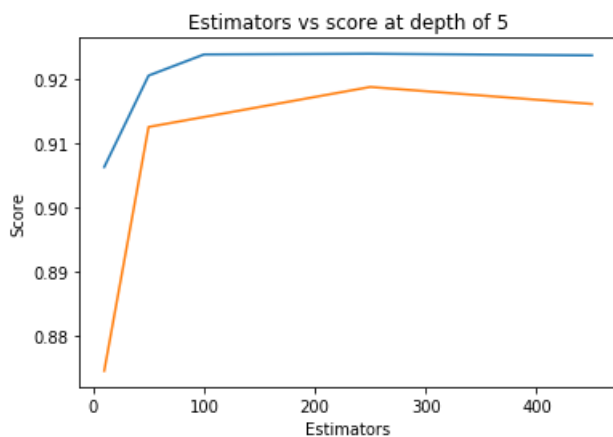
```
estimators = [10, 50, 100, 250, 450]
train_scores = []
test_scores = []
for i in estimators:
    clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=5, max_features='auto', max_leaf_nodes=None,
                                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, warm_start=False)
    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
```

Out[0]:

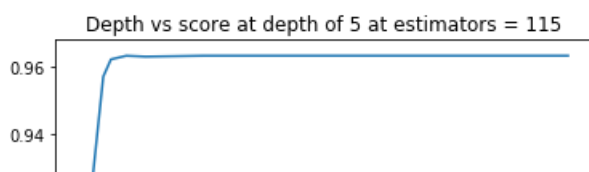
Text(0.5,1,'Estimators vs score at depth of 5')

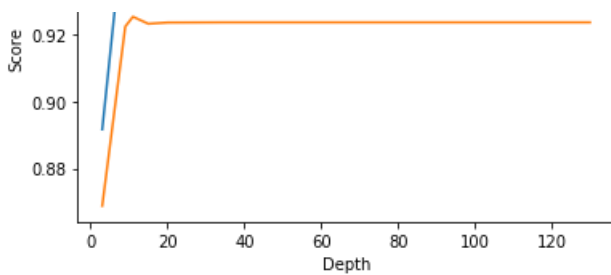


In [0]:

```
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, warm_start=False)
    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
```





In [0]:

```
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]
```

In [0]:

```
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]:

```
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)
```

In [0]:

```
clf.fit(df_final_train,y_train)
y_train_pred = clf.predict(df_final_train)
y_test_pred = clf.predict(df_final_test)
```

In [0]:

```
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
```

In [3]:

```
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")

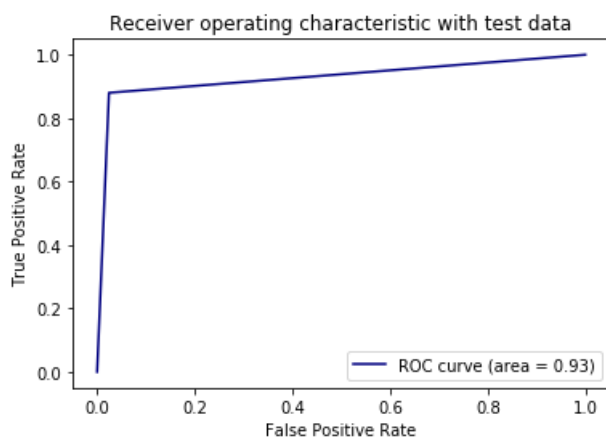
    plt.subplot(1, 3, 2)
    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Precision matrix")

    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()
```

In [0]:

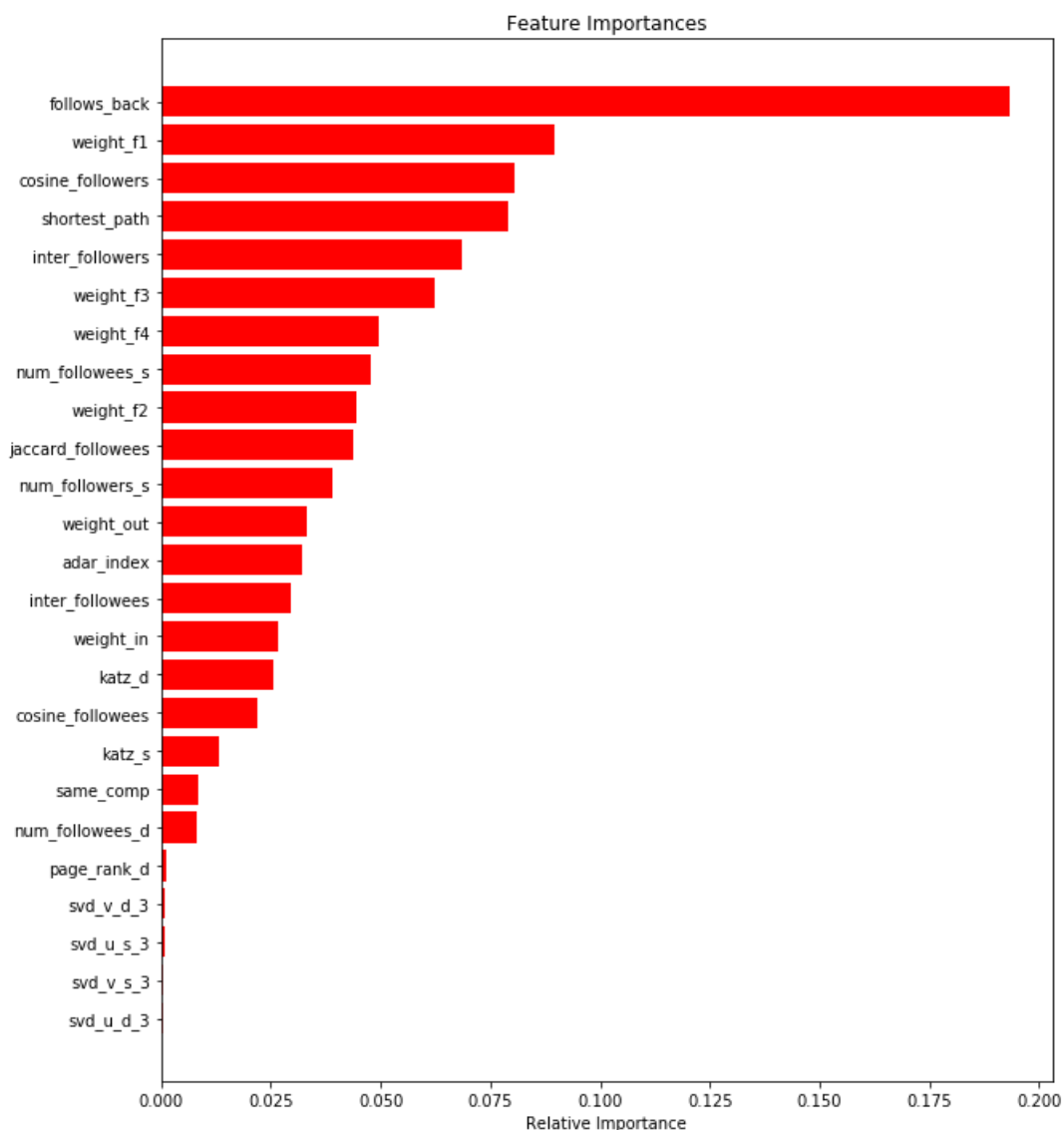
```
from sklearn.metrics import roc_curve, auc
fpr,tpr,ths = roc_curve(y_test,y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy',label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic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



In [0]:

```
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 Importances')
```

```
plt.xlabel('Relative Importance')
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/>
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
3. Tune hyperparameters for XG boost with all these features and check the error metric.

1. Preferential Attachment

In [6]:

```
if os.path.isfile('data/after_eda/train_pos_after_eda.csv'):
    train_graph=nx.read_edgelist('data/after_eda/train_pos_after_eda.csv',delimiter=',',create_using=nx.DiGraph(),nodetype=int)
    print(nx.info(train_graph))
else:
    print("please run the FB_EDA.ipynb or download the files from drive")
```

Name:
Type: DiGraph
Number of nodes: 1780722
Number of edges: 7550015
Average in degree: 4.2399

Average out degree: 4.2399

In [3]:

```
def preferential_attachment_for_followers(a,b):
    try:

        if len(set(train_graph.predecessors(a))) == 0 | len(set(train_graph.predecessors(b))) == 0
:
            return 0
        score = (len(set(train_graph.predecessors(a)))*(set(train_graph.predecessors(b))))

        return score
    except:
        return 0
```

In [4]:

```
def preferential_attachment_for_followees(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)))*(set(train_graph.predecessors(b))))

        return sim
    except:
        return 0
```

In [5]:

```
print(preferential_attachment_for_followers(273084,470294))
```

0

In [6]:

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

0

In [7]:

```
#mapping preferential followers to train and test data
df_final_train['preferential_attachment_followers'] = df_final_train.apply(lambda row:
                                                                           preferential_attachment_for_followers(row['source_node'
                                                                           ,row['destination_node']],axis=1)
df_final_test['preferential_attachment_followers'] = df_final_test.apply(lambda row:
                                                                           preferential_attachment_for_followers(row['source_node'
                                                                           ,row['destination_node']],axis=1)

#mapping preferential followees to train and test data
df_final_train['preferential_attachment_followees'] = df_final_train.apply(lambda row:
                                                                           preferential_attachment_for_followees(row['source_node'
                                                                           ,row['destination_node']],axis=1)
df_final_test['preferential_attachment_followees'] = df_final_test.apply(lambda row:
                                                                           preferential_attachment_for_followees(row['source_node'
                                                                           ,row['destination_node']],axis=1)
```

2. SVD_DOT (dot product of both source and destination)

In [8]:

```
def svd_dot(x,y):  
    return np.dot(x,y)
```

In [9]:

```
a=svd_dot(df_final_train.destination_node,df_final_train.source_node)
```

In [10]:

```
b=svd_dot(df_final_train.destination_node,df_final_train.source_node)
```

In [11]:

```
c=svd_dot(df_final_test.destination_node,df_final_test.source_node)
```

In [12]:

```
d=c=svd_dot(df_final_test.destination_node,df_final_test.source_node)
```

In [13]:

```
import warnings  
warnings.filterwarnings("ignore")  
  
df_final_train['svd_dot_U']=a
```

In [14]:

```
df_final_train['svd_dot_V.T']=b
```

In [15]:

```
df_final_test['svd_dot_U']=c
```

In [16]:

```
df_final_test['svd_dot_V.T']=d
```

In [17]:

```
y_train=df_final_train['indicator_link'].values
```

In [18]:

```
y_test=df_final_test['indicator_link'].values
```

In [19]:

```
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)
```

3. Hyperparameter Tuning of XG boost with all the above features and checking the error metric.

In [26]:

```
from sklearn.model_selection import GridSearchCV  
from xgboost import XGBClassifier  
  
model = XGBClassifier()  
param_grid = {"max_depth":[1, 2, 3], "n_estimators":[40, 80, 150], "subsample": [0.9,  
1.0],"learning_rate": [0.05, 0.1, 0.3],}
```

```
clf = GridSearchCV(model, param_grid, scoring='f1', cv = 5)

clf.fit(df_final_train,y_train)
```

```
train_auc_mean= clf.cv_results_['mean_train_score']
train_auc_test= clf.cv_results_['mean_test_score']
```

In [27]:

```
best_parameters=clf.best_params_
```

In [28]:

```
for param_name in sorted(best_parameters.keys()):
    print("%s: %r" % (param_name, best_parameters[param_name]))
```

```
learning_rate: 0.3
max_depth: 3
n_estimators: 150
subsample: 0.9
```

In [29]:

```
df_final_train.shape
```

Out[29]:

```
(100002, 55)
```

In [20]:

```
from xgboost import XGBClassifier
clf_xgb = XGBClassifier(max_depth=1,n_estimators=150,subsample=0.9,learning_rate=0.3)
clf_xgb.fit(df_final_train,y_train)
```

Out[20]:

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
               colsample_bytree=1, gamma=0, learning_rate=0.3, max_delta_step=0,
               max_depth=1, min_child_weight=1, missing=None, n_estimators=150,
               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=0.9)
```

In [33]:

```
y_train_pred = clf_xgb.predict(df_final_train)
y_test_pred = clf_xgb.predict(df_final_test)
```

In [23]:

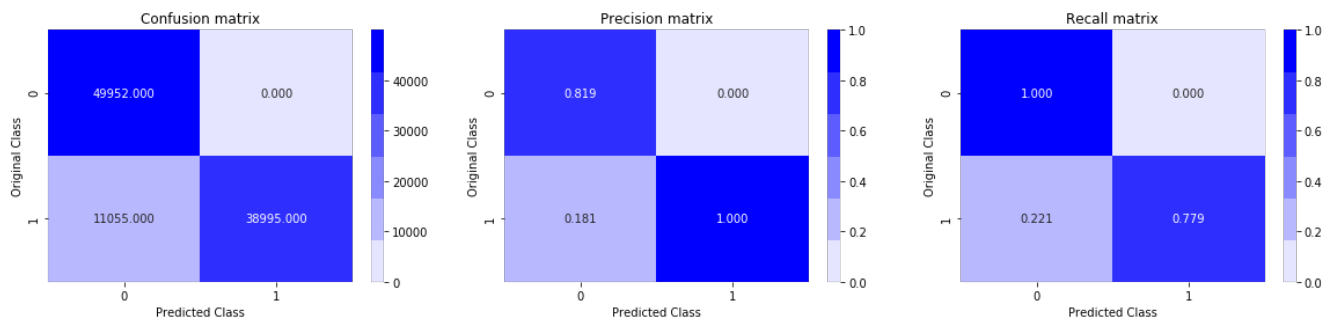
```
from sklearn.metrics import f1_score
print('Train f1 score',f1_score(y_train,y_train_pred))
print('Test f1 score',f1_score(y_test,y_test_pred))
```

```
Train f1 score 0.8758492896849908
Test f1 score 0.8713565339303078
```

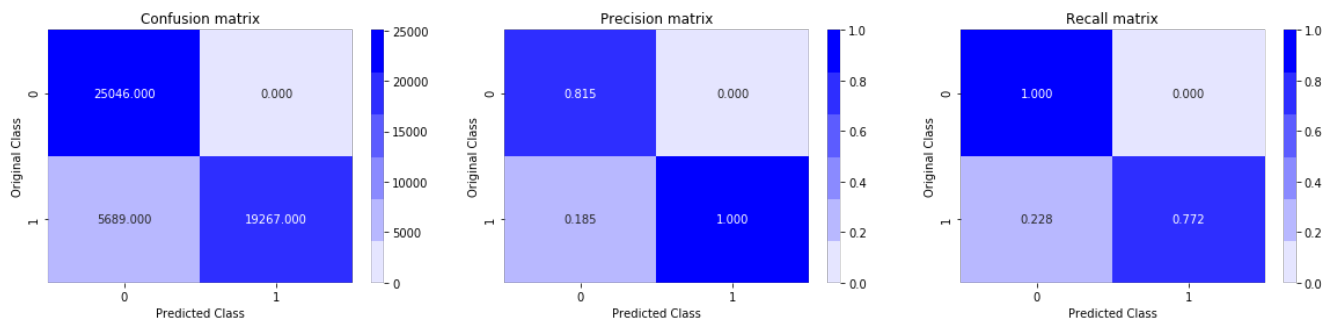
In [35]:

```
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
```

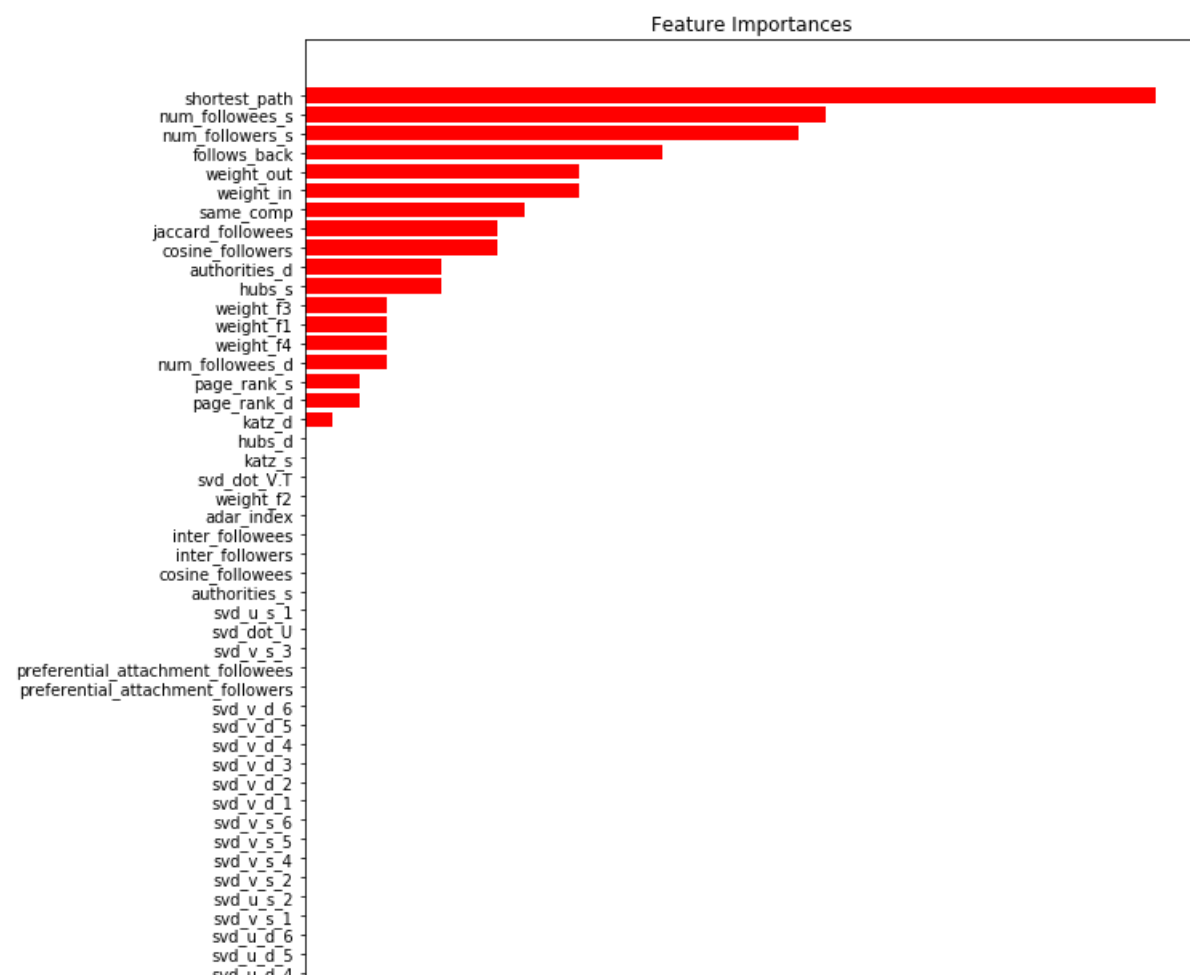



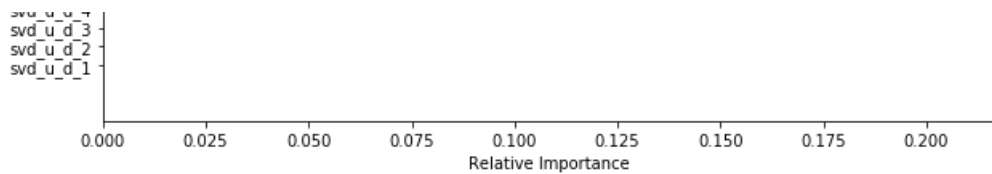
Test confusion_matrix



In [21]:

```
features = df_final_train.columns
importances = clf_xgb.feature_importances_
indices = (np.argsort(importances))[-50:]
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()
```





In [41]:

```
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["MODEL", "train f1-score", "test f1-score"]
x.add_row(["Random Forest", 0.965, 0.924 ])
x.add_row(["Xgboost", 0.875, 0.871])

print(x)
```

MODEL	train f1-score	test f1-score
Random Forest	0.965	0.924
Xgboost	0.875	0.871

Conclusion

1. In this case study test set is created in which there is no connection between edges and labelled as 0
2. The train set has been created in which there exists an edge between the vertices a and b and labelled as 1.
3. The featurization used here is jaccard distance in which the similarity of two sets is measured.
4. Another kind of featurization used is Preferential attachment which takes into consideration the product of two sets.
5. The techniques used here are svd and matrix factorization.
6. The models like Random Forest and Xgboost are used over the designed features.