# Social network Graph Link Prediction - Facebook Challenge

```python
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 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
```

```
In [2]:  #reading
         from pandas import read_hdf
         df_final_train = read_hdf('data/fea_sample/assignment_storage_sample_stage4.h5', 'train_df_assignment',mode=
         'r')
         df_final_test = read_hdf('data/fea_sample/assignment_storage_sample_stage4.h5', 'test_df_assignment',mode='r'
         )
```

```
In [3]:  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_dot_source', 'svd_dot_destination',
                'pa_score_followers', 'pa_score_followee', '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')
```
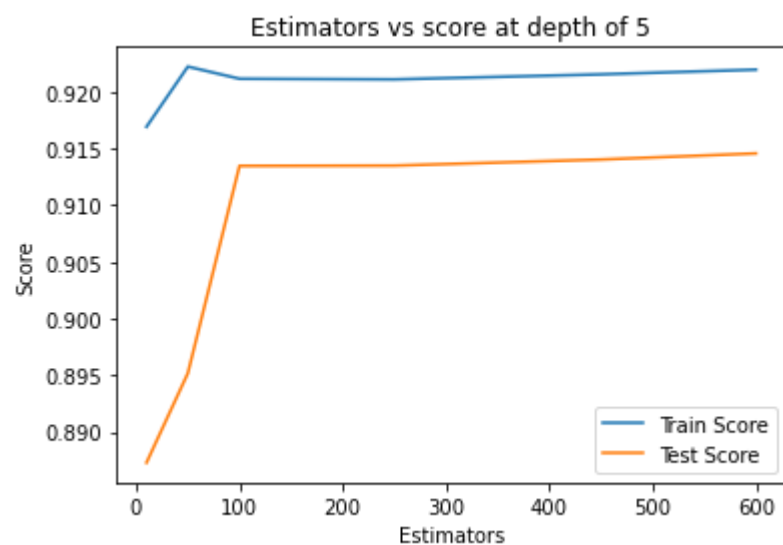
```
In [4]:  y_train = df_final_train.indicator_link
         y_test = df_final_test.indicator_link
```

```
In [5]:  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 [6]:
```python
estimators = [10,50,100,250,450,600]
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')
plt.legend()
```
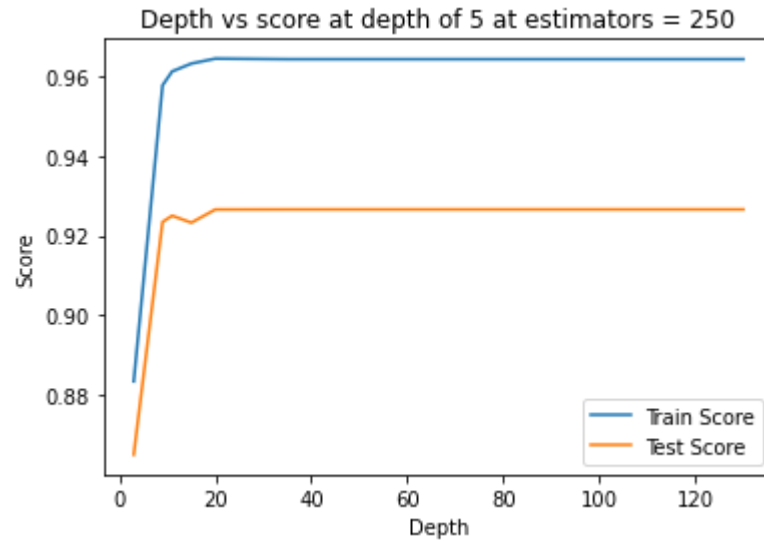
```
Estimators =  10 Train Score 0.9169208195776346 test Score 0.8872915097385129
Estimators =  50 Train Score 0.922227095211244 test Score 0.895179408177675
Estimators =  100 Train Score 0.9211555760098737 test Score 0.9134680134680134
Estimators =  250 Train Score 0.9211005335286118 test Score 0.913495537969355
Estimators =  450 Train Score 0.9215436101234051 test Score 0.914030767935686
Estimators =  600 Train Score 0.9219573400250941 test Score 0.9145694806834441
```

Out[6]: <matplotlib.legend.Legend at 0x221e3afd340>

In [7]:
```python
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=250, 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 = 250')
plt.legend()
plt.show()
```

```
depth =   3 Train Score 0.883489516902611 test Score 0.8651336233685519
depth =   9 Train Score 0.9575928135164622 test Score 0.9233913446546711
depth =  11 Train Score 0.9610838741176232 test Score 0.924957326196447
depth =  15 Train Score 0.9630464127312771 test Score 0.9232099911396143
depth =  20 Train Score 0.9642900571521219 test Score 0.926512332050716
depth =  35 Train Score 0.9641423764543723 test Score 0.9265349032800672
depth =  50 Train Score 0.9641423764543723 test Score 0.9265349032800672
depth =  70 Train Score 0.9641423764543723 test Score 0.9265349032800672
depth = 130 Train Score 0.9641423764543723 test Score 0.9265349032800672
```

In [11]:
```python
%%time

from sklearn.metrics import f1_score
from sklearn.ensemble import RandomForestClassifier
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.return_train_score=True # to get mean train score

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 [0.96179207 0.96175281 0.96017323 0.96171713 0.96311009]
mean train scores [0.9626882  0.96240305 0.96096108 0.96238224 0.96413017]
Wall time: 9min 7s
```

In [12]:
```python
print(rf_random.best_estimator_)
```

```
RandomForestClassifier(max_depth=14, min_samples_leaf=28, min_samples_split=111,
                       n_estimators=121, n_jobs=-1, random_state=25)
```

In [13]:
```python
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 [14]:  clf.fit(df_final_train,y_train)
          y_train_pred = clf.predict(df_final_train)
          y_test_pred = clf.predict(df_final_test)
```

```
In [15]:  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 0.964445347419748
Test f1 score 0.9266822617015448
```

```python
In [16]:   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()
```
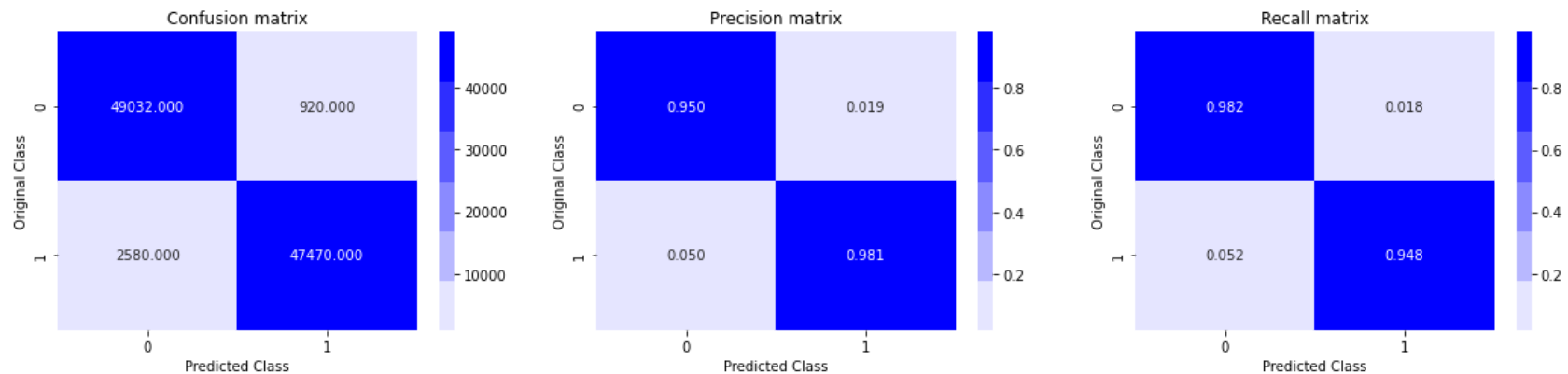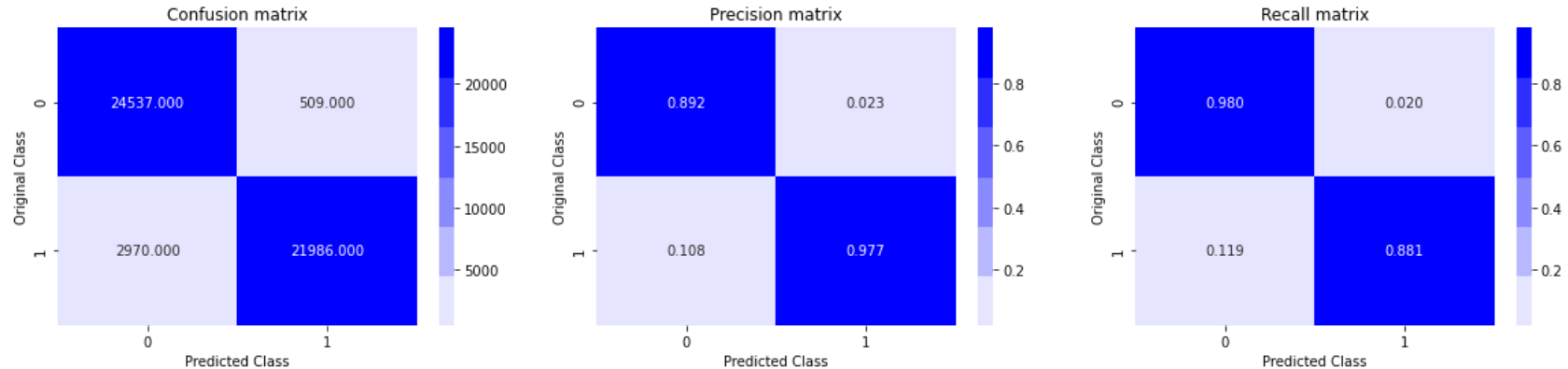
In [17]:
```python
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```
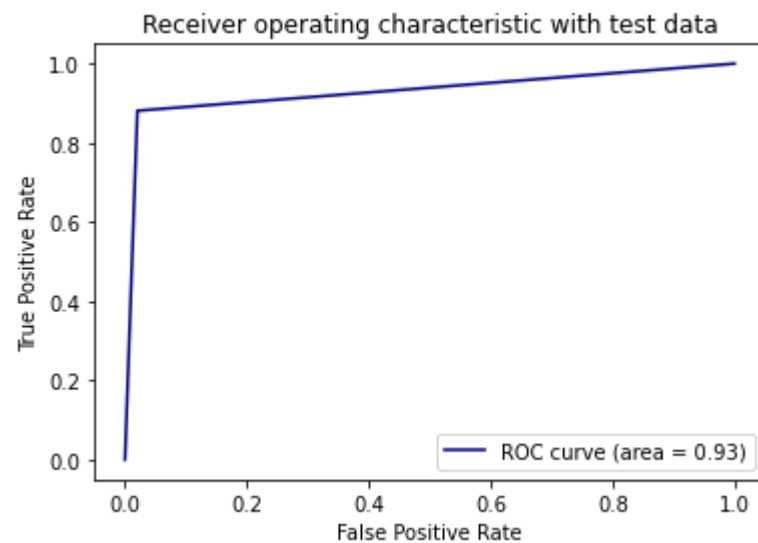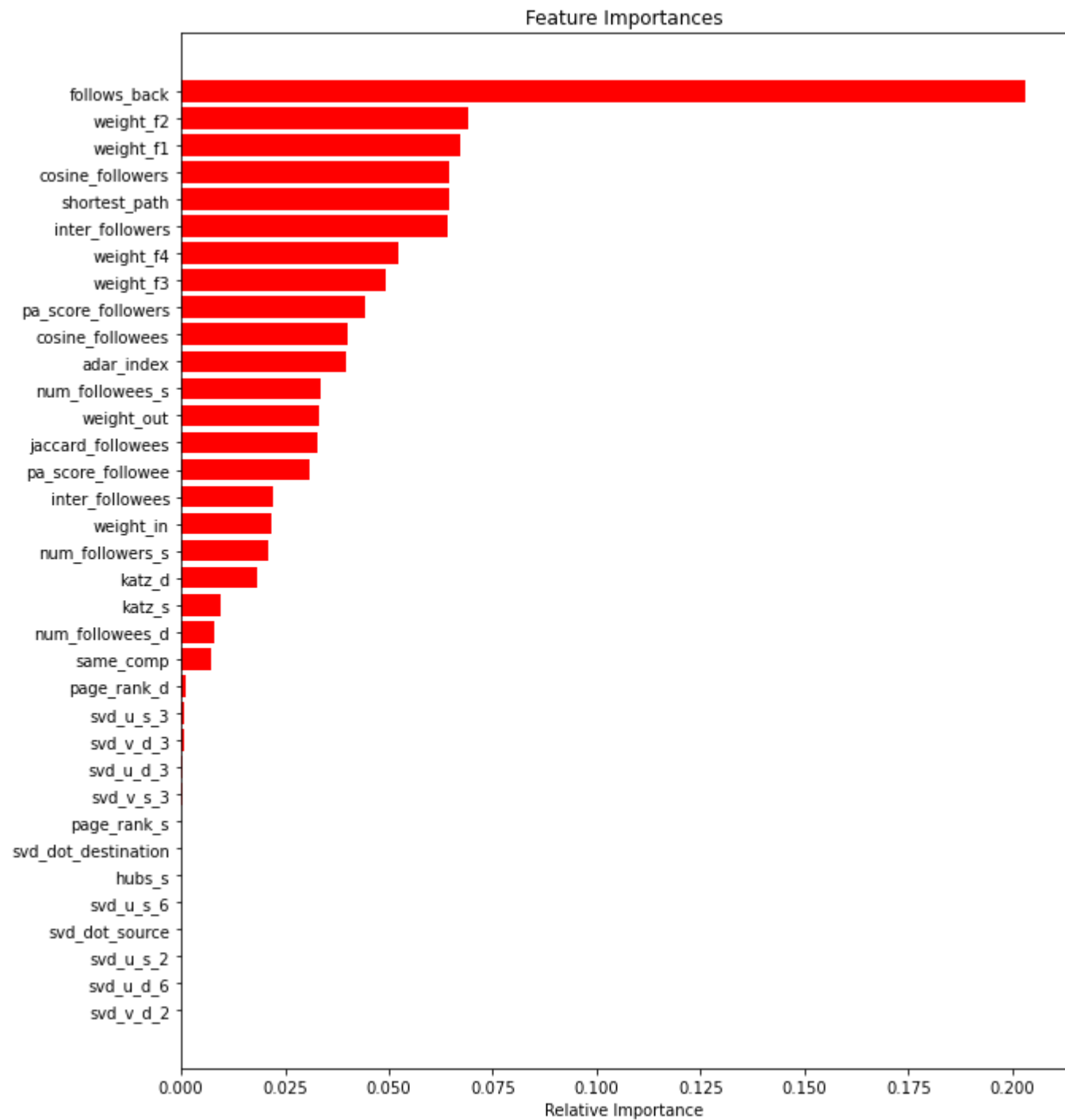
Train confusion_matrix



Test confusion_matrix

```
In [18]: 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 [21]:
```python
features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-35:]
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

## Summary:

1. Preferntial attachment is the 2nd most important feature as we can see in the plot

2. svd_dot feature has no more importance.

In [ ]:

In [ ]: