

# 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 [27]:

```
#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 [16]:

```
temp_train = pd.read_csv('file1.csv')
temp_test = pd.read_csv('file2.csv')
```

In [ ]:

```
temp_train = temp_train.drop(['Unnamed: 0'], axis=1)
temp_test = temp_test.drop(['Unnamed: 0'], axis=1)
```

In [26]:

```
temp_train.head(5)
```

Out[26]:

	PA	svd_dot_u	svd_dot_v
0	120	1.114958e-11	2.238775e-12
1	0	3.192812e-03	9.068719e-04

2	400A	1.797503e-35	2.467873e-86
3	1081	4.710376e-20	3.159386e-18
4	476	7.773952e-14	0.000000e+00

In [28]:

```
df_final_train['PA'] = temp_train['PA'].values
df_final_train['svd_dot_u'] = temp_train['svd_dot_u'].values
df_final_train['svd_dot_v'] = temp_train['svd_dot_v'].values

df_final_test['PA'] = temp_test['PA'].values
df_final_test['svd_dot_u'] = temp_test['svd_dot_u'].values
df_final_test['svd_dot_v'] = temp_test['svd_dot_v'].values
```

In [29]:

```
df_final_train.columns
```

Out[29]:

```
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', 'PA',
      'svd_dot_u'],
      dtype='object')
```

In [30]:

```
df_final_train.head(5)
```

Out[30]:

	source_node	destination_node	indicator_link	jaccard_followers	jaccard_followees	cosine_followers	cosine_followees
0	273084	1505602	1	0	0.000000	0.000000	0.000000
1	832016	1543415	1	0	0.187135	0.028382	0.343828
2	1325247	760242	1	0	0.369565	0.156957	0.566038
3	1368400	1006992	1	0	0.000000	0.000000	0.000000
4	140165	1708748	1	0	0.000000	0.000000	0.000000

5 rows × 56 columns



In [31]:

```
y_train = df_final_train.indicator_link
y_test = df_final_test.indicator_link
```

In [32]:

```
df_final_train.drop(['source_node', 'destination_node', 'indicator_link'], axis=1, inplace=True)
```

```
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)
```

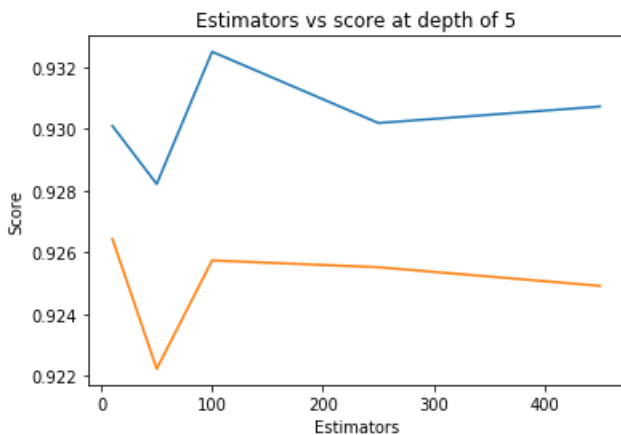
In [33]:

```
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 0.930099471427977 test Score 0.9264429305805156
Estimators = 50 Train Score 0.9282180547856373 test Score 0.9222374381122933
Estimators = 100 Train Score 0.9324989570296203 test Score 0.9257450749284392
Estimators = 250 Train Score 0.9301916799196451 test Score 0.9255258669698692
Estimators = 450 Train Score 0.9307253463732682 test Score 0.9249226494916966
```

Out[33]:

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

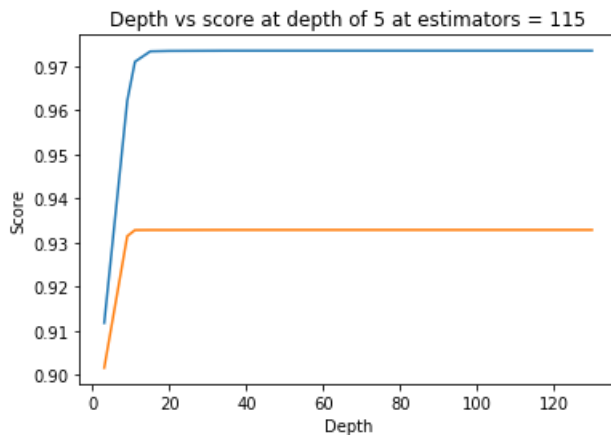


In [34]:

```
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,war
m_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 0.911745385968258 test Score 0.9015484119675325
depth = 9 Train Score 0.9622319961629913 test Score 0.931466470154753
depth = 11 Train Score 0.9710608977606647 test Score 0.9328281233554362
depth = 15 Train Score 0.9734193444178604 test Score 0.932851500126061
depth = 20 Train Score 0.9735144140960266 test Score 0.932848678404841
depth = 35 Train Score 0.9735585417318702 test Score 0.9328710999054524
depth = 50 Train Score 0.9735585417318702 test Score 0.9328710999054524
depth = 70 Train Score 0.9735585417318702 test Score 0.9328710999054524
depth = 130 Train Score 0.9735585417318702 test Score 0.9328710999054524
```



In [48]:

```
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'])
```

In [36]:

```
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 [37]:

```
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 [38]:

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

In [39]:

```
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 0.9752543193928629  
Test f1 score 0.9342074849229863

In [40]:

```
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")

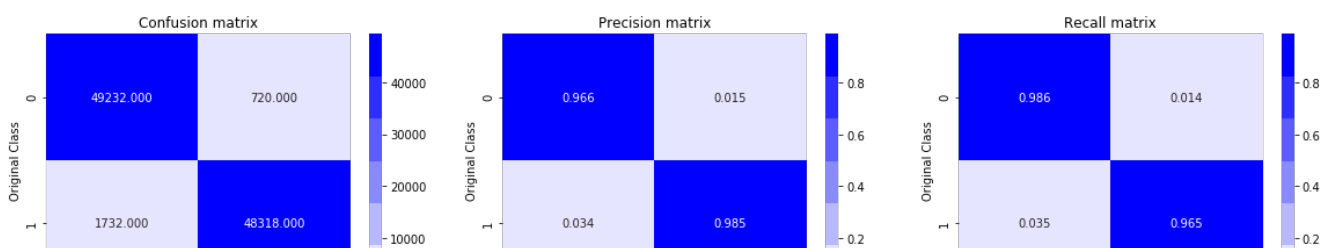
    plt.subplot(1, 3, 3)
    # representing B in heatmap format
    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Recall matrix")

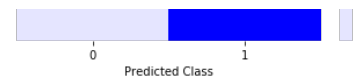
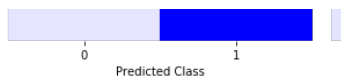
    plt.show()
```

In [41]:

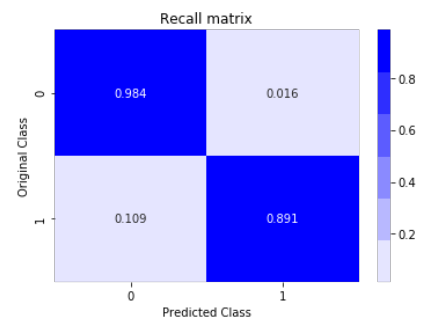
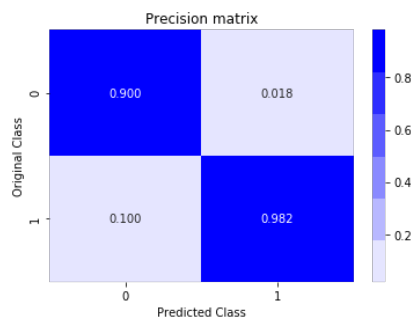
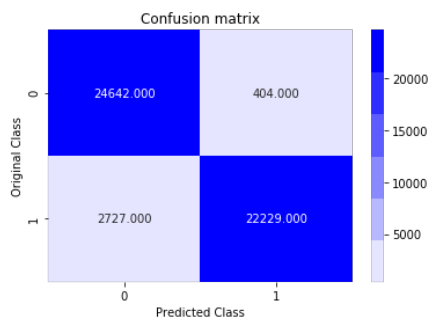
```
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

Train confusion\_matrix



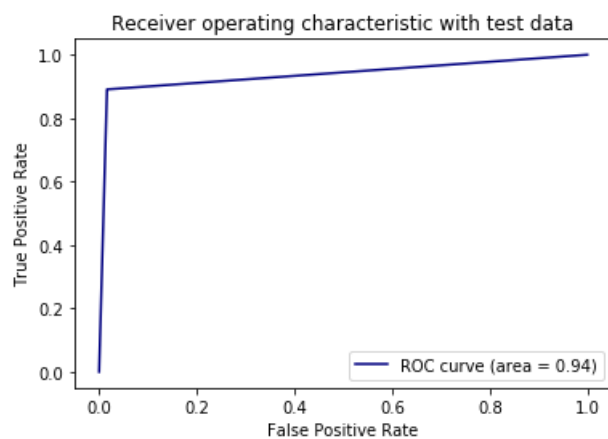


Test confusion\_matrix



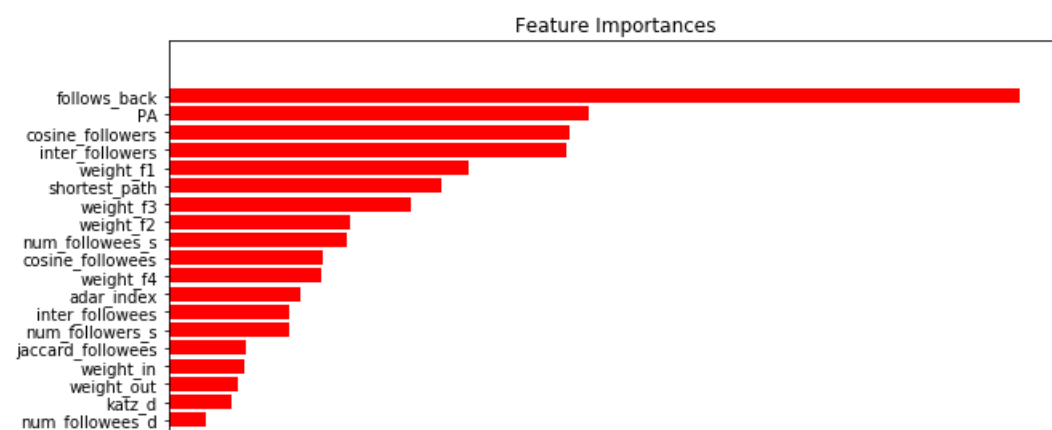
In [42]:

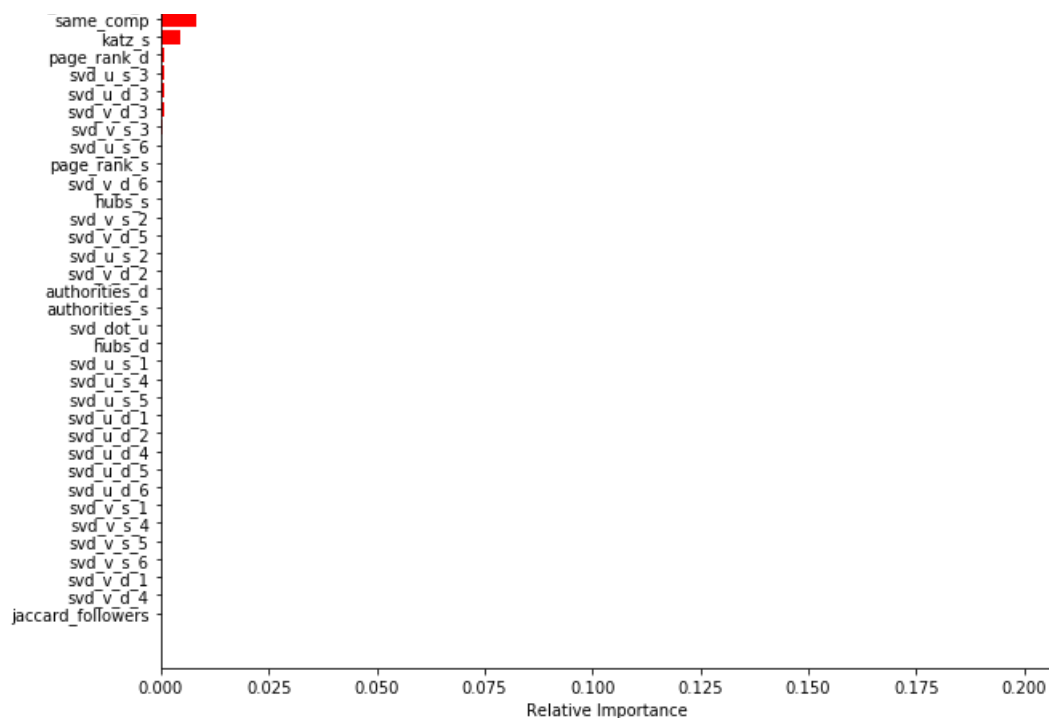
```
from sklearn.metrics import roc_curve, auc
fpr, tpr, ths = roc_curve(y_test, y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy', label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic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



In [47]:

```
features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-53:]
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 [50]:

```
import xgboost as xgb
x_cfl=xgb.XGBClassifier()

prams={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
    'n_estimators':[100,200,500,1000,2000],
    'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jobs=-1,)
random_cfl.fit(df_final_train,y_train)
```

Fitting 3 folds for each of 10 candidates, totalling 30 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done   5 tasks      | elapsed:   7.3min
[Parallel(n_jobs=-1)]: Done  10 tasks      | elapsed:  13.4min
[Parallel(n_jobs=-1)]: Done  17 tasks      | elapsed:  33.6min
[Parallel(n_jobs=-1)]: Done  27 out of  30 | elapsed:  47.5min remaining:   5.3min
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed:  54.6min finished
```

Out[50]:

```
RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
                  estimator=XGBClassifier(base_score=0.5, booster='gbtree',
                                          colsample_bylevel=1,
                                          colsample_bynode=1,
                                          colsample_bytree=1, gamma=0,
                                          learning_rate=0.1, max_delta_step=0,
                                          max_depth=3, min_child_weight=1,
                                          missing=None, n_estimators=100,
                                          n_jobs=1, nthread=None,
                                          objective='binary:logistic',
                                          random_state=0, reg_alpha=0.1,
                                          seed=None, silent=None, subsample=1,
                                          verbosity=1),
                  iid='warn', n_iter=10, n_jobs=-1,
                  param_distributions={'colsample_bytree': [0.1, 0.3, 0.5, 1],
                                      'learning_rate': [0.01, 0.03, 0.05, 0.1,
                                                         0.15, 0.2],
                                      'max_depth': [3, 5, 10],
                                      'n_estimators': [100, 200, 500, 1000,
                                                         2000],
                                      'subsample': [0.1, 0.3, 0.5, 1]},
                  pre_dispatch='2*n_jobs', random_state=None, refit=True,
                  return_train_score=False, scoring=None, verbose=10)
```

```
return_train_score=False, scoring=None, verbose=10,
```

In [52]:

```
clf = xgb.XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                        colsample_bynode=1, colsample_bytree=0.5, gamma=0,
                        learning_rate=0.2, max_delta_step=0, max_depth=5,
                        min_child_weight=1, missing=None, n_estimators=1000, n_jobs=-1,
                        nthread=None, objective='binary:logistic', random_state=0,
                        reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                        silent=None, subsample=1, verbosity=1)
```

In [51]:

```
print(random_cfl.best_estimator_)
```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=0.5, gamma=0,
              learning_rate=0.2, max_delta_step=0, max_depth=5,
              min_child_weight=1, missing=None, n_estimators=1000, n_jobs=1,
              nthread=None, objective='binary:logistic', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=None, subsample=1, verbosity=1)
```

In [53]:

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

In [54]:

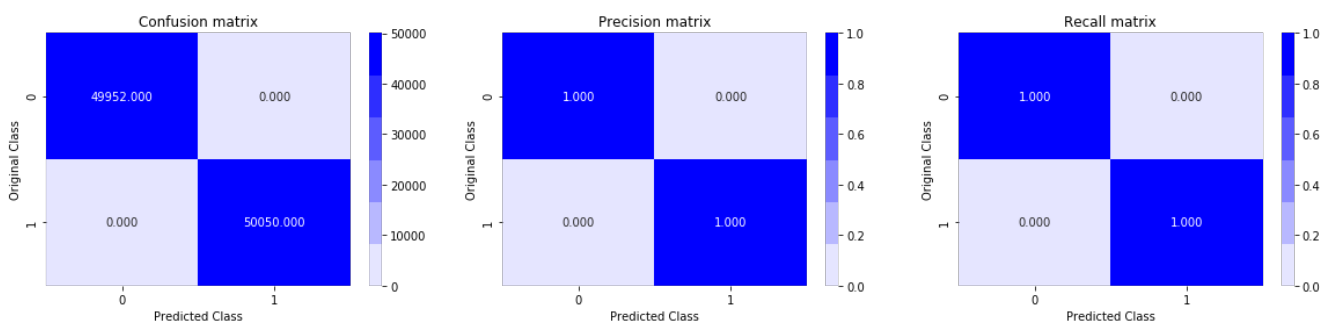
```
from sklearn.metrics import f1_score
print('Train f1 score',f1_score(y_train,y_train_pred))
print('Test f1 score',f1_score(y_test,y_test_pred))
```

```
Train f1 score 1.0
Test f1 score 0.9086618605356023
```

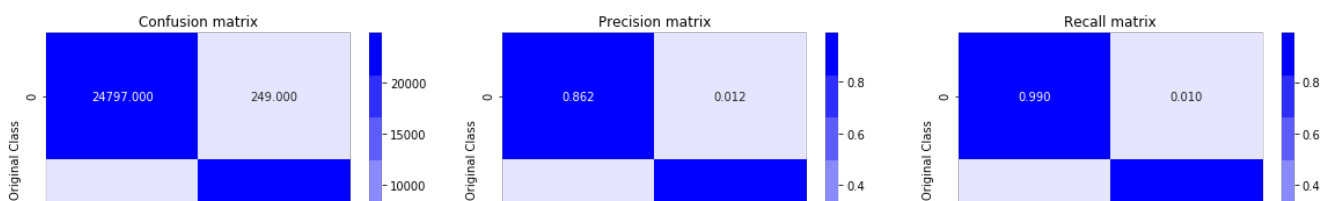
In [55]:

```
print('Train confusion_matrix')
plot_confusion_matrix(y_train,y_train_pred)
print('Test confusion_matrix')
plot_confusion_matrix(y_test,y_test_pred)
```

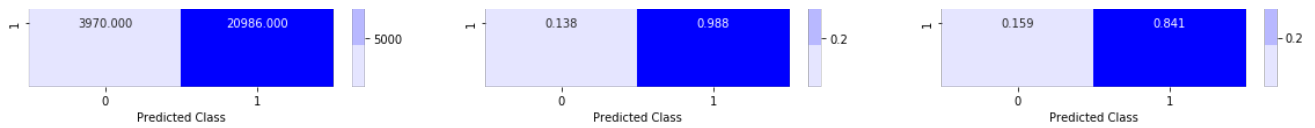
Train confusion\_matrix



Test confusion\_matrix

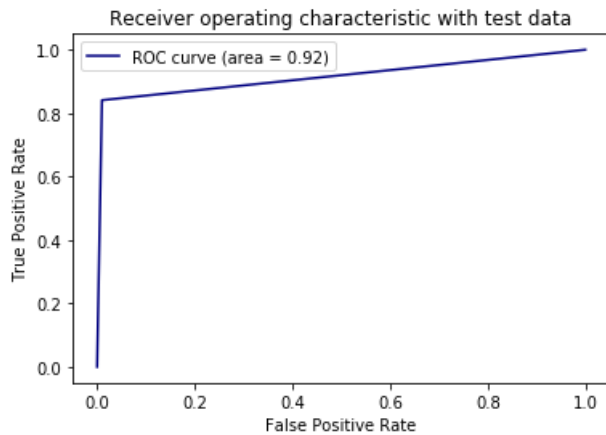






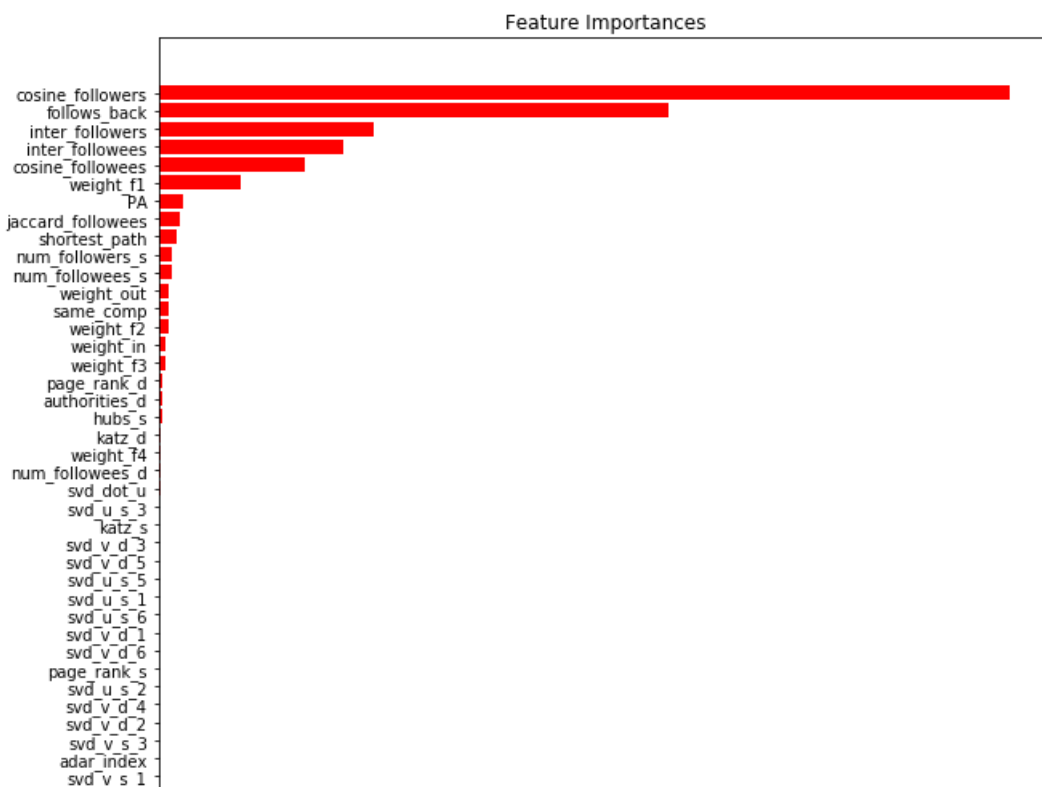
In [56]:

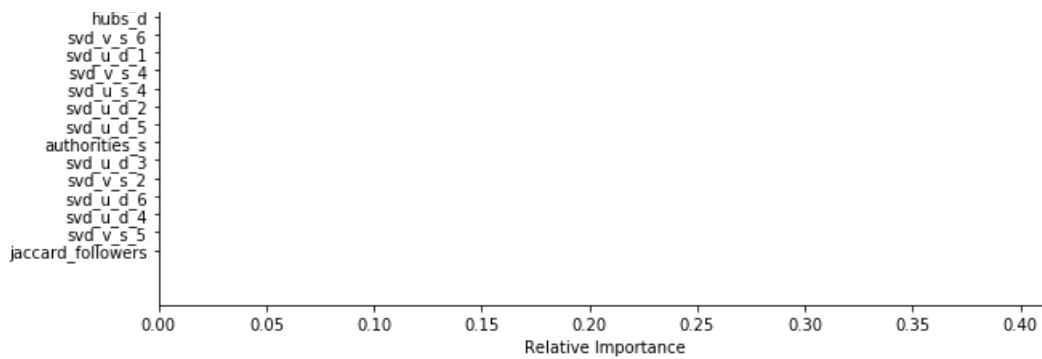
```
from sklearn.metrics import roc_curve, auc
fpr, tpr, ths = roc_curve(y_test, y_test_pred)
auc_sc = auc(fpr, tpr)
plt.plot(fpr, tpr, color='navy', label='ROC curve (area = %0.2f)' % auc_sc)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic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



In [57]:

```
features = df_final_train.columns
importances = clf.feature_importances_
indices = (np.argsort(importances))[-53:]
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()
```





## XG-BOOST is overfitting on Dtrain.

In [59]:

```
from prettytable import PrettyTable
```

```
x = PrettyTable()
x.field_names = ["Model", "ROC", "F1-Score"]
x.add_row(["Random Forest", 0.94, 0.93])
x.add_row(["XG Boost", 0.92, 0.90])
print(x)
```

```
+-----+-----+-----+
|      Model      | ROC  | F1-Score |
+-----+-----+-----+
| Random Forest   | 0.94 | 0.93     |
| XG Boost        | 0.92 | 0.9      |
+-----+-----+-----+
```

## 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](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.