# Social network Graph Link Prediction - Facebook Challenge

In [2]:
```python
#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 [7]:    1  #reading
           2  from pandas import read_hdf
           3  df_final_train = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'train_df',mode='r')
           4  df_final_test = read_hdf('data/fea_sample/storage_sample_stage4.h5', 'test_df',mode='r')
```

```
In [17]:   1  y_train = df_final_train.indicator_link
           2  y_test = df_final_test.indicator_link
```

```
In [ ]:    1
```

```
In [18]:   1  df_final_train.drop(['source_node', 'destination_node','indicator_link'],axis=1,inplace=True)
           2  df_final_test.drop(['source_node', 'destination_node','indicator_link'],axis=1,inplace=True)
```

```python
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=Fa
    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')
```
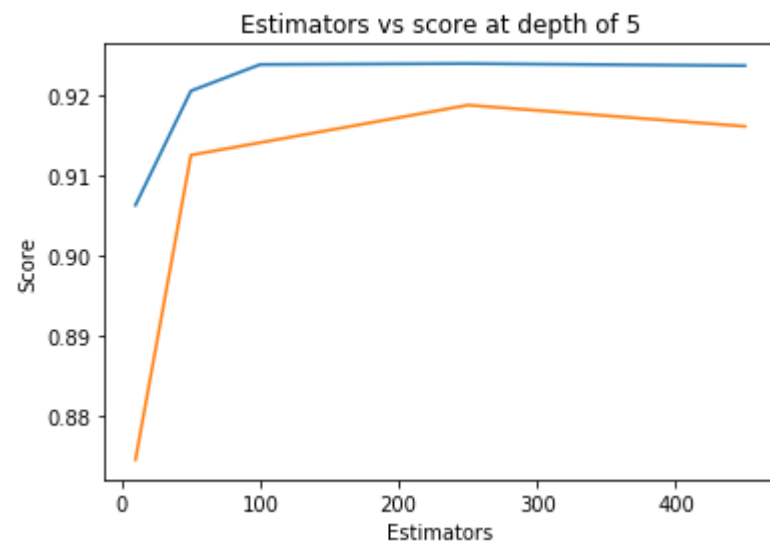
```
Estimators =  10 Train Score 0.9063252121775113 test Score 0.8745605278006858
Estimators =  50 Train Score 0.9205725512208812 test Score 0.9125653355634538
Estimators =  100 Train Score 0.9238690848446947 test Score 0.9141199714153599
Estimators =  250 Train Score 0.9239789348046863 test Score 0.9188007232664732
Estimators =  450 Train Score 0.9237190618658074 test Score 0.9161507685828595
```

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

Estimators vs score at depth of 5
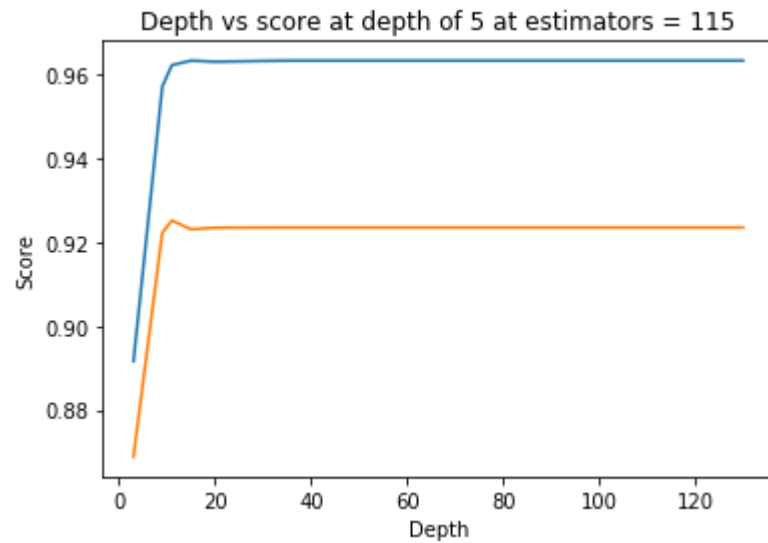
```
In [0]:  1  depths = [3,9,11,15,20,35,50,70,130]
         2  train_scores = []
         3  test_scores = []
         4  for i in depths:
         5      clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
         6              max_depth=i, max_features='auto', max_leaf_nodes=None,
         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_state=25,verbose=0,warm_start=
        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
```

Depth vs score at depth of 5 at estimators = 115



```
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=25)
          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]:
```python
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]:
```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 [0]:
```python
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]:
```python
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
```
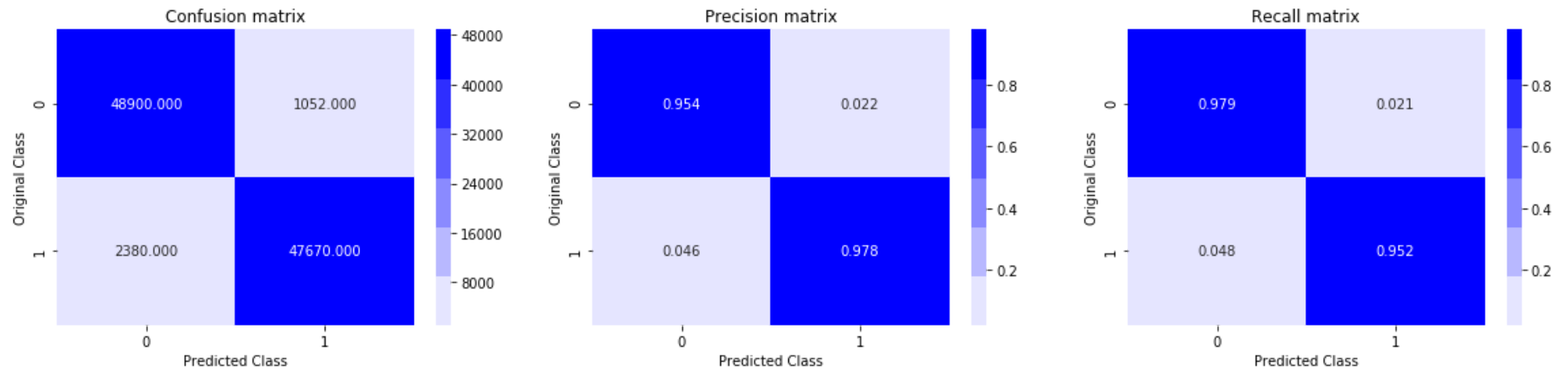
```
In [49]:   1  from sklearn.metrics import confusion_matrix
           2  def plot_confusion_matrix(test_y, predict_y):
           3      C = confusion_matrix(test_y, predict_y)
           4
           5      A =(((C.T)/(C.sum(axis=1))).T)
           6
           7      B =(C/C.sum(axis=0))
           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, yticklabels=labels)
          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, yticklabels=labels)
          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, yticklabels=labels)
          28      plt.xlabel('Predicted Class')
          29      plt.ylabel('Original Class')
          30      plt.title("Recall matrix")
          31
          32      plt.show()
```
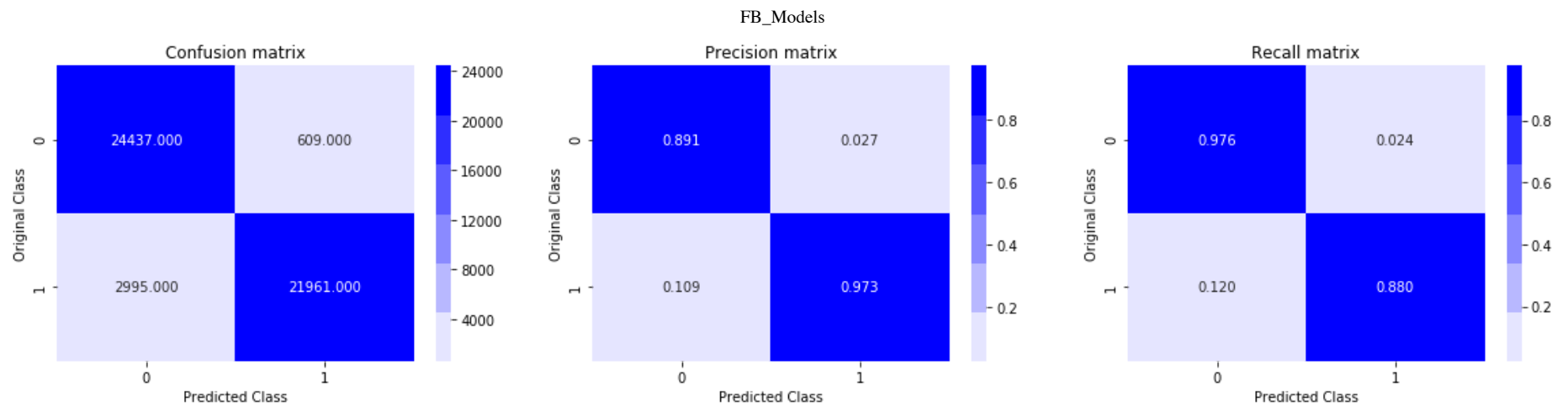
```
In [0]:    1  print('Train confusion_matrix')
           2  plot_confusion_matrix(y_train,y_train_pred)
           3  print('Test confusion_matrix')
           4  plot_confusion_matrix(y_test,y_test_pred)
```
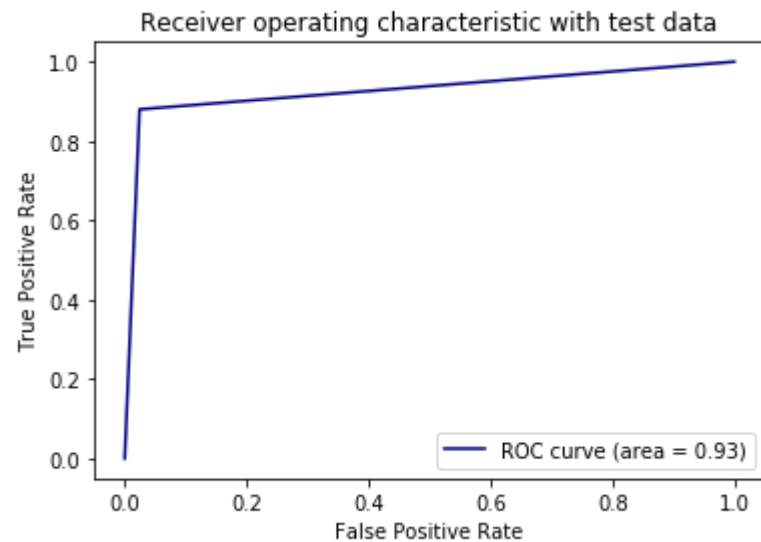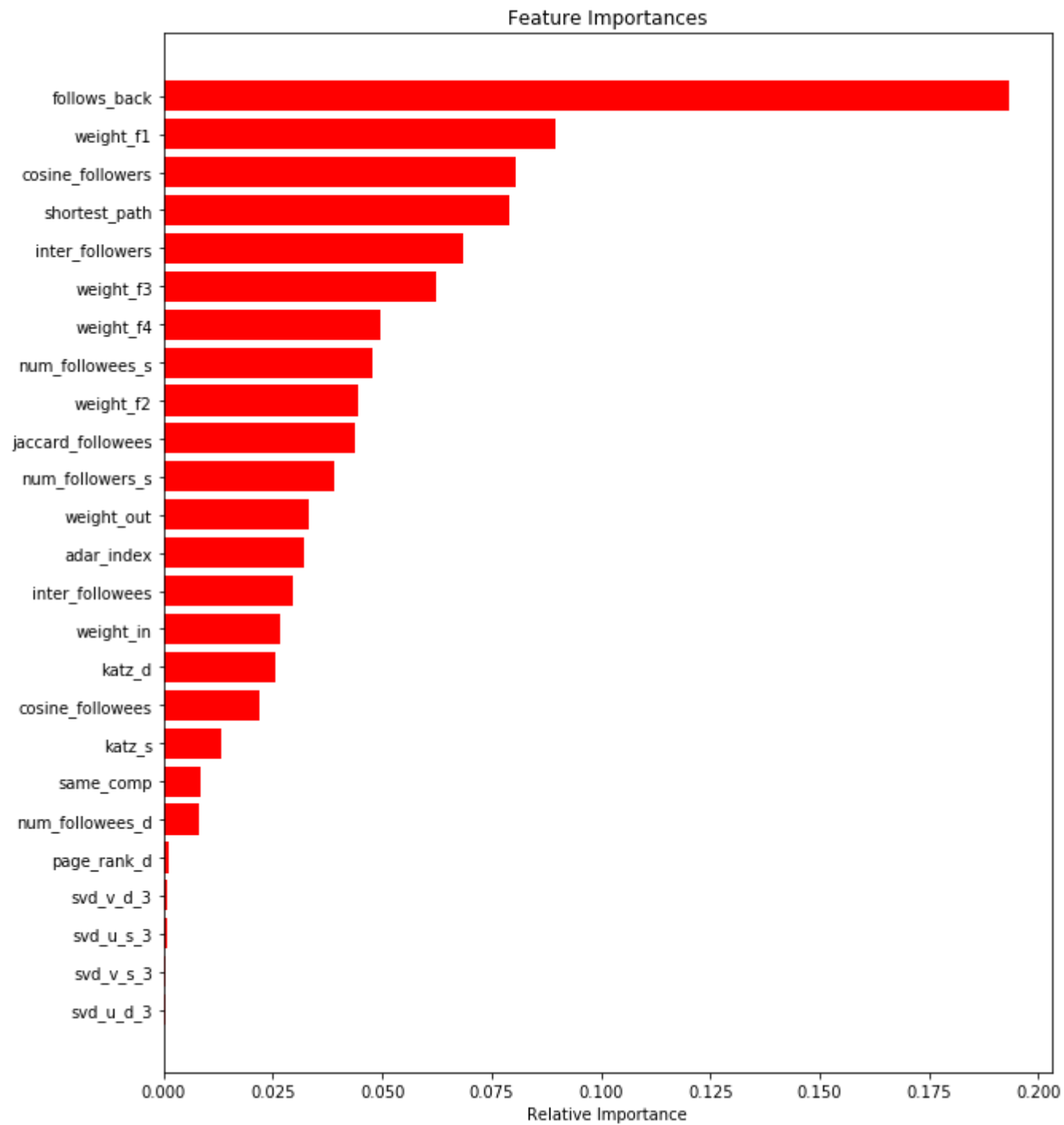
Train confusion_matrix



Test confusion_matrix

### Confusion matrix

|   | 0 | 1 |
|---|---|---|
| 0 | 24437.000 | 609.000 |
| 1 | 2995.000 | 21961.000 |

### Precision matrix

|   | 0 | 1 |
|---|---|---|
| 0 | 0.891 | 0.027 |
| 1 | 0.109 | 0.973 |

### Recall matrix

|   | 0 | 1 |
|---|---|---|
| 0 | 0.976 | 0.024 |
| 1 | 0.120 | 0.880 |

In [0]:
```python
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()
```



Receiver operating characteristic with test data — ROC curve (area = 0.93)

```
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()
```
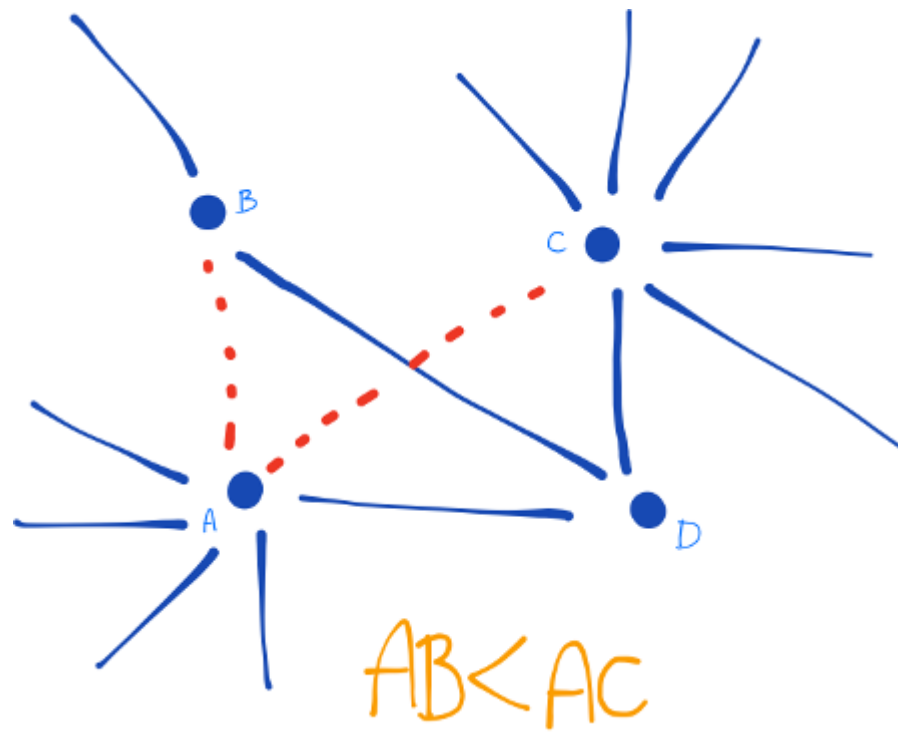
Feature Importances

# 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 sourse 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.

# Preferential Attachment :

One well-known concept in social networks is that users with many friends tend to create more connections in the future. This is due to the fact that in some social networks, like in finance, the rich get richer. We estimate how "rich" our two vertices are by calculating the multiplication between the number of friends ($|\Gamma(x)|$) or followers each vertex has. It may be noted that the similarity index does not require any node neighbor information; therefore, this similarity index has the lowest computational complexity.

In [19]:
```python
1  df_final_train.columns
2  # df_final_train["num_followees_d"]
```

Out[19]: 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',
        'num_followers_d'],
       dtype='object')

In [12]:
```python
1  # Kind of like friends of friends and mutual friends
2  # Perferntial Attachment is two vertices are by calculating the multiplication
3  # between the number of friends (|Γ(x)|) or followers each vertex has.
4  # So Source and Destination Followees multiplication
```

In [20]:
```
1  df_final_train['perferntial_attach'] = df_final_train["num_followees_s"] * df_final_train["num_followees_d"
2  df_final_train['perferntial_attach_follower'] = df_final_train["num_followers_s"] * df_final_train["num_fol
3
4
5  df_final_train.head()
```

Out[20]:

| 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 | num_followers_d | perferntial_attach | perferntial_attach |
|---|---|---|---|---|---|---|---|---|---|
| 1.719702e-14 | -1.355368e-12 | 4.675307e-13 | 1.128591e-06 | 6.616550e-14 | 9.771077e-13 | 4.159752e-14 | 6 | 120 | |
| 2.251737e-10 | 1.245101e-12 | -1.636948e-10 | -3.112650e-10 | 6.738902e-02 | 2.607801e-11 | 2.372904e-09 | 94 | 8662 | |
| 3.365389e-19 | -1.238370e-18 | 1.438175e-19 | -1.852863e-19 | -5.901864e-19 | 1.629341e-19 | -2.572452e-19 | 28 | 902 | |
| 4.498061e-13 | -9.818087e-10 | 3.454672e-11 | 5.213635e-08 | 9.595823e-13 | 3.047045e-10 | 1.246592e-13 | 11 | 35 | |
| 1.407670e-14 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 1 | 33 | |

In [21]:
```
1  df_final_test['perferntial_attach'] = df_final_test["num_followees_s"] * df_final_test["num_followees_d"]
2  df_final_test['perferntial_attach_follower'] = df_final_test["num_followers_s"] * df_final_test["num_follow
3  df_final_test.head()
```

Out[21]:

| inter_followees | adar_index | ... | 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 | num_followers_d | perferntial_attach |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.000000 | ... | 5.535503e-14 | -9.994076e-10 | 5.791910e-10 | 3.512364e-07 | 2.486658e-09 | 2.771146e-09 | 1.727694e-12 | 14 | 54 |
| 0 | 0.000000 | ... | 4.701436e-15 | -9.360516e-12 | 3.206809e-10 | 4.668696e-08 | 6.665777e-12 | 1.495979e-10 | 9.836670e-14 | 17 | 19 |
| 0 | 0.000000 | ... | 4.199834e-14 | -4.253075e-13 | 4.789463e-13 | 3.479824e-07 | 1.630549e-13 | 3.954708e-13 | 3.875785e-14 | 10 | 144 |
| 0 | 0.000000 | ... | 2.817657e-13 | -2.162590e-11 | 6.939194e-12 | 1.879861e-05 | 4.384816e-12 | 1.239414e-11 | 6.483485e-13 | 37 | 340 |
| 7 | 6.136433 | ... | 9.656662e-14 | -8.742904e-12 | 7.467370e-12 | 1.256880e-05 | 3.636983e-12 | 3.948463e-12 | 2.415863e-13 | 27 | 405 |

## Adding SVD Dot Feature that is product of source and destination SVDs

In [ ]:
```
1  # Dot product of Source and Destination Node SVD
```

In [25]:
```
1  source = ['svd_u_s_1', 'svd_u_s_2','svd_u_s_3', 'svd_u_s_4', 'svd_u_s_5', 'svd_u_s_6']
2  source_v = ['svd_v_s_1','svd_v_s_2', 'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6']
3  destination = ['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3', 'svd_u_d_4', 'svd_u_d_5','svd_u_d_6']
4  destination_v = ['svd_v_d_1', 'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5','svd_v_d_6']
5
```

```
In [22]:    1  # source_array = []
            2  # destination_array = []
            3  # svd_dot = []
            4
            5  # for num in tqdm(range(df_final_train.shape[0])):
            6  #      for i,j in (zip(source,destination)):
            7  #          source_array.append(np.array(df_final_train[i].iloc[num]))
            8  #          destination_array.append(np.array(df_final_train[j].iloc[num]))
            9  #      svd_dot.append(np.dot(source_array, destination_array))
           10  # df_final_train['svd_dot']=svd_dot
           11
```

```
In [32]:    1  # Corrections
            2  f = 0;
            3  for i,j in tqdm(zip(source,destination)):
            4      f = f + df_final_train[i] * df_final_train[j]
            5  df_final_train['svd_dot_source'] = f
            6
            7  f = 0;
            8  for i,j in tqdm(zip(source_v,destination_v)):
            9      f = f + df_final_train[i] * df_final_train[j]
           10  df_final_train['svd_dot_dest'] = f
           11
           12  f = 0;
           13  for i,j in tqdm(zip(source,destination)):
           14      f = f + df_final_test[i] * df_final_test[j]
           15  df_final_test['svd_dot_source'] = f
           16
           17  f = 0;
           18  for i,j in tqdm(zip(source_v,destination_v)):
           19      f = f + df_final_test[i] * df_final_test[j]
           20  df_final_test['svd_dot_dest'] = f
           21
```

```
6it [00:00, 482.16it/s]
6it [00:00, 480.27it/s]
6it [00:00, 623.44it/s]
6it [00:00, 337.19it/s]
```

In [33]: 
```
1  df_final_train.head()
```

Out[33]:

| vd_v_d_2 | svd_v_d_3 | svd_v_d_4 | svd_v_d_5 | svd_v_d_6 | num_followers_d | perferntial_attach | perferntial_attach_follower | svd_dot_source | svd |
|---|---|---|---|---|---|---|---|---|---|
| 307e-13 | 1.128591e-06 | 6.616550e-14 | 9.771077e-13 | 4.159752e-14 | 6 | 120 | 36 | 1.114958e-11 | 2.23 |
| 336948e-10 | -3.112650e-10 | 6.738902e-02 | 2.607801e-11 | 2.372904e-09 | 94 | 8662 | 8836 | 3.192812e-03 | 9.00 |
| 3175e-19 | -1.852863e-19 | -5.901864e-19 | 1.629341e-19 | -2.572452e-19 | 28 | 902 | 784 | 1.787503e-35 | 2.46 |
| 4672e-11 | 5.213635e-08 | 9.595823e-13 | 3.047045e-10 | 1.246592e-13 | 11 | 35 | 121 | 4.710376e-20 | 3.15 |
| 000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 1 | 33 | 1 | 7.773952e-14 | 0.00 |

In [35]: 
```
1  df_final_test.head()
```

Out[35]:

| | jaccard_followers | jaccard_followees | cosine_followers | cosine_followees | num_followers_s | num_followees_s | num_followees_d | inter_followers | inter |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 0.0 | 0.029161 | 0.000000 | 14 | 6 | 9 | 1 | |
| **1** | 0 | 0.0 | 0.000000 | 0.000000 | 17 | 1 | 19 | 0 | |
| **2** | 0 | 0.0 | 0.000000 | 0.000000 | 10 | 16 | 9 | 0 | |
| **3** | 0 | 0.0 | 0.000000 | 0.000000 | 37 | 10 | 34 | 0 | |
| **4** | 0 | 0.2 | 0.042767 | 0.347833 | 27 | 15 | 27 | 4 | |

5 rows × 56 columns

In [ ]: 
```
1
```

# Random Forest

```
In [37]:    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               }
           11
           12  clf = RandomForestClassifier(random_state=25,n_jobs=-1)
           13
           14  rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
           15                                 cv=2,scoring='f1')
           16
           17  rf_random.fit(df_final_train,y_train)
           18
```

```
Out[37]: RandomizedSearchCV(cv=2, error_score='raise-deprecating',
                   estimator=RandomForestClassifier(bootstrap=True,
                                                    class_weight=None,
                                                    criterion='gini',
                                                    max_depth=None,
                                                    max_features='auto',
                                                    max_leaf_nodes=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    n_estimators='warn',
                                                    n_jobs=-1, oob_scor...
                                                    random_state=25, verbose=0,
                                                    warm_start=False),
                   iid='warn', n_iter=10, n_jobs=None,
                   param_distributions={'max_depth': <scipy.stats._distn_infrastructure.rv_frozen object at 0
         x1a3f8af7f0>,
                                        'n_estimators': <scipy.stats._distn_infrastructure.rv_frozen object a
         t 0x1a3f9455c0>},
```

```
                         pre_dispatch='2*n_jobs', random_state=None, refit=True,
                         return_train_score=False, scoring='f1', verbose=0)
```

In [38]:
```
1  print('mean test scores',rf_random.cv_results_['mean_test_score'])
```

```
mean test scores [0.9663106  0.96309182 0.96826787 0.9609609  0.96657034 0.96474245
 0.96476606 0.96826012 0.96827655 0.96651507]
```

In [39]:
```
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=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, n_estimators=114,
                       n_jobs=-1, oob_score=False, random_state=25, verbose=0,
                       warm_start=False)
```

In [44]:
```
1  n_estimator_rf = 114
2  max_depth_rf = 14
3  clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
4                               max_depth=14, max_features='auto', max_leaf_nodes=None,
5                               min_impurity_decrease=0.0, min_impurity_split=None,
6                               min_samples_leaf=1, min_samples_split=2,
7                               min_weight_fraction_leaf=0.0, n_estimators=114,
8                               n_jobs=-1, oob_score=False, random_state=25, verbose=0,
9                               warm_start=False)
```

In [45]:
```
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)
```
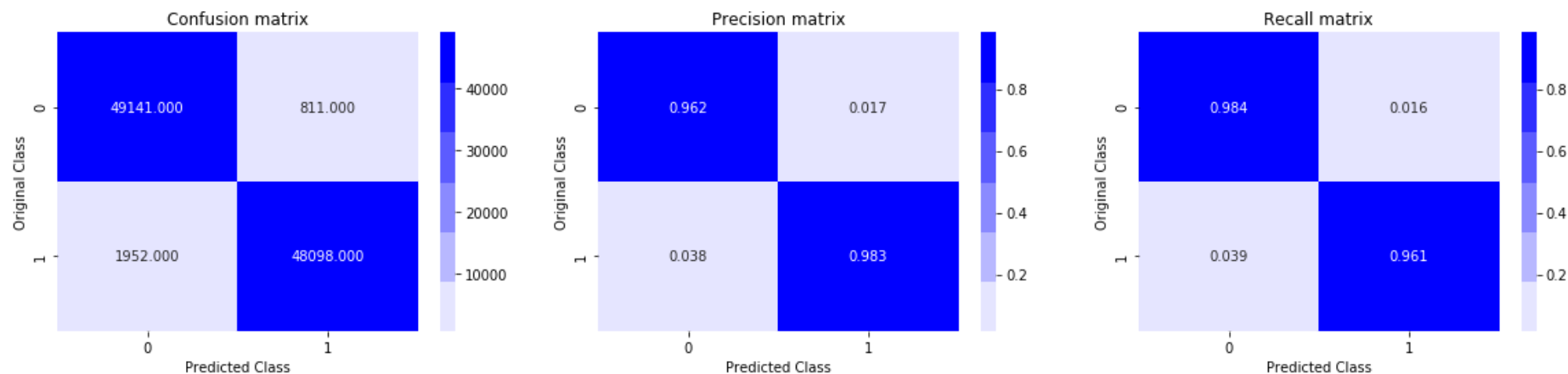
```python
In [46]:   1  from sklearn.metrics import f1_score
           2
           3  Train_f1_Score_rf = f1_score(y_train,y_train_pred)
           4  Test_f1_Score_rf = f1_score(y_test,y_test_pred)
           5
           6  print('Train f1 score',Train_f1_Score_rf)
           7  print('Test f1 score',Test_f1_Score_rf)
```

```
Train f1 score 0.9720793459917744
Test f1 score 0.9188237291012475
```
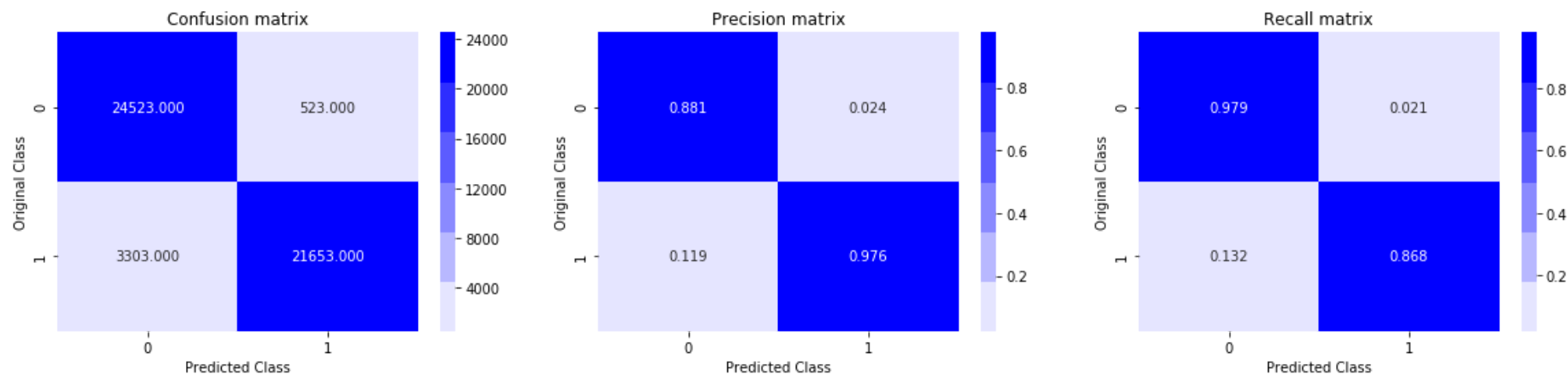
```
In [50]:   1  print('Train confusion_matrix')
           2  plot_confusion_matrix(y_train,y_train_pred)
           3  print('Test confusion_matrix')
           4  plot_confusion_matrix(y_test,y_test_pred)
```

Train confusion_matrix



Test confusion_matrix

```
In [ ]:  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()
```

# XGboost

```
In [40]:  1  from datetime import datetime
          2  from sklearn.metrics import f1_score
          3  from sklearn.model_selection import RandomizedSearchCV
          4  from xgboost import XGBClassifier
```

```
In [42]:   1  start = datetime.now()
           2  clf = XGBClassifier()
           3  params = {"n_estimators":sp_randint(105,125),
           4              "max_depth": sp_randint(10,15)
           5              }
           6  rf = RandomizedSearchCV(clf, param_distributions=params,cv=2,scoring='f1', n_jobs=-1)
           7
           8
           9
          10  rf.fit(df_final_train,y_train)
          11  print('mean test scores',rf.cv_results_['mean_test_score'])
          12  print(datetime.now() - start)
```

```
mean test scores [0.97926699 0.97903336 0.97881123 0.97901685 0.9791557  0.97929291
 0.97896439 0.97933791 0.97906177 0.97911945]
0:14:49.802030
```

In [43]:
```
1  print(rf.best_estimator_)
```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
              max_depth=12, min_child_weight=1, missing=None, n_estimators=118,
              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=1)
```
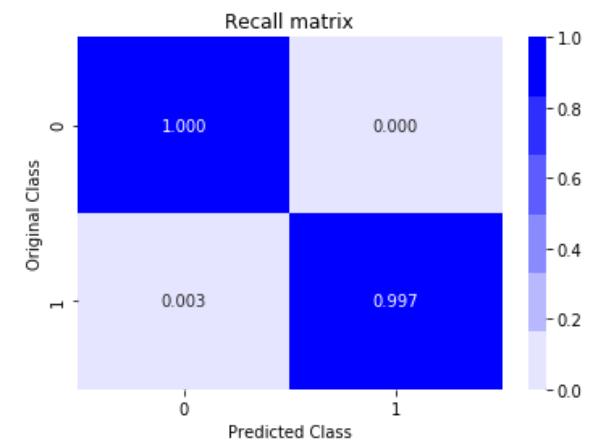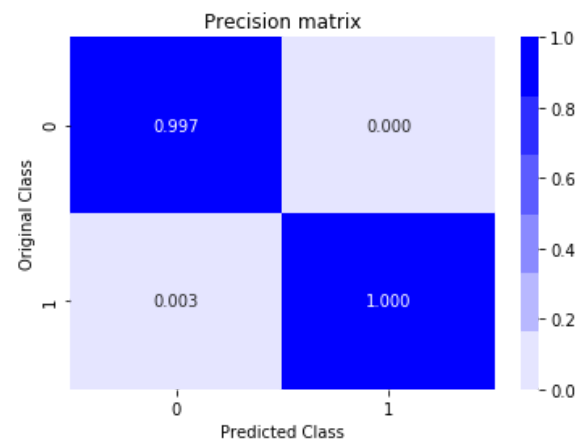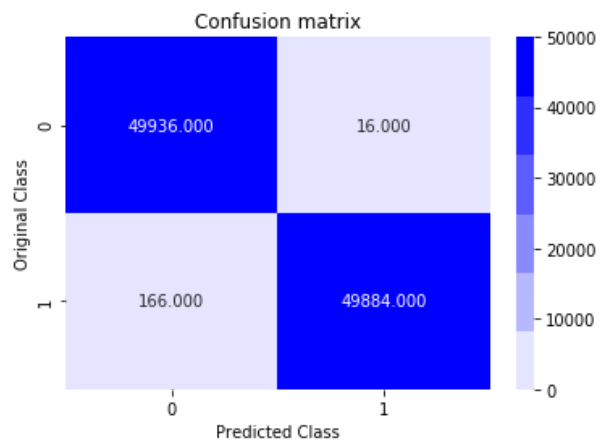
```python
1  n_estimator_xgb = 118
2  max_depth_xgb = 12
3  xgb =  XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
4              colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
5              max_depth=12, min_child_weight=1, missing=None, n_estimators=118,
6              n_jobs=1, nthread=None, objective='binary:logistic',
7              random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
8              seed=None, silent=True, subsample=1)
9  xgb.fit(df_final_train,y_train)
10 y_pred_train = xgb.predict(df_final_train)
11 y_pred_test = xgb.predict(df_final_test)
12
13 Train_f1_Score_xgb = f1_score(y_train,y_pred_train)
14 Test_f1_Score_xgb = f1_score(y_test,y_pred_test)
15
16 print('Train f1 score',Train_f1_Score_xgb)
17 print('Test f1 score',Test_f1_Score_xgb)
18
19 print('Train confusion_matrix')
20 plot_confusion_matrix(y_train,y_pred_train)
21 print('Test confusion_matrix')
22 plot_confusion_matrix(y_test,y_pred_test)
```
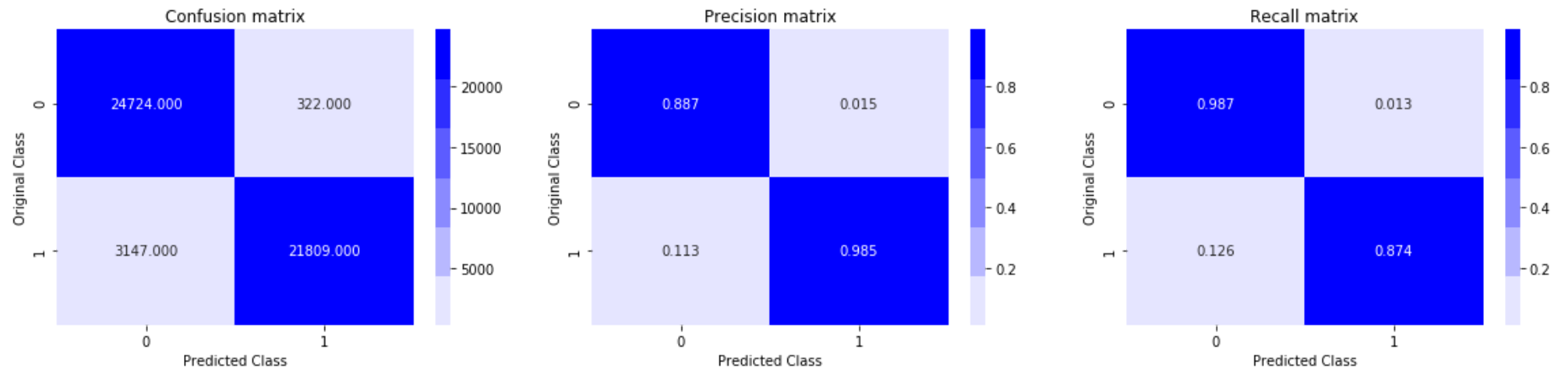
```
Train f1 score 0.9981790895447723
Test f1 score 0.9263278611931106
Train confusion_matrix
```
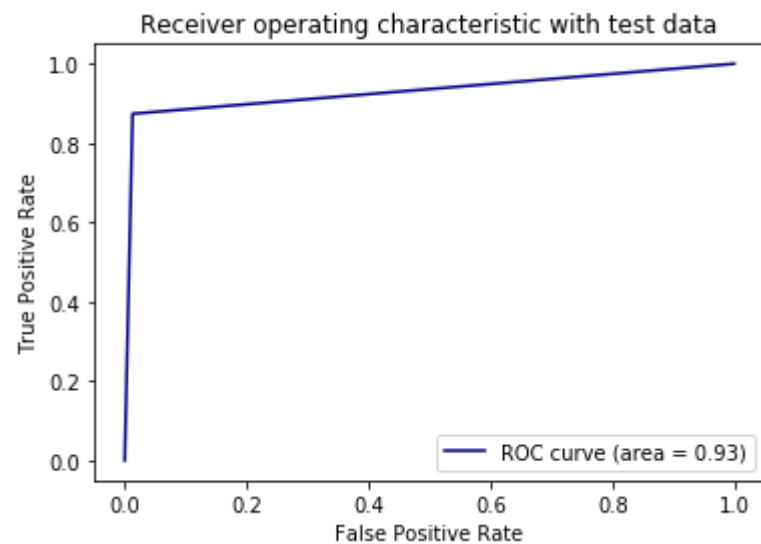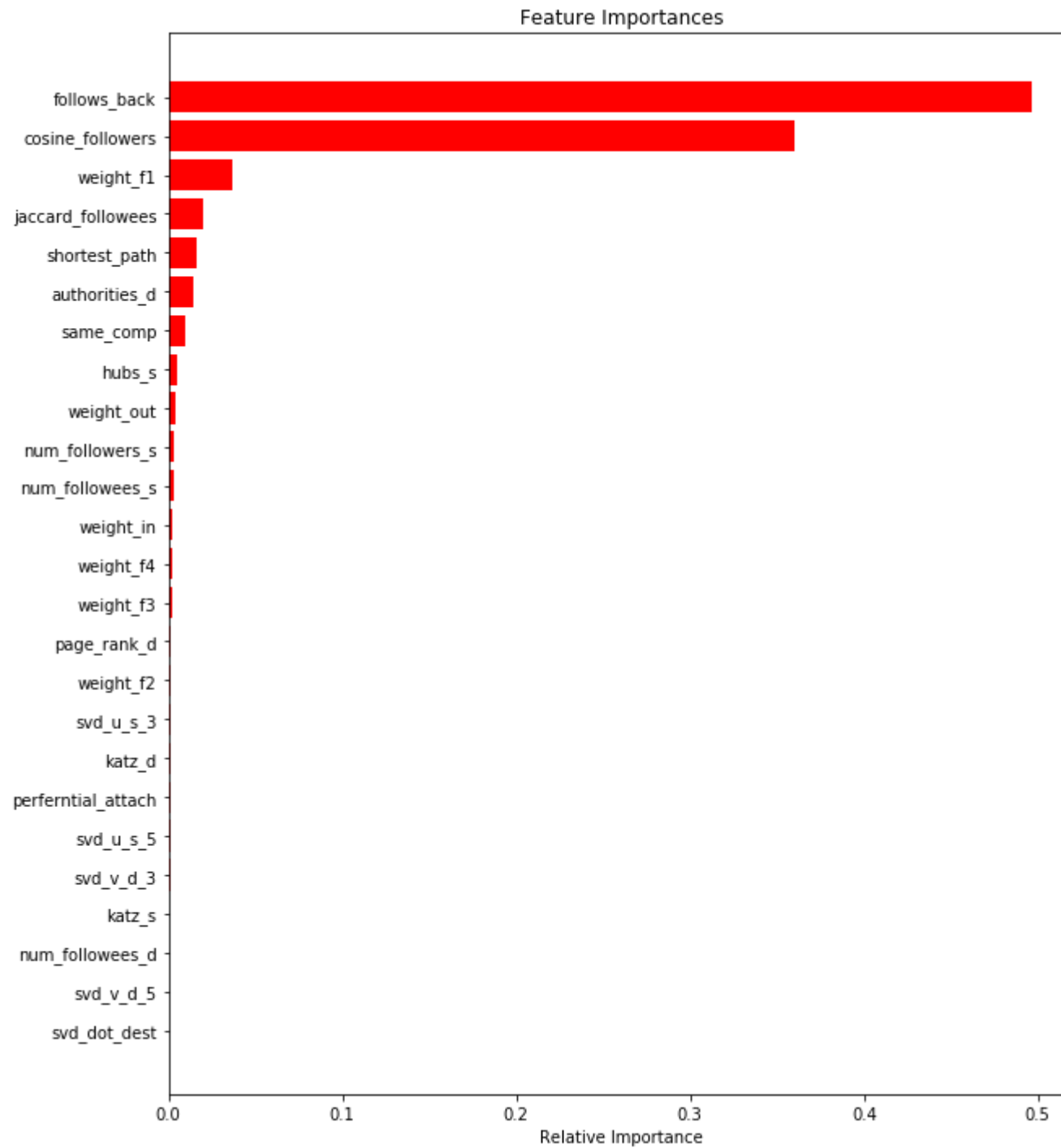


```
Test confusion_matrix
```

| | Confusion matrix | |
|---|---|---|
| 24724.000 | 322.000 |
| 3147.000 | 21809.000 |

| | Precision matrix | |
|---|---|---|
| 0.887 | 0.015 |
| 0.113 | 0.985 |

| | Recall matrix | |
|---|---|---|
| 0.987 | 0.013 |
| 0.126 | 0.874 |

In [52]:

```python
from sklearn.metrics import roc_curve, auc
fpr,tpr,ths = roc_curve(y_test,y_pred_test)
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()
```



Receiver operating characteristic with test data

ROC curve (area = 0.93)

```
In [53]:  1  features = df_final_train.columns
          2  importances = xgb.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()
```

## Feature Importances

```
In [54]:  1  from prettytable import PrettyTable
          2  x = PrettyTable()
          3  x.field_names = ["Model", "n_estimators", "max_depth", "Train_f1_Score","Test_f1_Score"]
          4  x.add_row(['Random Forest',n_estimator_rf,max_depth_rf,Train_f1_Score_rf,Test_f1_Score_rf])
          5  x.add_row(['XGBOOST',n_estimator_xgb,max_depth_xgb,Train_f1_Score_xgb,Test_f1_Score_xgb])
          6  print(x)
```

```
+---------------+--------------+-----------+--------------------+--------------------+
|     Model     | n_estimators | max_depth |   Train_f1_Score   |   Test_f1_Score    |
+---------------+--------------+-----------+--------------------+--------------------+
| Random Forest |     114      |     14    | 0.9720793459917744 | 0.9188237291012475 |
|    XGBOOST    |     118      |     12    | 0.9981790895447723 | 0.9263278611931106 |
+---------------+--------------+-----------+--------------------+--------------------+
```

# Conclusion :

- Facebook Friend Recommendation Case study is one of its kind as most of the Internet Companies like Facebook, Instagram, redit, Github etc has graph based features. Understanding various Graph feature is important and all through a fun ride.
- In EDA we made use of np.percentile and saw how we can even look at 99.1 percentile
- Explored different Graph Based Feature Engineering in which **Follow back & cosine followers** were important one
- Added new features as part of Assignment
- XGboost performed better than Random Forest

```
In [ ]:  1
```