

# 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_stage5.h5', 'train_df', mode='r')
df_final_test = read_hdf('data/fea_sample/storage_sample_stage5.h5', 'test_df', 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_followers_d',
      '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',
      'pref_attach_followers', 'pref_attach_followees', 'svd_u_dot',
      'svd_v_dot'],
      dtype='object', length=54)
```

```
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)
```

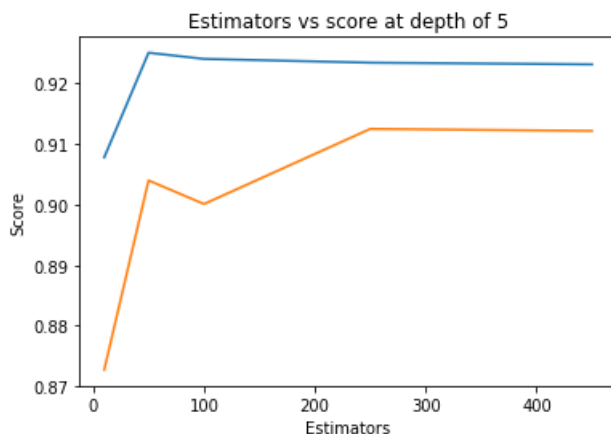
In [6]:

```
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 0.907812598524497 test Score 0.872724146084564
Estimators = 50 Train Score 0.9250816384388422 test Score 0.9039959147197754
Estimators = 100 Train Score 0.9240620226411935 test Score 0.900104036179112
Estimators = 250 Train Score 0.9234378600305829 test Score 0.9125084005376344
Estimators = 450 Train Score 0.9231606260986022 test Score 0.9121776173588232
```

Out[6]:

Text(0.5, 1.0, 'Estimators vs score at depth of 5')



In [7]:

```
depths = [3, 9, 11, 15, 20, 35, 50, 70, 130]
train_scores = []
test_scores = []
for i in depths:
```

```

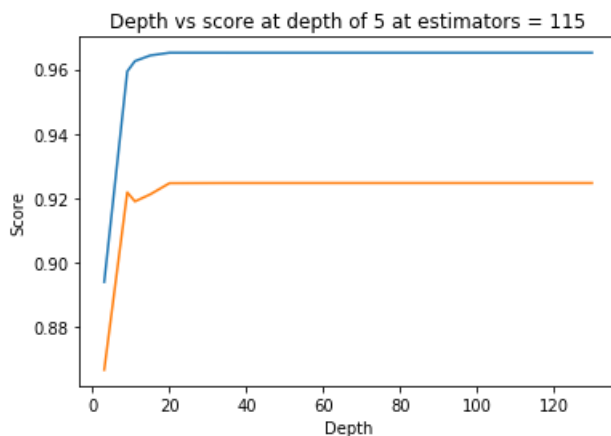
clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                             max_depth=i, max_features='auto', max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=52, min_samples_split=120,
                             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 = 115')
plt.show()

```

```

depth = 3 Train Score 0.8939495798319328 test Score 0.8666367885267943
depth = 9 Train Score 0.9593580594775045 test Score 0.9217771426169559
depth = 11 Train Score 0.9625519263663761 test Score 0.9189530914148438
depth = 15 Train Score 0.9643078991048971 test Score 0.9211613011966964
depth = 20 Train Score 0.9651459798304619 test Score 0.9246309246309246
depth = 35 Train Score 0.9651405262997284 test Score 0.9246756518453206
depth = 50 Train Score 0.9651405262997284 test Score 0.9246756518453206
depth = 70 Train Score 0.9651405262997284 test Score 0.9246756518453206
depth = 130 Train Score 0.9651405262997284 test Score 0.9246756518453206

```



In [8]:

```

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(100, 300),
              "max_depth": sp_randint(1, 20),
              "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, return_train_score=True)

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 [0.92229714 0.96370092 0.92583173 0.91811429 0.85030893]
mean train scores [0.92246108 0.96440014 0.92606435 0.91838134 0.85074738]

```

In [9]:

```
print(rf_random.best_estimator_)
```

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                        max_depth=13, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=33, min_samples_split=138,
                        min_weight_fraction_leaf=0.0, n_estimators=232,
                        n_jobs=-1, oob_score=False, random_state=25, verbose=0,
                        warm_start=False)
```

In [10]:

```
clf = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                             max_depth=13, max_features='auto', max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=33, min_samples_split=138,
                             min_weight_fraction_leaf=0.0, n_estimators=232,
                             n_jobs=-1, oob_score=False, random_state=25, verbose=0,
                             warm_start=False)
```

In [11]:

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

In [12]:

```
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 0.9642649829421049  
Test f1 score 0.9212394413666717

In [13]:

```
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')
```

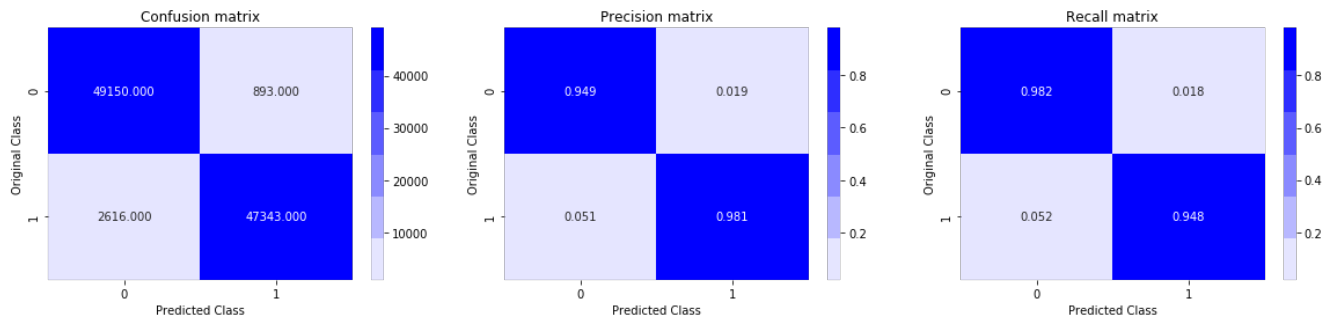
```
plt.ylabel('Original Class')
plt.title("Recall matrix")

plt.show()
```

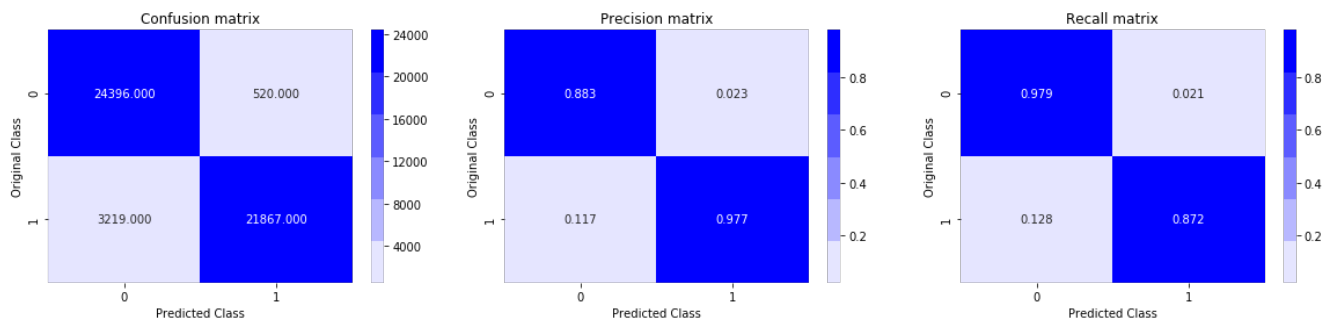
In [14]:

```
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

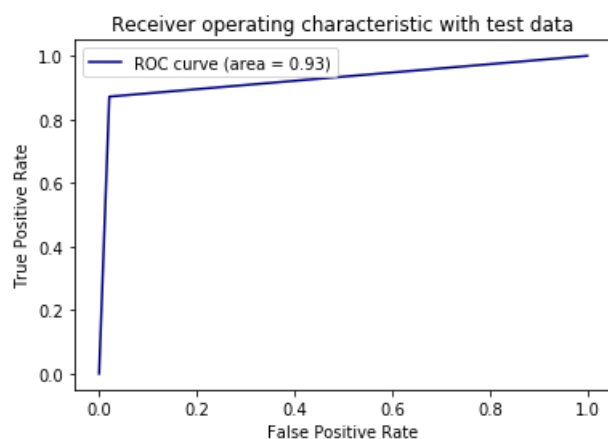


Test confusion\_matrix



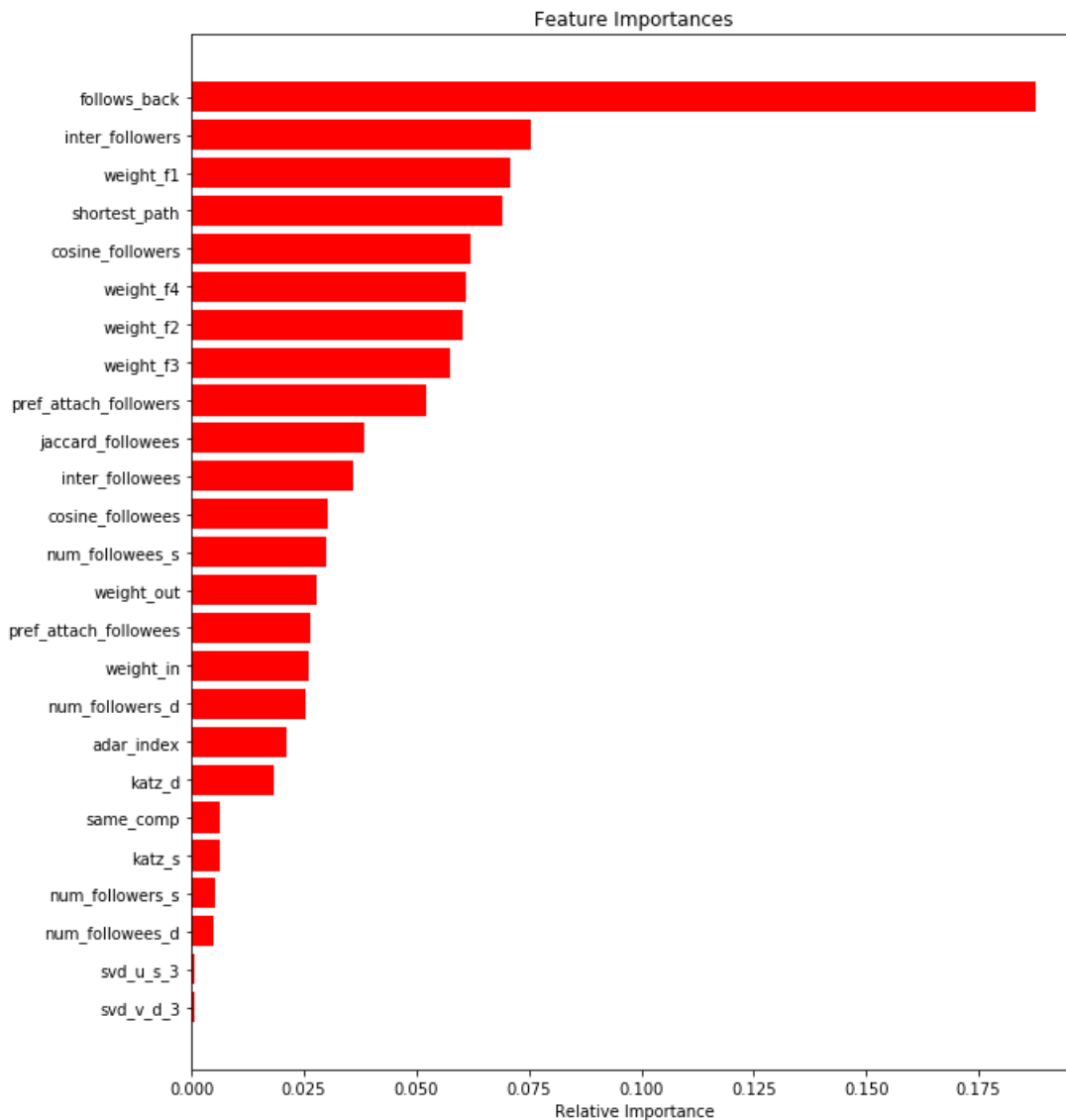
In [15]:

```
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 [16]:

```
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 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](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.

## XG Boost

In [17]:

```

estimators = [10,50,100,250,450]
train_scores = []
test_scores = []
for i in estimators:
    clf = xgb.XGBClassifier(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')

```

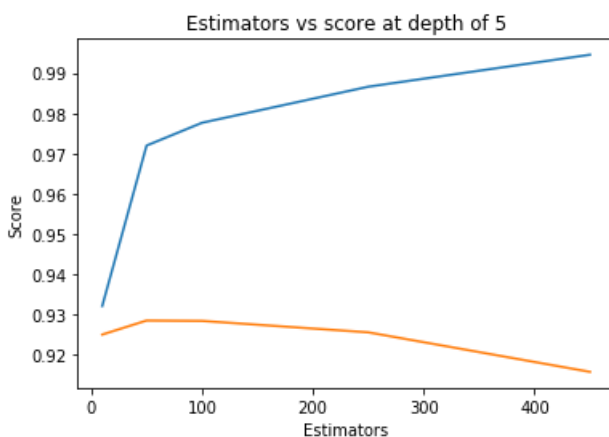
```

Estimators = 10 Train Score 0.9321499720673335 test Score 0.9250638360731845
Estimators = 50 Train Score 0.971956322073681 test Score 0.9285487590715782
Estimators = 100 Train Score 0.977614210110104 test Score 0.9284537737044389
Estimators = 250 Train Score 0.9865631772592107 test Score 0.925618785699175
Estimators = 450 Train Score 0.9944987652338025 test Score 0.9158132787165045

```

Out[17]:

Text(0.5, 1.0, 'Estimators vs score at depth of 5')



In [18]:

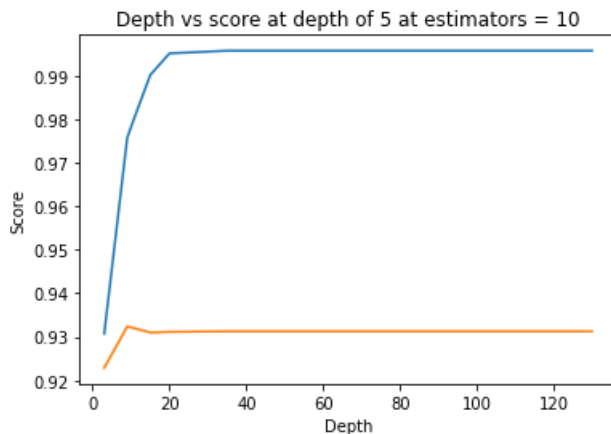
```

depths = [3,9,11,15,20,35,50,70,130]
train_scores = []
test_scores = []
for i in depths:
    clf = xgb.XGBClassifier(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=10, 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 = 10')

```

```
fig,ax=plt.subplots(1,2,figsize=(10,5))
plt.show()
```

```
depth = 3 Train Score 0.9307141352433116 test Score 0.9228822673499083
depth = 9 Train Score 0.9757445813984209 test Score 0.932388382999306
depth = 11 Train Score 0.9806776018328908 test Score 0.9319433351549035
depth = 15 Train Score 0.9901427189163038 test Score 0.9309910590605717
depth = 20 Train Score 0.9951005983695435 test Score 0.9311508393938759
depth = 35 Train Score 0.9957051397836514 test Score 0.9312612574875382
depth = 50 Train Score 0.9957051397836514 test Score 0.9312612574875382
depth = 70 Train Score 0.9957051397836514 test Score 0.9312612574875382
depth = 130 Train Score 0.9957051397836514 test Score 0.9312612574875382
```



In [19]:

```
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(1,100),
              "max_depth": sp_randint(1,20),
              "min_samples_split": sp_randint(110,190),
              "min_samples_leaf": sp_randint(25,65)}

clf = xgb.XGBClassifier(random_state=25,n_jobs=-1)

rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
                               n_iter=5,cv=10,scoring='f1',random_state=25, return_train_score=True
)

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 [0.9741056 0.97506212 0.92977995 0.92498501 0.93804115]
mean train scores [0.9746468 0.98471163 0.93038571 0.92529896 0.93832395]
```

In [20]:

```
print(rf_random.best_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=13,
              min_child_weight=1, min_samples_leaf=33, min_samples_split=138,
              missing=None, n_estimators=5, n_jobs=-1, nthread=None,
              objective='binary:logistic', random_state=25, reg_alpha=0,
              reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
              subsample=1, verbosity=1)
```



In [21]:

```
clf = xgb.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=13,
                        min_child_weight=1, min_samples_leaf=33, min_samples_split=138,
                        missing=None, n_estimators=5, n_jobs=-1, nthread=None,
                        objective='binary:logistic', random_state=25, reg_alpha=0,
                        reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
                        subsample=1, verbosity=1)
```

In [22]:

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

In [23]:

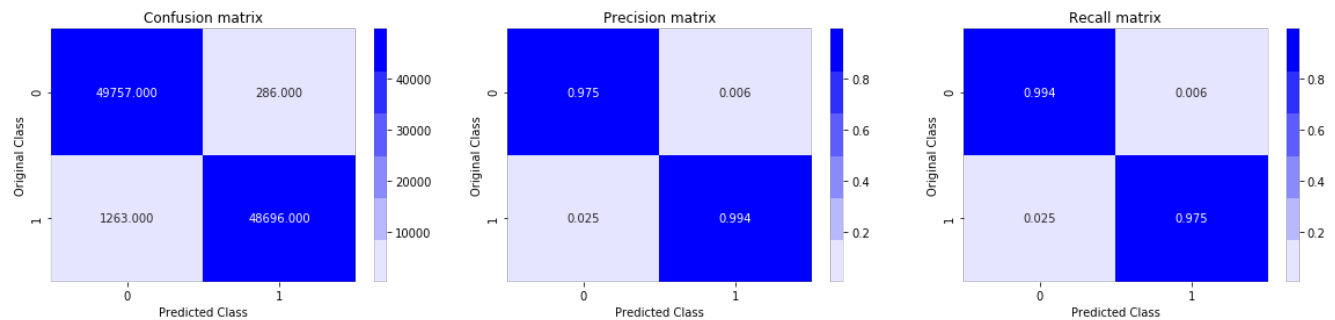
```
from sklearn.metrics import f1_score
print('Train f1 score',f1_score(y_train,y_train_pred))
print('Test f1 score',f1_score(y_test,y_test_pred))
```

Train f1 score 0.9843442051323517  
Test f1 score 0.9311907359234706

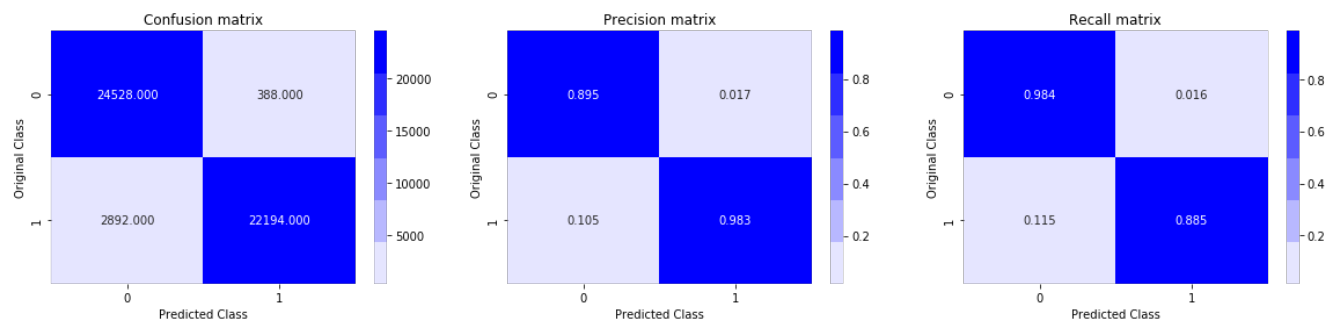
In [24]:

```
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



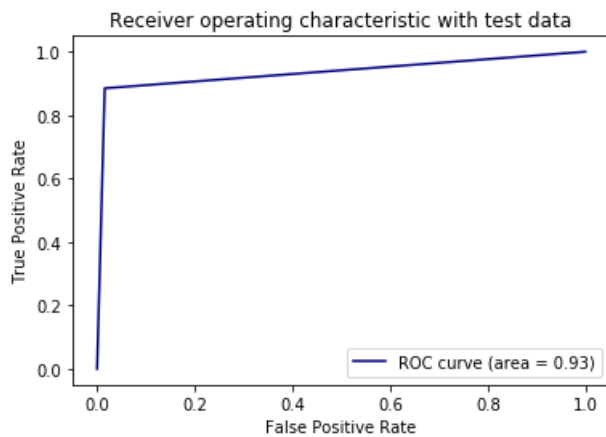
In [25]:

```
from sklearn.metrics import roc_curve, auc
for, tpr, ths = roc_curve(y_test,y_test_pred)
```

```

fpr, tpr, auc = roc_curve(y_test, y_hat)
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()

```

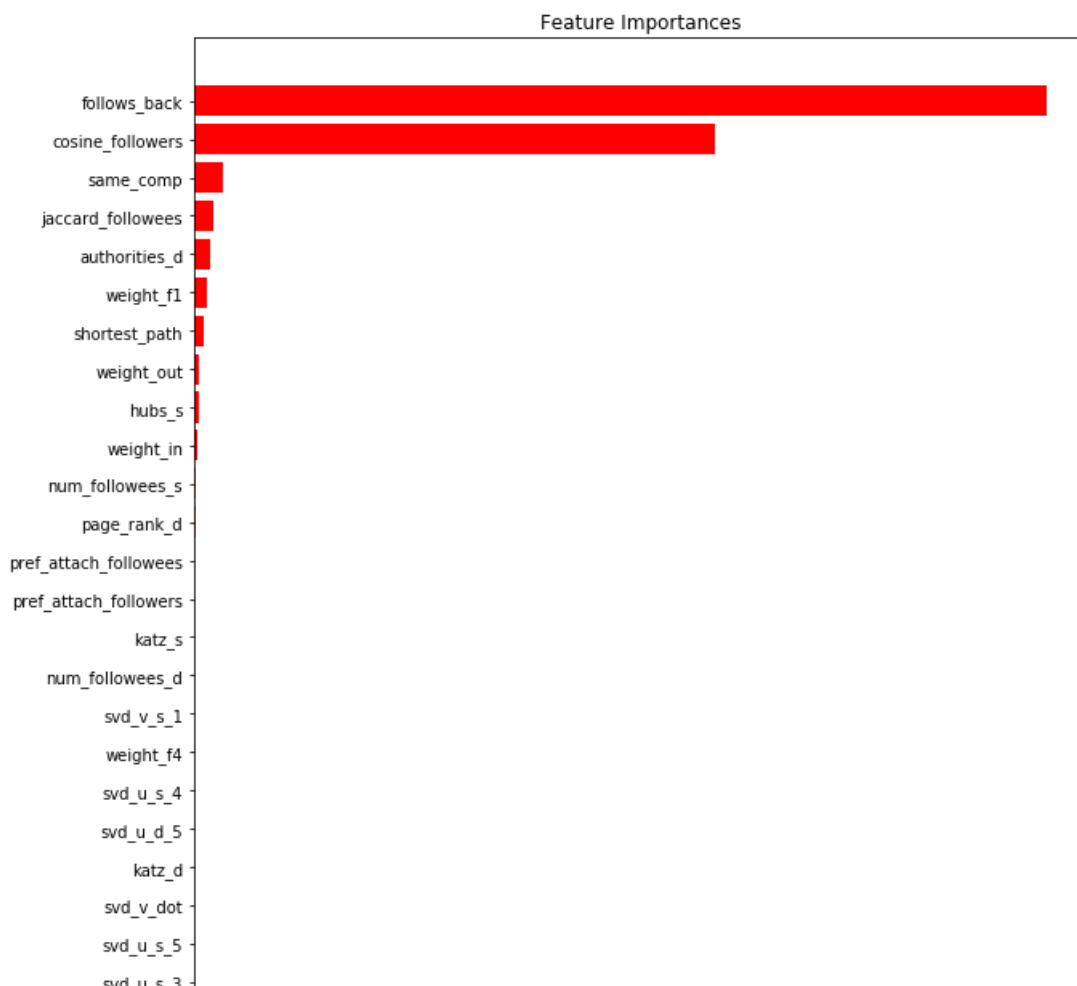


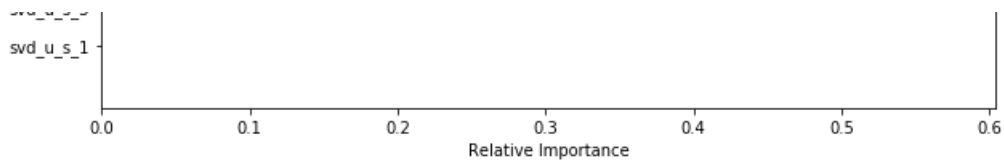
In [26]:

```

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 Importance')
plt.show()

```





## Conclusion

- Step 1. Perform EDA
- Step 2. Split the data into train and test data
- Step 3. Perform various Featurization method
- Step 4. Apply Two ML algorithm. (RandomForest and XGBoost)

In [27]:

```
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ['Model', 'n_estimators', 'max_depth', 'train f1 score', 'test f1 score']

x.add_row(['RandomForest', 232, 13, 0.96, 0.92])
x.add_row(['XGBoost', 5, 13, 0.98, 0.93])

print(x)
```

Model	n_estimators	max_depth	train f1 score	test f1 score
RandomForest	232	13	0.96	0.92
XGBoost	5	13	0.98	0.93

From the observation, recall matrix on test data using XGBoost has been improved than using RandomForest

In [ ]: