

```

In [7]: df_final_train['PA_followee'] = df_final_train.apply(lambda row:
PA_follower(row['source_node'],row['destination_node']),axis=1)
df_final_test['PA_followee'] = df_final_test.apply(lambda row:
PA_follower(row['source_node'],row['destination_node']),axis=1)

Creating another New feature SVD_DOT

In [8]: def svd_dot_fun(a,b):
temp=0
for i,j in zip(a,b):
temp=temp+1*j
return temp

In [9]: #Creating another Feature svd_dot
U=['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4',
svd_u_s_5', 'svd_u_s_6']
V=['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6']

In [10]: U1=['svd_v_s_1', 'svd_v_s_2',
'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6']
V1=['svd_v_d_1',
'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6']

In [11]: df_final_test['svd_dot_u']=df_final_test.apply(lambda row: svd_dot_fun(row[U],row[V]),axis=1)
df_final_test['svd_dot_v']=df_final_test.apply(lambda row: svd_dot_fun(row[U1],row[V1]),axis=1)

In [12]: df_final_train['svd_dot_u']=df_final_train.apply(lambda row: svd_dot_fun(row[U],row[V]),axis=1)
df_final_train['svd_dot_v']=df_final_train.apply(lambda row: svd_dot_fun(row[U1],row[V1]),axis=1)

In [13]: #writing the df_final_train, df_final_test into csv files with new features svd_dot and Preferential Attachment
df_final_train.to_csv('df_final_train.csv')
df_final_test.to_csv('df_final_test.csv')

In [14]: #by using these two files, we can avoid running above code
df_final_train=pd.read_csv('df_final_train.csv',index_col=None)
df_final_test=pd.read_csv('df_final_test.csv')

Creating a train,test,CV split

In [15]: X_train,X_cv=train_test_split(df_final_train,test_size=0.20)
y_train=X_train.indicator_link
y_cv=X_cv.indicator_link
y_test=df_final_test.indicator_link
X_test=df_final_test

In [16]: X_train.drop(['source_node', 'destination_node', 'indicator_link', 'Unnamed: 0'],inplace=True,axis=1)
X_cv.drop(['source_node', 'destination_node', 'indicator_link', 'Unnamed: 0'],inplace=True,axis=1)
X_test.drop(['source_node', 'destination_node', 'indicator_link', 'Unnamed: 0'],inplace=True,axis=1)

```

```
In [9]: #Creating another Feature Svd dot dot
Us['svd_u_s_1', 'svd_u_s_2', 'svd_u_s_3', 'svd_u_s_4',
     'svd_u_s_5', 'svd_u_s_6']
Vv=['svd_u_d_1', 'svd_u_d_2', 'svd_u_d_3',
     'svd_u_d_4', 'svd_u_d_5', 'svd_u_d_6']

In [10]:
U1=['svd_v_s_1', 'svd_v_s_2',
     'svd_v_s_3', 'svd_v_s_4', 'svd_v_s_5', 'svd_v_s_6']
V1=['svd_v_d_1',
     'svd_v_d_2', 'svd_v_d_3', 'svd_v_d_4', 'svd_v_d_5', 'svd_v_d_6']

In [11]:
df_final_test['svd_dot_u']=df_final_test.apply(lambda row: svd_dot_fun(row[U],row[V]),axis=1)
df_final_test['svd_dot_v']=df_final_test.apply(lambda row: svd_dot_fun(row[U1],row[V1]),axis=1)

In [12]:
df_final_train['svd_dot_u']=df_final_train.apply(lambda row: svd_dot_fun(row[U],row[V]),axis=1)
df_final_train['svd_dot_v']=df_final_train.apply(lambda row: svd_dot_fun(row[U1],row[V1]),axis=1)

In [13]:
#writing the df_final_train, df_final_test into csv files with new features svd_dot and Preferential Attachment
df_final_train.to_csv('df_final_train.csv')
df_final_test.to_csv('df_final_test.csv')

In [14]:
#by suing these two files, we can avoid running above code
df_final_train=pd.read_csv('df_final_train.csv',index_col=None)
df_final_test=pd.read_csv('df_final_train.csv')
```

## Creating a train,test,CV split

```
In [15]:
X_train,X_cv=train_test_split(df_final_train,test_size=0.20)
y_train=X_train.indicator_link
y_cv=X_cv.indicator_link
y_test=df_final_test.indicator_link
X_test=df_final_test

In [16]:
X_train.drop(['source_node','destination_node','indicator_link','Unnamed: 0'],inplace=True,axis=1)
X_cv.drop(['source_node','destination_node','indicator_link','Unnamed: 0'],inplace=True,axis=1)
X_test.drop(['source_node','destination_node','indicator_link','Unnamed: 0'],inplace=True,axis=1)

In [17]:
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.title('Original Class')
plt.xlabel('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()

```

## Building a model using with XGBoost

```

In [18]: alpha=[10,50,100,500,1000,2000]
cv_log_error_array=[]
for i in alpha:
    x_clf=XGBClassifier(n_estimators=1, nthread=1)
    cv_clf=fit(X_train,y_train,
               sig_clf = CalibratedClassifierCV(x_clf, method='sigmoid'))
    cv_log_clf=fit(x_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_clf.classes_, eps=1e-15))

for i in range(len(cv_log_error_array)):
    print ('log_loss for c = ',alpha[i], 'is',cv_log_error_array[i])

best_alpha = np.argmin(cv_log_error_array)

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title('Cross Validation Error for each alpha')
plt.xlabel('Alpha i's')
plt.ylabel('Error measure')
plt.show()

```

```

[17:32:27] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:29] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:30] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:31] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:32] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:33] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:34] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:35] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:42] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:47] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:52] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:32:58] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:33:09] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:33:22] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective
ve 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[17:33:54] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective

```

```
[*]
[7:37:52] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:37:53] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:37:53] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:41:31] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:43:03] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:43:03] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:46:07] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:47:39] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:51:41] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:53:03] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:55:04] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:57:37] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[7:59:08] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[8:02:43] WARNING: C:/Users/Administrator/workspace/ghobst-win64_release.1.5.1/src/learner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
log_loss for c = 10 is 0.6923834406586868
log_loss for c = 15 is 0.6842727740000000
log_loss for c = 100 is 0.6563434993146653973
log_loss for c = 500 is 0.60689933535649996
log_loss for c = 1000 is 0.6020478457866406294
log_loss for c = 2000 is 0.60257754635201091

Cross Validation Error for each alpha
(10.0, 0.93)
0.095
0.085
0.080
0.075
0.070
0.065
0.060
0.055
0.050
0.045
0.040
0.035
0.030
0.025
0.020
0.015
0.010
0.005
0.000
-0.005
-0.010
-0.015
-0.020
-0.025
-0.030
-0.035
-0.040
-0.045
-0.050
-0.055
-0.060
-0.065
-0.070
-0.075
-0.080
-0.085
-0.090
-0.095
-0.100
-0.105
-0.110
-0.115
-0.120
-0.125
-0.130
-0.135
-0.140
-0.145
-0.150
-0.155
-0.160
-0.165
-0.170
-0.175
-0.180
-0.185
-0.190
-0.195
-0.200
-0.205
-0.210
-0.215
-0.220
-0.225
-0.230
-0.235
-0.240
-0.245
-0.250
-0.255
-0.260
-0.265
-0.270
-0.275
-0.280
-0.285
-0.290
-0.295
-0.300
-0.305
-0.310
-0.315
-0.320
-0.325
-0.330
-0.335
-0.340
-0.345
-0.350
-0.355
-0.360
-0.365
-0.370
-0.375
-0.380
-0.385
-0.390
-0.395
-0.400
-0.405
-0.410
-0.415
-0.420
-0.425
-0.430
-0.435
-0.440
-0.445
-0.450
-0.455
-0.460
-0.465
-0.470
-0.475
-0.480
-0.485
-0.490
-0.495
-0.500
-0.505
-0.510
-0.515
-0.520
-0.525
-0.530
-0.535
-0.540
-0.545
-0.550
-0.555
-0.560
-0.565
-0.570
-0.575
-0.580
-0.585
-0.590
-0.595
-0.600
-0.605
-0.610
-0.615
-0.620
-0.625
-0.630
-0.635
-0.640
-0.645
-0.650
-0.655
-0.660
-0.665
-0.670
-0.675
-0.680
-0.685
-0.690
-0.695
-0.700
-0.705
-0.710
-0.715
-0.720
-0.725
-0.730
-0.735
-0.740
-0.745
-0.750
-0.755
-0.760
-0.765
-0.770
-0.775
-0.780
-0.785
-0.790
-0.795
-0.800
-0.805
-0.810
-0.815
-0.820
-0.825
-0.830
-0.835
-0.840
-0.845
-0.850
-0.855
-0.860
-0.865
-0.870
-0.875
-0.880
-0.885
-0.890
-0.895
-0.900
-0.905
-0.910
-0.915
-0.920
-0.925
-0.930
-0.935
-0.940
-0.945
-0.950
-0.955
-0.960
-0.965
-0.970
-0.975
-0.980
-0.985
-0.990
-0.995
1.000
1.005
1.010
1.015
1.020
1.025
1.030
1.035
1.040
1.045
1.050
1.055
1.060
1.065
1.070
1.075
1.080
1.085
1.090
1.095
1.100
1.105
1.110
1.115
1.120
1.125
1.130
1.135
1.140
1.145
1.150
1.155
1.160
1.165
1.170
1.175
1.180
1.185
1.190
1.195
1.200
1.205
1.210
1.215
1.220
1.225
1.230
1.235
1.240
1.245
1.250
1.255
1.260
1.265
1.270
1.275
1.280
1.285
1.290
1.295
1.300
1.305
1.310
1.315
1.320
1.325
1.330
1.335
1.340
1.345
1.350
1.355
1.360
1.365
1.370
1.375
1.380
1.385
1.390
1.395
1.400
1.405
1.410
1.415
1.420
1.425
1.430
1.435
1.440
1.445
1.450
1.455
1.460
1.465
1.470
1.475
1.480
1.485
1.490
1.495
1.500
1.505
1.510
1.515
1.520
1.525
1.530
1.535
1.540
1.545
1.550
1.555
1.560
1.565
1.570
1.575
1.580
1.585
1.590
1.595
1.600
1.605
1.610
1.615
1.620
1.625
1.630
1.635
1.640
1.645
1.650
1.655
1.660
1.665
1.670
1.675
1.680
1.685
1.690
1.695
1.700
1.705
1.710
1.715
1.720
1.725
1.730
1.735
1.740
1.745
1.750
1.755
1.760
1.765
1.770
1.775
1.780
1.785
1.790
1.795
1.800
1.805
1.810
1.815
1.820
1.825
1.830
1.835
1.840
1.845
1.850
1.855
1.860
1.865
1.870
1.875
1.88
```

```

ve "binary:logistic" was changed from "error" to "logloss". Explicitly set eval_metric if you'd like to restore the old behavior.
[18:05:44] WARNING: C:/Users/Administrator/workspace/sgxboost-win64-release-1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from "error" to "logloss". Explicitly set eval_metric if you'd like to restore the old behavior.
[18:05:55] WARNING: C:/Users/Administrator/workspace/sgxboost-win64-release-1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from "error" to "logloss". Explicitly set eval_metric if you'd like to restore the old behavior.
[18:06:06] WARNING: C:/Users/Administrator/workspace/sgxboost-win64-release-1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from "error" to "logloss". Explicitly set eval_metric if you'd like to restore the old behavior.
[18:06:16] WARNING: C:/Users/Administrator/workspace/sgxboost-win64-release-1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from "error" to "logloss". Explicitly set eval_metric if you'd like to restore the old behavior.
For values of best alpha = 100 The train log loss is: 0.02794699022779936
For values of best alpha = 100 The cross validation log loss is: 0.05634093148653973
For values of best alpha = 100 The test log loss is: 0.0256257211729765

Confusion matrix
Precision matrix
Recall matrix

```

```

The Best f1 score for train is 0.9969531991409821
The Best f1 score for test is 0.9936145684772995
The Best f1 score for Cv is 0.9881311144738267

In [20]:
clf=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
clf.fit(X_train,y_train)

[18:06:30] WARNING: C:/Users/Administrator/workspace/sgxboost-win64-release-1.5.1/src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from "error" to "logloss". Explicitly set eval_metric if you'd like to restore the old behavior.

Out[20]:
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bytree=1,
              colsample_bynode=1, colsample_bytree=1, enable_categorical=False,
              gamma=0, gpu_id=-1, importance_type=None,
              interaction_constraints='', learning_rate=0.300000012,
              max_delta_step=0, max_depth=6, min_child_weight=1, missing=nan,
              monotone_constraints={}, n_estimators=100, n_jobs=12,
              nthread=-1, num_parallel_trees=1, predictor='auto', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
              tree_method='exact', validate_parameters=1, verbosity=None)

Feature importance

In [28]:
features = X_train.columns
importances = x_rf.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')

```

Feature Importances

Feature	Importance (approx.)
cosine_followers	0.95
follows_back	0.85
weight_f1	0.35
jaccard_followees	0.15
rum_followers_3	0.10
shortest_path	0.08
weight_in	0.07
same_comp	0.06
rum_followees_3	0.05
weight_out	0.04