

In [1]:

```
import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import sqlite3
from sqlalchemy import create_engine # database connection
import csv
import os

import datetime as dt
import numpy as np

from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
#from sklearn.cross_validation import StratifiedKFold
from sklearn.model_selection import StratifiedKFold
from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier

from sklearn.model_selection import cross_val_score
#from sklearn.model_selection import c
from sklearn.linear_model import SGDClassifier
#from mlxtend.classifier import StackingClassifier

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import precision_recall_curve, auc, roc_curve
```

## 4. Machine Learning Models

### 4.1 Reading data from file and storing into sql table

In [2]:

```
#Creating db file from csv
if not os.path.isfile('train.db'):
    disk_engine = create_engine('sqlite:///train.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in pd.read_csv('final_features_train.csv', names=['Unnamed: 0', 'id', 'is_duplicate', 'cwc_min', 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq', 'first_word_eq', 'abs_len_dif', 'mean_len', 'token_set_ratio', 'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio', 'longest_substr_ratio', 'freq_qid1', 'freq_qid2', 'q1len', 'q2len', 'q1_n_words', 'q2_n_words', 'word_Common', 'word_Total', 'word_share', 'freq_a1+a2', 'freq_a1-
```

```

q2','0_x','1_x','2_x','3_x','4_x','5_x','6_x','7_x','8_x','9_x','10_x','11_x','12_x','13_x','14_x','
'15_x','16_x','17_x','18_x','19_x','20_x','21_x','22_x','23_x','24_x','25_x','26_x','27_x','28_x','
29_x','30_x','31_x','32_x','33_x','34_x','35_x','36_x','37_x','38_x','39_x','40_x','41_x','42_x','4
3_x','44_x','45_x','46_x','47_x','48_x','49_x','50_x','51_x','52_x','53_x','54_x','55_x','56_x','57
_x','58_x','59_x','60_x','61_x','62_x','63_x','64_x','65_x','66_x','67_x','68_x','69_x','70_x','71_
_x','72_x','73_x','74_x','75_x','76_x','77_x','78_x','79_x','80_x','81_x','82_x','83_x','84_x','85_x
','86_x','87_x','88_x','89_x','90_x','91_x','92_x','93_x','94_x','95_x','96_x','97_x','98_x','99_x'
,'100_x','101_x','102_x','103_x','104_x','105_x','106_x','107_x','108_x','109_x','110_x','111_x','
112_x','113_x','114_x','115_x','116_x','117_x','118_x','119_x','120_x','121_x','122_x','123_x','12
4_x','125_x','126_x','127_x','128_x','129_x','130_x','131_x','132_x','133_x','134_x','135_x','136_
_x','137_x','138_x','139_x','140_x','141_x','142_x','143_x','144_x','145_x','146_x','147_x','148_x'
,'149_x','150_x','151_x','152_x','153_x','154_x','155_x','156_x','157_x','158_x','159_x','160_x','
161_x','162_x','163_x','164_x','165_x','166_x','167_x','168_x','169_x','170_x','171_x','172_x','17
3_x','174_x','175_x','176_x','177_x','178_x','179_x','180_x','181_x','182_x','183_x','184_x','185_
_x','186_x','187_x','188_x','189_x','190_x','191_x','192_x','193_x','194_x','195_x','196_x','197_x'
,'198_x','199_x','200_x','201_x','202_x','203_x','204_x','205_x','206_x','207_x','208_x','209_x','
210_x','211_x','212_x','213_x','214_x','215_x','216_x','217_x','218_x','219_x','220_x','221_x','22
2_x','223_x','224_x','225_x','226_x','227_x','228_x','229_x','230_x','231_x','232_x','233_x','234_
_x','235_x','236_x','237_x','238_x','239_x','240_x','241_x','242_x','243_x','244_x','245_x','246_x'
,'247_x','248_x','249_x','250_x','251_x','252_x','253_x','254_x','255_x','256_x','257_x','258_x','
259_x','260_x','261_x','262_x','263_x','264_x','265_x','266_x','267_x','268_x','269_x','270_x','27
1_x','272_x','273_x','274_x','275_x','276_x','277_x','278_x','279_x','280_x','281_x','282_x','283_
_x','284_x','285_x','286_x','287_x','288_x','289_x','290_x','291_x','292_x','293_x','294_x','295_x'
,'296_x','297_x','298_x','299_x','300_x','301_x','302_x','303_x','304_x','305_x','306_x','307_x','
308_x','309_x','310_x','311_x','312_x','313_x','314_x','315_x','316_x','317_x','318_x','319_x','32
0_x','321_x','322_x','323_x','324_x','325_x','326_x','327_x','328_x','329_x','330_x','331_x','332_
_x','333_x','334_x','335_x','336_x','337_x','338_x','339_x','340_x','341_x','342_x','343_x','344_x'
,'345_x','346_x','347_x','348_x','349_x','350_x','351_x','352_x','353_x','354_x','355_x','356_x','
357_x','358_x','359_x','360_x','361_x','362_x','363_x','364_x','365_x','366_x','367_x','368_x','36
9_x','370_x','371_x','372_x','373_x','374_x','375_x','376_x','377_x','378_x','379_x','380_x','381_
_x','382_x','383_x','0_y','1_y','2_y','3_y','4_y','5_y','6_y','7_y','8_y','9_y','10_y','11_y','12_y'
,'13_y','14_y','15_y','16_y','17_y','18_y','19_y','20_y','21_y','22_y','23_y','24_y','25_y','26_y','
27_y','28_y','29_y','30_y','31_y','32_y','33_y','34_y','35_y','36_y','37_y','38_y','39_y','40_y','
41_y','42_y','43_y','44_y','45_y','46_y','47_y','48_y','49_y','50_y','51_y','52_y','53_y','54_y','5
5_y','56_y','57_y','58_y','59_y','60_y','61_y','62_y','63_y','64_y','65_y','66_y','67_y','68_y','69
_y','70_y','71_y','72_y','73_y','74_y','75_y','76_y','77_y','78_y','79_y','80_y','81_y','82_y','83_
_y','84_y','85_y','86_y','87_y','88_y','89_y','90_y','91_y','92_y','93_y','94_y','95_y','96_y','97_y
','98_y','99_y','100_y','101_y','102_y','103_y','104_y','105_y','106_y','107_y','108_y','109_y','11
0_y','111_y','112_y','113_y','114_y','115_y','116_y','117_y','118_y','119_y','120_y','121_y','122_
_y','123_y','124_y','125_y','126_y','127_y','128_y','129_y','130_y','131_y','132_y','133_y','134_y'
,'135_y','136_y','137_y','138_y','139_y','140_y','141_y','142_y','143_y','144_y','145_y','146_y','
147_y','148_y','149_y','150_y','151_y','152_y','153_y','154_y','155_y','156_y','157_y','158_y','15
9_y','160_y','161_y','162_y','163_y','164_y','165_y','166_y','167_y','168_y','169_y','170_y','171_
_y','172_y','173_y','174_y','175_y','176_y','177_y','178_y','179_y','180_y','181_y','182_y','183_y'
,'184_y','185_y','186_y','187_y','188_y','189_y','190_y','191_y','192_y','193_y','194_y','195_y','
196_y','197_y','198_y','199_y','200_y','201_y','202_y','203_y','204_y','205_y','206_y','207_y','20
8_y','209_y','210_y','211_y','212_y','213_y','214_y','215_y','216_y','217_y','218_y','219_y','220_
_y','221_y','222_y','223_y','224_y','225_y','226_y','227_y','228_y','229_y','230_y','231_y','232_y'
,'233_y','234_y','235_y','236_y','237_y','238_y','239_y','240_y','241_y','242_y','243_y','244_y','
245_y','246_y','247_y','248_y','249_y','250_y','251_y','252_y','253_y','254_y','255_y','256_y','25
7_y','258_y','259_y','260_y','261_y','262_y','263_y','264_y','265_y','266_y','267_y','268_y','269_
_y','270_y','271_y','272_y','273_y','274_y','275_y','276_y','277_y','278_y','279_y','280_y','281_y'
,'282_y','283_y','284_y','285_y','286_y','287_y','288_y','289_y','290_y','291_y','292_y','293_y','
294_y','295_y','296_y','297_y','298_y','299_y','300_y','301_y','302_y','303_y','304_y','305_y','30
6_y','307_y','308_y','309_y','310_y','311_y','312_y','313_y','314_y','315_y','316_y','317_y','318_y'
,'319_y','320_y','321_y','322_y','323_y','324_y','325_y','326_y','327_y','328_y','329_y','330_y'
,'331_y','332_y','333_y','334_y','335_y','336_y','337_y','338_y','339_y','340_y','341_y','342_y','
343_y','344_y','345_y','346_y','347_y','348_y','349_y','350_y','351_y','352_y','353_y','354_y','35
5_y','356_y','357_y','358_y','359_y','360_y','361_y','362_y','363_y','364_y','365_y','366_y','367_
_y','368_y','369_y','370_y','371_y','372_y','373_y','374_y','375_y','376_y','377_y','378_y','379_y'
,'380_y','381_y','382_y','383_y'], chunksize=chunksize, iterator=True, encoding='utf-8', ):
    df.index += index_start
    j+=1
    print('{} rows'.format(j*chunksize))
    df.to_sql('data', disk_engine, if_exists='append')
    index_start = df.index[-1] + 1

```

C:\Users\nnagari\AppData\Local\Continuum\anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3058: DtypeWarning: Columns (1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,

```
interactivity=interactivity, compiler=compiler, result=result)
```

In [3]:

```

#Creating db file from csv
if not os.path.isfile('test.db'):
    disk_engine = create_engine('sqlite:///test.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in pd.read_csv('final_features_test.csv', names=['Unnamed: 0', 'id', 'is_duplicate', 'cwc_m
in', 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq', 'first_word_eq', 'abs_len_diff',
'mean_len', 'token_set_ratio', 'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio', 'longest_substr
ratio', 'freq_qid1', 'freq_qid2', 'qlen', 'q2len', 'q1_n_words', 'q2_n_words', 'word_Common', 'word_Total',
'word_share', 'freq_q1+q2', 'freq_q1-
q2', '0_x', '1_x', '2_x', '3_x', '4_x', '5_x', '6_x', '7_x', '8_x', '9_x', '10_x', '11_x', '12_x', '13_x', '14_x',
'15_x', '16_x', '17_x', '18_x', '19_x', '20_x', '21_x', '22_x', '23_x', '24_x', '25_x', '26_x', '27_x', '28_x', '
29_x', '30_x', '31_x', '32_x', '33_x', '34_x', '35_x', '36_x', '37_x', '38_x', '39_x', '40_x', '41_x', '42_x', '4
3_x', '44_x', '45_x', '46_x', '47_x', '48_x', '49_x', '50_x', '51_x', '52_x', '53_x', '54_x', '55_x', '56_x', '57
_x', '58_x', '59_x', '60_x', '61_x', '62_x', '63_x', '64_x', '65_x', '66_x', '67_x', '68_x', '69_x', '70_x', '71
_x', '72_x', '73_x', '74_x', '75_x', '76_x', '77_x', '78_x', '79_x', '80_x', '81_x', '82_x', '83_x', '84_x', '85_x
', '86_x', '87_x', '88_x', '89_x', '90_x', '91_x', '92_x', '93_x', '94_x', '95_x', '96_x', '97_x', '98_x', '99_x',
'100_x', '101_x', '102_x', '103_x', '104_x', '105_x', '106_x', '107_x', '108_x', '109_x', '110_x', '111_x', '
112_x', '113_x', '114_x', '115_x', '116_x', '117_x', '118_x', '119_x', '120_x', '121_x', '122_x', '123_x', '12
4_x', '125_x', '126_x', '127_x', '128_x', '129_x', '130_x', '131_x', '132_x', '133_x', '134_x', '135_x', '136
_x', '137_x', '138_x', '139_x', '140_x', '141_x', '142_x', '143_x', '144_x', '145_x', '146_x', '147_x', '148_x',
'149_x', '150_x', '151_x', '152_x', '153_x', '154_x', '155_x', '156_x', '157_x', '158_x', '159_x', '160_x', '1
61_x', '162_x', '163_x', '164_x', '165_x', '166_x', '167_x', '168_x', '169_x', '170_x', '171_x', '172_x', '17
3_x', '174_x', '175_x', '176_x', '177_x', '178_x', '179_x', '180_x', '181_x', '182_x', '183_x', '184_x', '185_
x', '186_x', '187_x', '188_x', '189_x', '190_x', '191_x', '192_x', '193_x', '194_x', '195_x', '196_x', '197_x',
'198_x', '199_x', '200_x', '201_x', '202_x', '203_x', '204_x', '205_x', '206_x', '207_x', '208_x', '209_x', '
210_x', '211_x', '212_x', '213_x', '214_x', '215_x', '216_x', '217_x', '218_x', '219_x', '220_x', '221_x', '22
2_x', '223_x', '224_x', '225_x', '226_x', '227_x', '228_x', '229_x', '230_x', '231_x', '232_x', '233_x', '234
_x', '235_x', '236_x', '237_x', '238_x', '239_x', '240_x', '241_x', '242_x', '243_x', '244_x', '245_x', '246_x',
'247_x', '248_x', '249_x', '250_x', '251_x', '252_x', '253_x', '254_x', '255_x', '256_x', '257_x', '258_x', '1
259_x', '260_x', '261_x', '262_x', '263_x', '264_x', '265_x', '266_x', '267_x', '268_x', '269_x', '270_x', '27
1_x', '272_x', '273_x', '274_x', '275_x', '276_x', '277_x', '278_x', '279_x', '280_x', '281_x', '282_x', '283
_x', '284_x', '285_x', '286_x', '287_x', '288_x', '289_x', '290_x', '291_x', '292_x', '293_x', '294_x', '295_x',
'296_x', '297_x', '298_x', '299_x', '300_x', '301_x', '302_x', '303_x', '304_x', '305_x', '306_x', '307_x', '
308_x', '309_x', '310_x', '311_x', '312_x', '313_x', '314_x', '315_x', '316_x', '317_x', '318_x', '319_x', '32
0_x', '321_x', '322_x', '323_x', '324_x', '325_x', '326_x', '327_x', '328_x', '329_x', '330_x', '331_x', '332
_x', '333_x', '334_x', '335_x', '336_x', '337_x', '338_x', '339_x', '340_x', '341_x', '342_x', '343_x', '344_x',
'345_x', '346_x', '347_x', '348_x', '349_x', '350_x', '351_x', '352_x', '353_x', '354_x', '355_x', '356_x', '
357_x', '358_x', '359_x', '360_x', '361_x', '362_x', '363_x', '364_x', '365_x', '366_x', '367_x', '368_x', '36
9_x', '370_x', '371_x', '372_x', '373_x', '374_x', '375_x', '376_x', '377_x', '378_x', '379_x', '380_x', '381
_x', '382_x', '383_x', '0_y', '1_y', '2_y', '3_y', '4_y', '5_y', '6_y', '7_y', '8_y', '9_y', '10_y', '11_y', '12_y',
'13_y', '14_y', '15_y', '16_y', '17_y', '18_y', '19_y', '20_y', '21_y', '22_y', '23_y', '24_y', '25_y', '26_y',
'27_y', '28_y', '29_y', '30_y', '31_y', '32_y', '33_y', '34_y', '35_y', '36_y', '37_y', '38_y', '39_y', '40_y', '
41_y', '42_y', '43_y', '44_y', '45_y', '46_y', '47_y', '48_y', '49_y', '50_y', '51_y', '52_y', '53_y', '54_y', '5
5_y', '56_y', '57_y', '58_y', '59_y', '60_y', '61_y', '62_y', '63_y', '64_y', '65_y', '66_y', '67_y', '68_y', '69
_y', '70_y', '71_y', '72_y', '73_y', '74_y', '75_y', '76_y', '77_y', '78_y', '79_y', '80_y', '81_y', '82_y', '83
_y', '84_y', '85_y', '86_y', '87_y', '88_y', '89_y', '90_y', '91_y', '92_y', '93_y', '94_y', '95_y', '96_y', '97_y
', '98_y', '99_y', '100_y', '101_y', '102_y', '103_y', '104_y', '105_y', '106_y', '107_y', '108_y', '109_y', '11
0_y', '111_y', '112_y', '113_y', '114_y', '115_y', '116_y', '117_y', '118_y', '119_y', '120_y', '121_y', '122
_y', '123_y', '124_y', '125_y', '126_y', '127_y', '128_y', '129_y', '130_y', '131_y', '132_y', '133_y', '134_y',
'135_y', '136_y', '137_y', '138_y', '139_y', '140_y', '141_y', '142_y', '143_y', '144_y', '145_y', '146_y', '1
47_y', '148_y', '149_y', '150_y', '151_y', '152_y', '153_y', '154_y', '155_y', '156_y', '157_y', '158_y', '15
9_y', '160_y', '161_y', '162_y', '163_y', '164_y', '165_y', '166_y', '167_y', '168_y', '169_y', '170_y', '171
_y', '172_y', '173_y', '174_y', '175_y', '176_y', '177_y', '178_y', '179_y', '180_y', '181_y', '182_y', '183_y',
'184_y', '185_y', '186_y', '187_y', '188_y', '189_y', '190_y', '191_y', '192_y', '193_y', '194_y', '195_y', '
196_y', '197_y', '198_y', '199_y', '200_y', '201_y', '202_y', '203_y', '204_y', '205_y', '206_y', '207_y', '20
8_y', '209_y', '210_y', '211_y', '212_y', '213_y', '214_y', '215_y', '216_y', '217_y', '218_y', '219_y', '220
_y', '221_y', '222_y', '223_y', '224_y', '225_y', '226_y', '227_y', '228_y', '229_y', '230_y', '231_y', '232_y',
'233_y', '234_y', '235_y', '236_y', '237_y', '238_y', '239_y', '240_y', '241_y', '242_y', '243_y', '244_y', '
245_y', '246_y', '247_y', '248_y', '249_y', '250_y', '251_y', '252_y', '253_y', '254_y', '255_y', '256_y', '25
7_y', '258_y', '259_y', '260_y', '261_y', '262_y', '263_y', '264_y', '265_y', '266_y', '267_y', '268_y', '269
_y', '270_y', '271_y', '272_y', '273_y', '274_y', '275_y', '276_y', '277_y', '278_y', '279_y', '280_y', '281_y',
'282_y', '283_y', '284_y', '285_y', '286_y',
```

```

0_y', '307_y', '308_y', '309_y', '310_y', '311_y', '312_y', '313_y', '314_y', '315_y', '316_y', '317_y', '318_
y', '319_y', '320_y', '321_y', '322_y', '323_y', '324_y', '325_y', '326_y', '327_y', '328_y', '329_y', '330_y',
'331_y', '332_y', '333_y', '334_y', '335_y', '336_y', '337_y', '338_y', '339_y', '340_y', '341_y', '342_y', '
343_y', '344_y', '345_y', '346_y', '347_y', '348_y', '349_y', '350_y', '351_y', '352_y', '353_y', '354_y', '35
5_y', '356_y', '357_y', '358_y', '359_y', '360_y', '361_y', '362_y', '363_y', '364_y', '365_y', '366_y', '367_
y', '368_y', '369_y', '370_y', '371_y', '372_y', '373_y', '374_y', '375_y', '376_y', '377_y', '378_y', '379_y'
, '380_y', '381_y', '382_y', '383_y'], chunksize=chunksize, iterator=True, encoding='utf-8', ):
    df.index += index_start
    j+=1
    print('{} rows'.format(j*chunksize))
    df.to_sql('data', disk_engine, if_exists='append')
    index_start = df.index[-1] + 1

```

180000 rows

In [14]:

```

#http://www.sqlitetutorial.net/sqlite-python/create-tables/
def create_connection(db_file):
    """ create a database connection to the SQLite database
        specified by db_file
    :param db_file: database file
    :return: Connection object or None
    """
    try:
        conn = sqlite3.connect(db_file)
        return conn
    except Error as e:
        print(e)

    return None

def checkTableExists(dbcon):
    cursr = dbcon.cursor()
    str = "select name from sqlite_master where type='table'"
    table_names = cursr.execute(str)
    print("Tables in the database:")
    tables = table_names.fetchall()
    print(tables[0][0])
    return len(tables)

```

In [15]:

```

read_db = 'train.db'
conn_r = create_connection(read_db)
checkTableExists(conn_r)
conn_r.close()

```

Tables in the database:  
data

In [16]:

```

# try to sample data according to the computing power you have
if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        # for selecting first 1M rows
        # data = pd.read_sql_query("""SELECT * FROM data LIMIT 100001;""", conn_r)

        # for selecting random points
        data = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 100001;", conn_r)
        conn_r.commit()
        conn_r.close()

```

In [17]:

```

read_db = 'test.db'
conn_r = create_connection(read_db)
checkTableExists(conn_r)
conn_r.close()

```

Tables in the database:  
data

In [18]:

```
# try to sample data according to the computing power you have
if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        # for selecting first 1M rows
        # data = pd.read_sql_query("""SELECT * FROM data LIMIT 100001;""", conn_r)

        # for selecting random points
        data_test = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 50001;", conn_r)
        conn_r.commit()
        conn_r.close()
```

In [19]:

```
# remove the first row
data.drop(data.index[0], inplace=True)
y_true = data['is_duplicate']
data.drop(['Unnamed: 0', 'id', 'index', 'is_duplicate'], axis=1, inplace=True)
```

In [20]:

```
data.head()
```

Out[20]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	l
1	1	0.599988000239995	0.599988000239995	0.66664444518516	0.66664444518516	0.624992187597655	0.62499
2	0	0.0	0.0	0.199996000079998	0.124998437519531	0.0714280612281341	0.05882
3	1	0.499987500312492	0.499987500312492	0.999966667777741	0.749981250468738	0.714275510349852	0.62499
4	0	0.222219753113854	0.105262603881032	0.818174380232907	0.52940865053735	0.43999824000704	0.23404
5	0	0.66664444518516	0.499987500312492	0.499987500312492	0.499987500312492	0.571420408279882	0.49999

5 rows × 794 columns

In [22]:

```
# remove the first row
data_test.drop(data_test.index[0], inplace=True)
y_test_true = data_test['is_duplicate']
data_test.drop(['Unnamed: 0', 'id', 'index', 'is_duplicate'], axis=1, inplace=True)
```

In [23]:

```
data_test.head()
```

Out[23]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	las
1	1	0.833319444675922	0.416663194473379	0.499987500312492	0.249996875039062	0.699993000069999	0.304346
2	0	0.999950002499875	0.399992000159997	0.499987500312492	0.399992000159997	0.666655555740738	0.399996
3	1	0.333329629670781	0.29999700003	0.285710204139941	0.285710204139941	0.294115916965194	0.294115
4	0	0.749981250468738	0.499991666805553	0.499987500312492	0.333327777870369	0.454541322351615	0.454541
5	1	0.66664444518516	0.66664444518516	0.249993750156246	0.249993750156246	0.428565306209911	0.374995

5 rows × 794 columns

## 4.2 Converting strings to numerics

In [25]:

```
# https://stackoverflow.com/questions/7368789/convert-all-strings-in-a-list-to-int
y_true = map(int, y_true.values)
```

In [26]:

```
# https://stackoverflow.com/questions/7368789/convert-all-strings-in-a-list-to-int
y_test_true = map(int, y_test_true.values)
```

In [27]:

```
# after we read from sql table each entry was read it as a string
# we convert all the features into numeric before we apply any model
cols = list(data.columns)
for i in cols:
    data[i] = data[i].apply(pd.to_numeric,errors='coerce')
    print(i)
```

```
cwc_min
cwc_max
csc_min
csc_max
ctc_min
ctc_max
last_word_eq
first_word_eq
abs_len_diff
mean_len
token_set_ratio
token_sort_ratio
fuzz_ratio
fuzz_partial_ratio
longest_substr_ratio
freq_qid1
freq_qid2
q1len
q2len
q1_n_words
q2_n_words
word_Common
word_Total
word_share
freq_q1+q2
freq_q1-q2
0_x
1_x
2_x
3_x
4_x
5_x
6_x
7_x
8_x
9_x
10_x
11_x
12_x
13_x
14_x
15_x
16_x
17_x
18_x
19_x
20_x
21_x
22_x
23_x
```

24\_x  
25\_x  
26\_x  
27\_x  
28\_x  
29\_x  
30\_x  
31\_x  
32\_x  
33\_x  
34\_x  
35\_x  
36\_x  
37\_x  
38\_x  
39\_x  
40\_x  
41\_x  
42\_x  
43\_x  
44\_x  
45\_x  
46\_x  
47\_x  
48\_x  
49\_x  
50\_x  
51\_x  
52\_x  
53\_x  
54\_x  
55\_x  
56\_x  
57\_x  
58\_x  
59\_x  
60\_x  
61\_x  
62\_x  
63\_x  
64\_x  
65\_x  
66\_x  
67\_x  
68\_x  
69\_x  
70\_x  
71\_x  
72\_x  
73\_x  
74\_x  
75\_x  
76\_x  
77\_x  
78\_x  
79\_x  
80\_x  
81\_x  
82\_x  
83\_x  
84\_x  
85\_x  
86\_x  
87\_x  
88\_x  
89\_x  
90\_x  
91\_x  
92\_x  
93\_x  
94\_x  
95\_x  
96\_x  
97\_x  
98\_x  
99\_x  
100\_x

101\_x  
102\_x  
103\_x  
104\_x  
105\_x  
106\_x  
107\_x  
108\_x  
109\_x  
110\_x  
111\_x  
112\_x  
113\_x  
114\_x  
115\_x  
116\_x  
117\_x  
118\_x  
119\_x  
120\_x  
121\_x  
122\_x  
123\_x  
124\_x  
125\_x  
126\_x  
127\_x  
128\_x  
129\_x  
130\_x  
131\_x  
132\_x  
133\_x  
134\_x  
135\_x  
136\_x  
137\_x  
138\_x  
139\_x  
140\_x  
141\_x  
142\_x  
143\_x  
144\_x  
145\_x  
146\_x  
147\_x  
148\_x  
149\_x  
150\_x  
151\_x  
152\_x  
153\_x  
154\_x  
155\_x  
156\_x  
157\_x  
158\_x  
159\_x  
160\_x  
161\_x  
162\_x  
163\_x  
164\_x  
165\_x  
166\_x  
167\_x  
168\_x  
169\_x  
170\_x  
171\_x  
172\_x  
173\_x  
174\_x  
175\_x  
176\_x  
177\_x



178\_x  
179\_x  
180\_x  
181\_x  
182\_x  
183\_x  
184\_x  
185\_x  
186\_x  
187\_x  
188\_x  
189\_x  
190\_x  
191\_x  
192\_x  
193\_x  
194\_x  
195\_x  
196\_x  
197\_x  
198\_x  
199\_x  
200\_x  
201\_x  
202\_x  
203\_x  
204\_x  
205\_x  
206\_x  
207\_x  
208\_x  
209\_x  
210\_x  
211\_x  
212\_x  
213\_x  
214\_x  
215\_x  
216\_x  
217\_x  
218\_x  
219\_x  
220\_x  
221\_x  
222\_x  
223\_x  
224\_x  
225\_x  
226\_x  
227\_x  
228\_x  
229\_x  
230\_x  
231\_x  
232\_x  
233\_x  
234\_x  
235\_x  
236\_x  
237\_x  
238\_x  
239\_x  
240\_x  
241\_x  
242\_x  
243\_x  
244\_x  
245\_x  
246\_x  
247\_x  
248\_x  
249\_x  
250\_x  
251\_x  
252\_x  
253\_x  
254\_x

255\_x  
256\_x  
257\_x  
258\_x  
259\_x  
260\_x  
261\_x  
262\_x  
263\_x  
264\_x  
265\_x  
266\_x  
267\_x  
268\_x  
269\_x  
270\_x  
271\_x  
272\_x  
273\_x  
274\_x  
275\_x  
276\_x  
277\_x  
278\_x  
279\_x  
280\_x  
281\_x  
282\_x  
283\_x  
284\_x  
285\_x  
286\_x  
287\_x  
288\_x  
289\_x  
290\_x  
291\_x  
292\_x  
293\_x  
294\_x  
295\_x  
296\_x  
297\_x  
298\_x  
299\_x  
300\_x  
301\_x  
302\_x  
303\_x  
304\_x  
305\_x  
306\_x  
307\_x  
308\_x  
309\_x  
310\_x  
311\_x  
312\_x  
313\_x  
314\_x  
315\_x  
316\_x  
317\_x  
318\_x  
319\_x  
320\_x  
321\_x  
322\_x  
323\_x  
324\_x  
325\_x  
326\_x  
327\_x  
328\_x  
329\_x  
330\_x  
331\_x

332\_x  
333\_x  
334\_x  
335\_x  
336\_x  
337\_x  
338\_x  
339\_x  
340\_x  
341\_x  
342\_x  
343\_x  
344\_x  
345\_x  
346\_x  
347\_x  
348\_x  
349\_x  
350\_x  
351\_x  
352\_x  
353\_x  
354\_x  
355\_x  
356\_x  
357\_x  
358\_x  
359\_x  
360\_x  
361\_x  
362\_x  
363\_x  
364\_x  
365\_x  
366\_x  
367\_x  
368\_x  
369\_x  
370\_x  
371\_x  
372\_x  
373\_x  
374\_x  
375\_x  
376\_x  
377\_x  
378\_x  
379\_x  
380\_x  
381\_x  
382\_x  
383\_x  
0\_y  
1\_y  
2\_y  
3\_y  
4\_y  
5\_y  
6\_y  
7\_y  
8\_y  
9\_y  
10\_y  
11\_y  
12\_y  
13\_y  
14\_y  
15\_y  
16\_y  
17\_y  
18\_y  
19\_y  
20\_y  
21\_y  
22\_y  
23\_y  
24\_y

25\_y  
26\_y  
27\_y  
28\_y  
29\_y  
30\_y  
31\_y  
32\_y  
33\_y  
34\_y  
35\_y  
36\_y  
37\_y  
38\_y  
39\_y  
40\_y  
41\_y  
42\_y  
43\_y  
44\_y  
45\_y  
46\_y  
47\_y  
48\_y  
49\_y  
50\_y  
51\_y  
52\_y  
53\_y  
54\_y  
55\_y  
56\_y  
57\_y  
58\_y  
59\_y  
60\_y  
61\_y  
62\_y  
63\_y  
64\_y  
65\_y  
66\_y  
67\_y  
68\_y  
69\_y  
70\_y  
71\_y  
72\_y  
73\_y  
74\_y  
75\_y  
76\_y  
77\_y  
78\_y  
79\_y  
80\_y  
81\_y  
82\_y  
83\_y  
84\_y  
85\_y  
86\_y  
87\_y  
88\_y  
89\_y  
90\_y  
91\_y  
92\_y  
93\_y  
94\_y  
95\_y  
96\_y  
97\_y  
98\_y  
99\_y  
100\_y  
101 v

102\_y  
103\_y  
104\_y  
105\_y  
106\_y  
107\_y  
108\_y  
109\_y  
110\_y  
111\_y  
112\_y  
113\_y  
114\_y  
115\_y  
116\_y  
117\_y  
118\_y  
119\_y  
120\_y  
121\_y  
122\_y  
123\_y  
124\_y  
125\_y  
126\_y  
127\_y  
128\_y  
129\_y  
130\_y  
131\_y  
132\_y  
133\_y  
134\_y  
135\_y  
136\_y  
137\_y  
138\_y  
139\_y  
140\_y  
141\_y  
142\_y  
143\_y  
144\_y  
145\_y  
146\_y  
147\_y  
148\_y  
149\_y  
150\_y  
151\_y  
152\_y  
153\_y  
154\_y  
155\_y  
156\_y  
157\_y  
158\_y  
159\_y  
160\_y  
161\_y  
162\_y  
163\_y  
164\_y  
165\_y  
166\_y  
167\_y  
168\_y  
169\_y  
170\_y  
171\_y  
172\_y  
173\_y  
174\_y  
175\_y  
176\_y  
177\_y  
178 v

179\_y  
180\_y  
181\_y  
182\_y  
183\_y  
184\_y  
185\_y  
186\_y  
187\_y  
188\_y  
189\_y  
190\_y  
191\_y  
192\_y  
193\_y  
194\_y  
195\_y  
196\_y  
197\_y  
198\_y  
199\_y  
200\_y  
201\_y  
202\_y  
203\_y  
204\_y  
205\_y  
206\_y  
207\_y  
208\_y  
209\_y  
210\_y  
211\_y  
212\_y  
213\_y  
214\_y  
215\_y  
216\_y  
217\_y  
218\_y  
219\_y  
220\_y  
221\_y  
222\_y  
223\_y  
224\_y  
225\_y  
226\_y  
227\_y  
228\_y  
229\_y  
230\_y  
231\_y  
232\_y  
233\_y  
234\_y  
235\_y  
236\_y  
237\_y  
238\_y  
239\_y  
240\_y  
241\_y  
242\_y  
243\_y  
244\_y  
245\_y  
246\_y  
247\_y  
248\_y  
249\_y  
250\_y  
251\_y  
252\_y  
253\_y  
254\_y  
255\_y

256\_y  
257\_y  
258\_y  
259\_y  
260\_y  
261\_y  
262\_y  
263\_y  
264\_y  
265\_y  
266\_y  
267\_y  
268\_y  
269\_y  
270\_y  
271\_y  
272\_y  
273\_y  
274\_y  
275\_y  
276\_y  
277\_y  
278\_y  
279\_y  
280\_y  
281\_y  
282\_y  
283\_y  
284\_y  
285\_y  
286\_y  
287\_y  
288\_y  
289\_y  
290\_y  
291\_y  
292\_y  
293\_y  
294\_y  
295\_y  
296\_y  
297\_y  
298\_y  
299\_y  
300\_y  
301\_y  
302\_y  
303\_y  
304\_y  
305\_y  
306\_y  
307\_y  
308\_y  
309\_y  
310\_y  
311\_y  
312\_y  
313\_y  
314\_y  
315\_y  
316\_y  
317\_y  
318\_y  
319\_y  
320\_y  
321\_y  
322\_y  
323\_y  
324\_y  
325\_y  
326\_y  
327\_y  
328\_y  
329\_y  
330\_y  
331\_y  
332\_y

```
332_y
333_y
334_y
335_y
336_y
337_y
338_y
339_y
340_y
341_y
342_y
343_y
344_y
345_y
346_y
347_y
348_y
349_y
350_y
351_y
352_y
353_y
354_y
355_y
356_y
357_y
358_y
359_y
360_y
361_y
362_y
363_y
364_y
365_y
366_y
367_y
368_y
369_y
370_y
371_y
372_y
373_y
374_y
375_y
376_y
377_y
378_y
379_y
380_y
381_y
382_y
383_y
```

In [28]:

```
# after we read from sql table each entry was read it as a string
# we convert all the features into numeric before we apply any model
cols = list(data_test.columns)
for i in cols:
    data_test[i] = data_test[i].apply(pd.to_numeric,errors='coerce')
    print(i)
```

```
cwc_min
cwc_max
csc_min
csc_max
ctc_min
ctc_max
last_word_eq
first_word_eq
abs_len_diff
mean_len
token_set_ratio
token_sort_ratio
fuzz_ratio
fuzz_partial_ratio
```



longest\_substr\_ratio  
freq\_qid1  
freq\_qid2  
q1len  
q2len  
q1\_n\_words  
q2\_n\_words  
word\_Common  
word\_Total  
word\_share  
freq\_q1+q2  
freq\_q1-q2  
0\_x  
1\_x  
2\_x  
3\_x  
4\_x  
5\_x  
6\_x  
7\_x  
8\_x  
9\_x  
10\_x  
11\_x  
12\_x  
13\_x  
14\_x  
15\_x  
16\_x  
17\_x  
18\_x  
19\_x  
20\_x  
21\_x  
22\_x  
23\_x  
24\_x  
25\_x  
26\_x  
27\_x  
28\_x  
29\_x  
30\_x  
31\_x  
32\_x  
33\_x  
34\_x  
35\_x  
36\_x  
37\_x  
38\_x  
39\_x  
40\_x  
41\_x  
42\_x  
43\_x  
44\_x  
45\_x  
46\_x  
47\_x  
48\_x  
49\_x  
50\_x  
51\_x  
52\_x  
53\_x  
54\_x  
55\_x  
56\_x  
57\_x  
58\_x  
59\_x  
60\_x  
61\_x  
62\_x  
63\_x  
64\_x

65\_x  
66\_x  
67\_x  
68\_x  
69\_x  
70\_x  
71\_x  
72\_x  
73\_x  
74\_x  
75\_x  
76\_x  
77\_x  
78\_x  
79\_x  
80\_x  
81\_x  
82\_x  
83\_x  
84\_x  
85\_x  
86\_x  
87\_x  
88\_x  
89\_x  
90\_x  
91\_x  
92\_x  
93\_x  
94\_x  
95\_x  
96\_x  
97\_x  
98\_x  
99\_x  
100\_x  
101\_x  
102\_x  
103\_x  
104\_x  
105\_x  
106\_x  
107\_x  
108\_x  
109\_x  
110\_x  
111\_x  
112\_x  
113\_x  
114\_x  
115\_x  
116\_x  
117\_x  
118\_x  
119\_x  
120\_x  
121\_x  
122\_x  
123\_x  
124\_x  
125\_x  
126\_x  
127\_x  
128\_x  
129\_x  
130\_x  
131\_x  
132\_x  
133\_x  
134\_x  
135\_x  
136\_x  
137\_x  
138\_x  
139\_x  
140\_x  
141\_x

142\_x  
143\_x  
144\_x  
145\_x  
146\_x  
147\_x  
148\_x  
149\_x  
150\_x  
151\_x  
152\_x  
153\_x  
154\_x  
155\_x  
156\_x  
157\_x  
158\_x  
159\_x  
160\_x  
161\_x  
162\_x  
163\_x  
164\_x  
165\_x  
166\_x  
167\_x  
168\_x  
169\_x  
170\_x  
171\_x  
172\_x  
173\_x  
174\_x  
175\_x  
176\_x  
177\_x  
178\_x  
179\_x  
180\_x  
181\_x  
182\_x  
183\_x  
184\_x  
185\_x  
186\_x  
187\_x  
188\_x  
189\_x  
190\_x  
191\_x  
192\_x  
193\_x  
194\_x  
195\_x  
196\_x  
197\_x  
198\_x  
199\_x  
200\_x  
201\_x  
202\_x  
203\_x  
204\_x  
205\_x  
206\_x  
207\_x  
208\_x  
209\_x  
210\_x  
211\_x  
212\_x  
213\_x  
214\_x  
215\_x  
216\_x  
217\_x  
218\_x

219\_x  
220\_x  
221\_x  
222\_x  
223\_x  
224\_x  
225\_x  
226\_x  
227\_x  
228\_x  
229\_x  
230\_x  
231\_x  
232\_x  
233\_x  
234\_x  
235\_x  
236\_x  
237\_x  
238\_x  
239\_x  
240\_x  
241\_x  
242\_x  
243\_x  
244\_x  
245\_x  
246\_x  
247\_x  
248\_x  
249\_x  
250\_x  
251\_x  
252\_x  
253\_x  
254\_x  
255\_x  
256\_x  
257\_x  
258\_x  
259\_x  
260\_x  
261\_x  
262\_x  
263\_x  
264\_x  
265\_x  
266\_x  
267\_x  
268\_x  
269\_x  
270\_x  
271\_x  
272\_x  
273\_x  
274\_x  
275\_x  
276\_x  
277\_x  
278\_x  
279\_x  
280\_x  
281\_x  
282\_x  
283\_x  
284\_x  
285\_x  
286\_x  
287\_x  
288\_x  
289\_x  
290\_x  
291\_x  
292\_x  
293\_x  
294\_x  
295\_x

296\_x  
297\_x  
298\_x  
299\_x  
300\_x  
301\_x  
302\_x  
303\_x  
304\_x  
305\_x  
306\_x  
307\_x  
308\_x  
309\_x  
310\_x  
311\_x  
312\_x  
313\_x  
314\_x  
315\_x  
316\_x  
317\_x  
318\_x  
319\_x  
320\_x  
321\_x  
322\_x  
323\_x  
324\_x  
325\_x  
326\_x  
327\_x  
328\_x  
329\_x  
330\_x  
331\_x  
332\_x  
333\_x  
334\_x  
335\_x  
336\_x  
337\_x  
338\_x  
339\_x  
340\_x  
341\_x  
342\_x  
343\_x  
344\_x  
345\_x  
346\_x  
347\_x  
348\_x  
349\_x  
350\_x  
351\_x  
352\_x  
353\_x  
354\_x  
355\_x  
356\_x  
357\_x  
358\_x  
359\_x  
360\_x  
361\_x  
362\_x  
363\_x  
364\_x  
365\_x  
366\_x  
367\_x  
368\_x  
369\_x  
370\_x  
371\_x  
372\_x

373\_x  
374\_x  
375\_x  
376\_x  
377\_x  
378\_x  
379\_x  
380\_x  
381\_x  
382\_x  
383\_x  
0\_y  
1\_y  
2\_y  
3\_y  
4\_y  
5\_y  
6\_y  
7\_y  
8\_y  
9\_y  
10\_y  
11\_y  
12\_y  
13\_y  
14\_y  
15\_y  
16\_y  
17\_y  
18\_y  
19\_y  
20\_y  
21\_y  
22\_y  
23\_y  
24\_y  
25\_y  
26\_y  
27\_y  
28\_y  
29\_y  
30\_y  
31\_y  
32\_y  
33\_y  
34\_y  
35\_y  
36\_y  
37\_y  
38\_y  
39\_y  
40\_y  
41\_y  
42\_y  
43\_y  
44\_y  
45\_y  
46\_y  
47\_y  
48\_y  
49\_y  
50\_y  
51\_y  
52\_y  
53\_y  
54\_y  
55\_y  
56\_y  
57\_y  
58\_y  
59\_y  
60\_y  
61\_y  
62\_y  
63\_y  
64\_y  
65\_v

66\_y  
67\_y  
68\_y  
69\_y  
70\_y  
71\_y  
72\_y  
73\_y  
74\_y  
75\_y  
76\_y  
77\_y  
78\_y  
79\_y  
80\_y  
81\_y  
82\_y  
83\_y  
84\_y  
85\_y  
86\_y  
87\_y  
88\_y  
89\_y  
90\_y  
91\_y  
92\_y  
93\_y  
94\_y  
95\_y  
96\_y  
97\_y  
98\_y  
99\_y  
100\_y  
101\_y  
102\_y  
103\_y  
104\_y  
105\_y  
106\_y  
107\_y  
108\_y  
109\_y  
110\_y  
111\_y  
112\_y  
113\_y  
114\_y  
115\_y  
116\_y  
117\_y  
118\_y  
119\_y  
120\_y  
121\_y  
122\_y  
123\_y  
124\_y  
125\_y  
126\_y  
127\_y  
128\_y  
129\_y  
130\_y  
131\_y  
132\_y  
133\_y  
134\_y  
135\_y  
136\_y  
137\_y  
138\_y  
139\_y  
140\_y  
141\_y  
142\_y

143\_y  
144\_y  
145\_y  
146\_y  
147\_y  
148\_y  
149\_y  
150\_y  
151\_y  
152\_y  
153\_y  
154\_y  
155\_y  
156\_y  
157\_y  
158\_y  
159\_y  
160\_y  
161\_y  
162\_y  
163\_y  
164\_y  
165\_y  
166\_y  
167\_y  
168\_y  
169\_y  
170\_y  
171\_y  
172\_y  
173\_y  
174\_y  
175\_y  
176\_y  
177\_y  
178\_y  
179\_y  
180\_y  
181\_y  
182\_y  
183\_y  
184\_y  
185\_y  
186\_y  
187\_y  
188\_y  
189\_y  
190\_y  
191\_y  
192\_y  
193\_y  
194\_y  
195\_y  
196\_y  
197\_y  
198\_y  
199\_y  
200\_y  
201\_y  
202\_y  
203\_y  
204\_y  
205\_y  
206\_y  
207\_y  
208\_y  
209\_y  
210\_y  
211\_y  
212\_y  
213\_y  
214\_y  
215\_y  
216\_y  
217\_y  
218\_y  
219\_y



219\_y  
220\_y  
221\_y  
222\_y  
223\_y  
224\_y  
225\_y  
226\_y  
227\_y  
228\_y  
229\_y  
230\_y  
231\_y  
232\_y  
233\_y  
234\_y  
235\_y  
236\_y  
237\_y  
238\_y  
239\_y  
240\_y  
241\_y  
242\_y  
243\_y  
244\_y  
245\_y  
246\_y  
247\_y  
248\_y  
249\_y  
250\_y  
251\_y  
252\_y  
253\_y  
254\_y  
255\_y  
256\_y  
257\_y  
258\_y  
259\_y  
260\_y  
261\_y  
262\_y  
263\_y  
264\_y  
265\_y  
266\_y  
267\_y  
268\_y  
269\_y  
270\_y  
271\_y  
272\_y  
273\_y  
274\_y  
275\_y  
276\_y  
277\_y  
278\_y  
279\_y  
280\_y  
281\_y  
282\_y  
283\_y  
284\_y  
285\_y  
286\_y  
287\_y  
288\_y  
289\_y  
290\_y  
291\_y  
292\_y  
293\_y  
294\_y  
295\_y  
296\_y

296\_y  
297\_y  
298\_y  
299\_y  
300\_y  
301\_y  
302\_y  
303\_y  
304\_y  
305\_y  
306\_y  
307\_y  
308\_y  
309\_y  
310\_y  
311\_y  
312\_y  
313\_y  
314\_y  
315\_y  
316\_y  
317\_y  
318\_y  
319\_y  
320\_y  
321\_y  
322\_y  
323\_y  
324\_y  
325\_y  
326\_y  
327\_y  
328\_y  
329\_y  
330\_y  
331\_y  
332\_y  
333\_y  
334\_y  
335\_y  
336\_y  
337\_y  
338\_y  
339\_y  
340\_y  
341\_y  
342\_y  
343\_y  
344\_y  
345\_y  
346\_y  
347\_y  
348\_y  
349\_y  
350\_y  
351\_y  
352\_y  
353\_y  
354\_y  
355\_y  
356\_y  
357\_y  
358\_y  
359\_y  
360\_y  
361\_y  
362\_y  
363\_y  
364\_y  
365\_y  
366\_y  
367\_y  
368\_y  
369\_y  
370\_y  
371\_y  
372\_y  
373\_y

```
373_y
374_y
375_y
376_y
377_y
378_y
379_y
380_y
381_y
382_y
383_y
```

## 4.3 Random train test split( 70:30)

In [2]:

```
data_train=pd.read_csv('final_features_train.csv')
data_test=pd.read_csv('final_features_test.csv')
data_train.drop(data_train.index[0], inplace=True)
y_true = data_train['is_duplicate'].values
data_train.drop(['Unnamed: 0_x', 'Unnamed: 0', 'id', 'is_duplicate'], axis=1, inplace=True)
data_test.drop(data_test.index[0], inplace=True)
y_test_true = data_test['is_duplicate'].values
data_test.drop(['Unnamed: 0_x', 'Unnamed: 0', 'id', 'is_duplicate'], axis=1, inplace=True)
```

In [3]:

```
#training data got from final_features_test
X_train=data_train[:70000].values
X_test=data_test[:30000].values
y_train=y_true[:70000]
y_test=y_test_true[:30000]
```

In [4]:

```
print("Shpae of train data",X_train.shape,y_train.shape)
print("Shape of test data",X_test.shape,y_test.shape)
```

```
Shpae of train data (70000, 219) (70000,)
Shape of test data (30000, 219) (30000,)
```

In [5]:

```
print("-"*10, "Distribution of output variable in train data", "-"*10)
train_distr = Counter(y_train)
train_len = len(y_train)
print("Class 0: ",int(train_distr[0])/train_len,"Class 1: ", int(train_distr[1])/train_len)
print("-"*10, "Distribution of output variable in train data", "-"*10)
test_distr = Counter(y_test)
test_len = len(y_test)
print("Class 0: ",int(test_distr[1])/test_len, "Class 1: ",int(test_distr[1])/test_len)
```

```
----- Distribution of output variable in train data -----
Class 0:  0.6335571428571428 Class 1:  0.3664428571428571
----- Distribution of output variable in train data -----
Class 0:  0.36986666666666667 Class 1:  0.3698666666666667
```

In [25]:

```
# This function plots the confusion matrices given y_i, y_i_hat.
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j

    A = ((C.T) / (C.sum(axis=1))).T
    #divid each element of the confusion matrix with the sum of elements in that column

    # C = [[1, 2],
    #       [3, 4]]
    # C.T = [[1, 3],
```

```

#         [2, 4]]
# C.sum(axis = 1)  axis=0 corresponds to columns and axis=1 corresponds to rows in two
dimensional array
# C.sum(axix =1) = [[3, 7]]
# ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
#                             [2/3, 4/7]]

# ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
#                               [3/7, 4/7]]
# sum of row elements = 1

B=(C/C.sum(axis=0))
#divid each element of the confusion matrix with the sum of elements in that row
# C = [[1, 2],
#       [3, 4]]
# C.sum(axis = 0)  axis=0 corresponds to columns and axis=1 corresponds to rows in two
dimensional array
# C.sum(axix =0) = [[4, 6]]
# (C/C.sum(axis=0)) = [[1/4, 2/6],
#                       [3/4, 4/6]]
plt.figure(figsize=(20,4))

labels = [1,2]
# 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")
print("Confusion matrix")
df_cm = pd.DataFrame(C,range(2),range(2))
df_cm.columns = ['Predicted NO','Predicted YES']
df_cm = df_cm.rename({0: 'Actual NO', 1: 'Actual YES'})
print(df_cm)

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")
print("Precision matrix")
df_cm = pd.DataFrame(B,range(2),range(2))
df_cm.columns = ['Predicted NO','Predicted YES']
df_cm = df_cm.rename({0: 'Actual NO', 1: 'Actual YES'})
print(df_cm)

plt.subplot(1, 3, 3)
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")
print("Recall matrix")
df_cm = pd.DataFrame(B,range(2),range(2))
df_cm.columns = ['Predicted NO','Predicted YES']
df_cm = df_cm.rename({0: 'Actual NO', 1: 'Actual YES'})
print(df_cm)
plt.show()

```

## 4.4 Building a random model (Finding worst-case log-loss)

In [124]:

```

# we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to generate 9 numbers and divide each of the numbers by their sum
# ref: https://stackoverflow.com/a/18662466/4084039
# we create a output array that has exactly same size as the CV data
predicted_y = np.zeros((test_len,2))
for i in range(test_len):
    rand_probs = np.random.rand(1,2)
    predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y, eps=1e-15))

```

```
predicted_y = np.argmax(predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y)
```

Log loss on Test Data using Random Model 0.8879982259662449

Confusion matrix

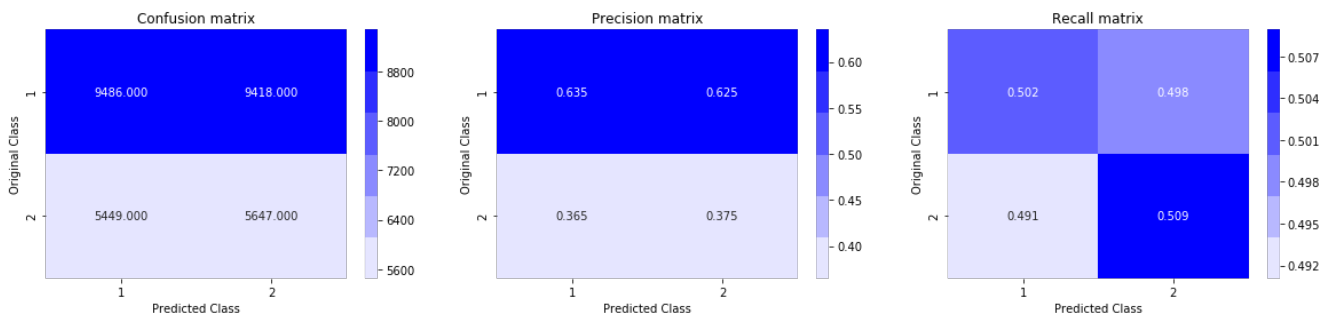
	Predicted NO	Predicted YES
Actual NO	9486	9418
Actual YES	5449	5647

Precision matrix

	Predicted NO	Predicted YES
Actual NO	0.635152	0.625158
Actual YES	0.364848	0.374842

Recall matrix

	Predicted NO	Predicted YES
Actual NO	0.635152	0.625158
Actual YES	0.364848	0.374842



## 4.4 Logistic Regression with hyperparameter tuning

In [138]:

```
%%time
#warnings.simplefilter("ignore", 'FutureWarning')
alpha = [10 ** x for x in range(-7, 7)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', class_weight='balanced', n_jobs=-1, learning_rate='adaptive', eta0=5)
    clf.fit(X_train, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="isotonic", cv=2)
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_test)
    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 2)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
```

```

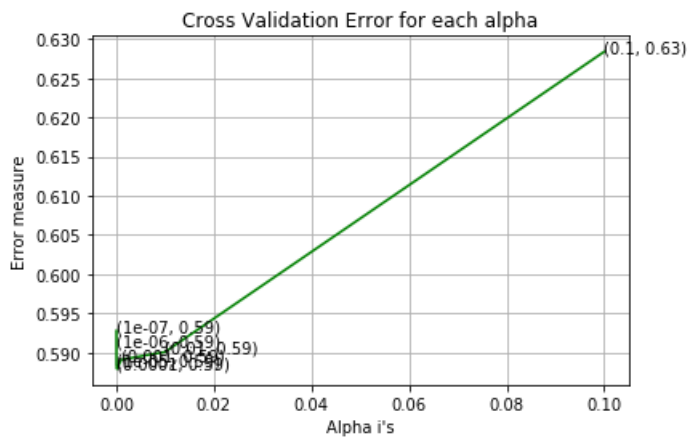
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2',
loss='log',class_weight='balanced',n_jobs=-1,learning_rate='adaptive',eta0 =5)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="isotonic",cv=2)
sig_clf.fit(X_train, y_train)

predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, p
redict_y, labels=clf.classes_, eps=1e-15))
predicted_y =np.argmax(predict_y,axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)

```

For values of alpha = 1e-07 The log loss is: 0.5927904576378324  
 For values of alpha = 1e-06 The log loss is: 0.5909443528510991  
 For values of alpha = 1e-05 The log loss is: 0.5884413243175798  
 For values of alpha = 0.0001 The log loss is: 0.5879319408614725  
 For values of alpha = 0.001 The log loss is: 0.5891032254351823  
 For values of alpha = 0.01 The log loss is: 0.5901007570784752  
 For values of alpha = 0.1 The log loss is: 0.628373793051189



For values of best alpha = 0.0001 The train log loss is: 0.5808928707970528  
 For values of best alpha = 0.0001 The test log loss is: 0.5908179212138209  
 Total number of data points : 30000

Confusion matrix

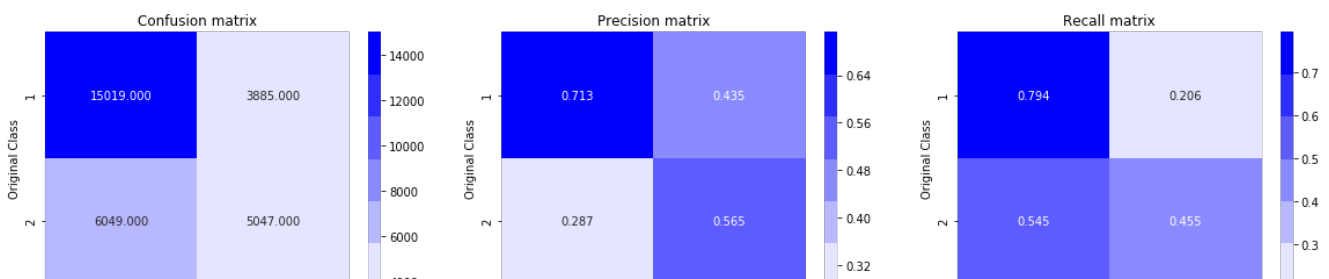
	Predicted NO	Predicted YES
Actual NO	15019	3885
Actual YES	6049	5047

Precision matrix

	Predicted NO	Predicted YES
Actual NO	0.712882	0.434953
Actual YES	0.287118	0.565047

Recall matrix

	Predicted NO	Predicted YES
Actual NO	0.712882	0.434953
Actual YES	0.287118	0.565047





Wall time: 4min 59s

## 4.5 Linear SVM with hyperparameter tuning

In [142]:

```
alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-
learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_i
ter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0
=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----
#alpha=i, penalty='l2', loss='log',class_weight='balanced',n_jobs=-1,learning_rate
='adaptive',eta0 =5

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='hinge',class_weight='balanced',n_jobs=-1,learn
ing_rate='adaptive',eta0 =5)
    clf.fit(X_train, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="isotonic",cv=2)
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_test)
    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_test, predict_y, labels=clf.cl
asses_, eps=1e-15))

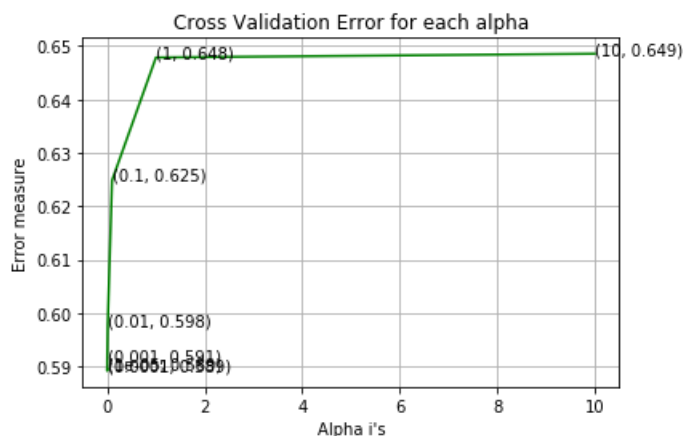
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array,c='g')
for i, txt in enumerate(np.round(log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")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

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2',
loss='hinge',class_weight='balanced',n_jobs=-1,learning_rate='adaptive',eta0 =5)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="isotonic",cv=2)
sig_clf.fit(X_train, y_train)

predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, p
redict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y,axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)
```

For values of alpha = 1e-05 The log loss is: 0.5894821727136259  
For values of alpha = 0.0001 The log loss is: 0.5891677952937698  
For values of alpha = 0.001 The log loss is: 0.5911466047808938

For values of alpha = 0.0001 The log loss is: 0.5832719717267644  
 For values of alpha = 0.01 The log loss is: 0.5976186549930591  
 For values of alpha = 0.1 The log loss is: 0.6248539574278599  
 For values of alpha = 1 The log loss is: 0.6478512912941948  
 For values of alpha = 10 The log loss is: 0.6485579315566874



For values of best alpha = 0.0001 The train log loss is: 0.5832719717267644  
 For values of best alpha = 0.0001 The test log loss is: 0.5931414884160295  
 Total number of data points : 30000

Confusion matrix

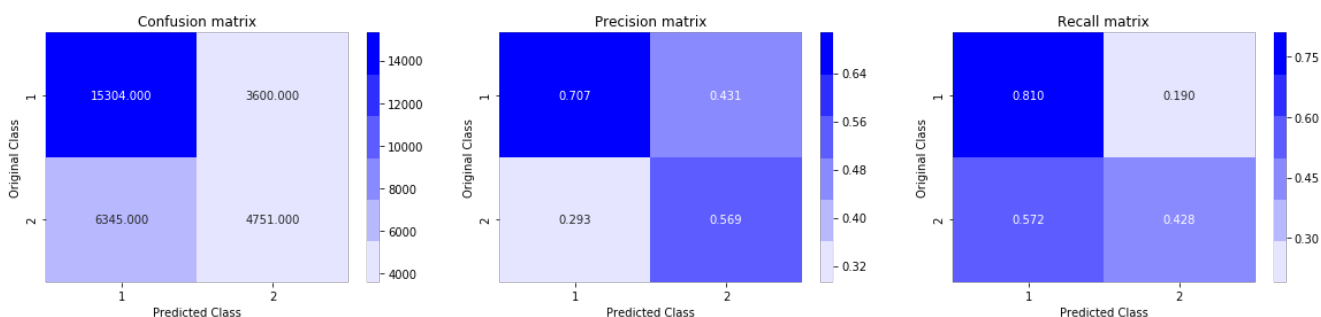
	Predicted NO	Predicted YES
Actual NO	15304	3600
Actual YES	6345	4751

Precision matrix

	Predicted NO	Predicted YES
Actual NO	0.706915	0.431086
Actual YES	0.293085	0.568914

Recall matrix

	Predicted NO	Predicted YES
Actual NO	0.706915	0.431086
Actual YES	0.293085	0.568914



## 5. Assignments

### 1. Try out models (Logistic regression, Linear-SVM) with simple TF-IDF vectors instead of TD\_IDF weighted word2Vec.

In [6]:

```
#Reading the
data_train=pd.read_csv('final_features_train.csv')
data_test=pd.read_csv('final_features_test.csv')
data_train.drop(data_train.index[0], inplace=True)
y_train_tfidf = data_train['is_duplicate']
data_train.drop(['Unnamed: 0_x', 'Unnamed: 0', 'is_duplicate'], axis=1, inplace=True)
data_test.drop(data_test.index[0], inplace=True)
y_test_tfidf = data_test['is_duplicate']
data_test.drop(['Unnamed: 0_x', 'Unnamed: 0', 'is_duplicate'], axis=1, inplace=True)
```

In [7]:



```
#removing the column which not useful from train data set.
data_train.drop(['0_x','1_x','2_x','3_x','4_x','5_x','6_x','7_x','8_x','9_x','10_x','11_x','12_x','13_x','14_x','15_x','16_x','17_x','18_x','19_x','20_x','21_x','22_x','23_x','24_x','25_x','26_x','27_x','28_x','29_x','30_x','31_x','32_x','33_x','34_x','35_x','36_x','37_x','38_x','39_x','40_x','41_x','42_x','43_x','44_x','45_x','46_x','47_x','48_x','49_x','50_x','51_x','52_x','53_x','54_x','55_x','56_x','57_x','58_x','59_x','60_x','61_x','62_x','63_x','64_x','65_x','66_x','67_x','68_x','69_x','70_x','71_x','72_x','73_x','74_x','75_x','76_x','77_x','78_x','79_x','80_x','81_x','82_x','83_x','84_x','85_x','86_x','87_x','88_x','89_x','90_x','91_x','92_x','93_x','94_x','95_x','0_y','1_y','2_y','3_y','4_y','5_y','6_y','7_y','8_y','9_y','10_y','11_y','12_y','13_y','14_y','15_y','16_y','17_y','18_y','19_y','20_y','21_y','22_y','23_y','24_y','25_y','26_y','27_y','28_y','29_y','30_y','31_y','32_y','33_y','34_y','35_y','36_y','37_y','38_y','39_y','40_y','41_y','42_y','43_y','44_y','45_y','46_y','47_y','48_y','49_y','50_y','51_y','52_y','53_y','54_y','55_y','56_y','57_y','58_y','59_y','60_y','61_y','62_y','63_y','64_y','65_y','66_y','67_y','68_y','69_y','70_y','71_y','72_y','73_y','74_y','75_y','76_y','77_y','78_y','79_y','80_y','81_y','82_y','83_y','84_y','85_y','86_y','87_y','88_y','89_y','90_y','91_y','92_y','93_y','94_y','95_y'], axis=1, inplace=True)
#removing the column which not useful from train data set.
data_test.drop(['0_x','1_x','2_x','3_x','4_x','5_x','6_x','7_x','8_x','9_x','10_x','11_x','12_x','13_x','14_x','15_x','16_x','17_x','18_x','19_x','20_x','21_x','22_x','23_x','24_x','25_x','26_x','27_x','28_x','29_x','30_x','31_x','32_x','33_x','34_x','35_x','36_x','37_x','38_x','39_x','40_x','41_x','42_x','43_x','44_x','45_x','46_x','47_x','48_x','49_x','50_x','51_x','52_x','53_x','54_x','55_x','56_x','57_x','58_x','59_x','60_x','61_x','62_x','63_x','64_x','65_x','66_x','67_x','68_x','69_x','70_x','71_x','72_x','73_x','74_x','75_x','76_x','77_x','78_x','79_x','80_x','81_x','82_x','83_x','84_x','85_x','86_x','87_x','88_x','89_x','90_x','91_x','92_x','93_x','94_x','95_y'], axis=1, inplace=True)
```

In [8]:

```
data_train[:2]
```

Out[8]:

	id	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	...	freq_
1	93043	0.499988	0.499988	0.249994	0.249994	0.374995	0.374995	1.0	1.0	0.0	...	2
2	261901	0.833319	0.714276	0.749981	0.749981	0.799992	0.727266	1.0	1.0	1.0	...	6

2 rows × 28 columns



In [9]:

```
data_test[:2]
```

Out[9]:

	id	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	...	freq_
1	344312	0.749981	0.599988	0.666644	0.499988	0.624992	0.624992	0.0	0.0	0.0	...	1
2	290852	0.999900	0.499975	0.999950	0.999950	0.999967	0.749981	1.0	1.0	1.0	...	1

2 rows × 28 columns



In [10]:

```
#distribution of data set.
y_train_tfidf.value_counts()
```

Out[10]:

```
0    178654
1     104348
Name: is duplicate, dtype: int64
```

In [11]:

```
y_test_tfidf.value_counts()
```

Out[11]:

```
0    76371
1    44915
```

Name: is\_duplicate, dtype: int64

### Imbalance dataset.

In [12]:

```
#Reading nlp_features_train and removing some columns which is already there.
df=pd.read_csv('nlp_features_train.csv',encoding='latin-1')
df.drop(['is_duplicate','cwc_min','cwc_max','csc_min','csc_max','ctc_min','ctc_max','last_word_eq',
'first_word_eq','abs_len_diff','mean_len','token_set_ratio','token_sort_ratio','fuzz_ratio','fuzz_
partial_ratio','longest_substr_ratio'],axis=1,inplace=True)
data_train=df.merge(data_train,on='id', how='right')
#data looks like
# 'Unnamed:0','id', 'qid1', 'qid2', 'question1', 'question2', 'index', 'cwc_min',
# 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq',
# 'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_ratio',
# 'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio',
# 'longest_substr_ratio', 'freq_qid1', 'freq_qid2', 'q1len', 'q2len',
# 'q1_n_words', 'q2_n_words', 'word_Common', 'word_Total', 'word_share',
# 'freq_q1+q2', 'freq_q1-q2'
```

In [13]:

```
data_train.columns
```

Out[13]:

```
Index(['Unnamed: 0', 'id', 'qid1', 'qid2', 'question1', 'question2', 'cwc_min',
'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq',
'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_ratio',
'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio',
'longest_substr_ratio', 'Unnamed: 0_y', 'freq_qid1', 'freq_qid2',
'q1len', 'q2len', 'q1_n_words', 'q2_n_words', 'word_Common',
'word_Total', 'word_share', 'freq_q1+q2', 'freq_q1-q2'],
dtype='object')
```

In [14]:

```
#Reading nlp_features_train and removing some columns which is already there.
df_test=pd.read_csv('nlp_features_test.csv',encoding='latin-1')
df_test.drop(['is_duplicate','cwc_min','cwc_max','csc_min','csc_max','ctc_min','ctc_max','last_wor
d_eq','first_word_eq','abs_len_diff','mean_len','token_set_ratio','token_sort_ratio','fuzz_ratio',
'fuzz_partial_ratio','longest_substr_ratio'],axis=1,inplace=True)
data_test=df_test.merge(data_test,on='id', how='right')
#data looks like
# 'Unnamed:0','id', 'qid1', 'qid2', 'question1', 'question2', 'index', 'cwc_min',
# 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq',
# 'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_ratio',
# 'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio',
# 'longest_substr_ratio', 'freq_qid1', 'freq_qid2', 'q1len', 'q2len',
# 'q1_n_words', 'q2_n_words', 'word_Common', 'word_Total', 'word_share',
# 'freq_q1+q2', 'freq_q1-q2'
```

In [15]:

```
data_test.columns
```

Out[15]:

```
Index(['Unnamed: 0', 'id', 'qid1', 'qid2', 'question1', 'question2', 'cwc_min',
```

```

'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq',
'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_ratio',
'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio',
'longest_substr_ratio', 'Unnamed: 0_y', 'freq_qid1', 'freq_qid2',
'qlen', 'q2len', 'q1_n_words', 'q2_n_words', 'word_Common',
'word_Total', 'word_share', 'freq_q1+q2', 'freq_q1-q2'],
dtype='object')

```

In [16]:

```

from sklearn.feature_extraction.text import TfidfVectorizer
tfidf_vectorizer_q1 = TfidfVectorizer(min_df=30,max_features=1000)
#Fitting train data and transforming train ,cv and test vector shape should be same.
tr_q1_tfidf=tfidf_vectorizer_q1.fit_transform(data_train['question1'].values.astype('U'))
te_q1_tfidf=tfidf_vectorizer_q1.transform(data_test['question1'].values.astype('U'))

print("Shape of matrix TFIDF Vectorizer on question1 \nTrain data-{},\nTest
data-{}".format(tr_q1_tfidf.shape,te_q1_tfidf.shape))

```

```

Shape of matrix TFIDF Vectorizer on question1
Train data-(283002, 1000),
Test data-(121286, 1000)

```

In [17]:

```

tfidf_vectorizer_q2 = TfidfVectorizer(min_df=30,max_features=1000)
#Fitting train data and transforming train ,cv and test vector shape should be same.
tr_q2_tfidf=tfidf_vectorizer_q2.fit_transform(data_train['question2'].values.astype('U'))
te_q2_tfidf=tfidf_vectorizer_q2.transform(data_test['question2'].values.astype('U'))

print("Shape of matrix TFIDF Vectorizer on question2 \nTrain data-{},\nTest
data-{}".format(tr_q2_tfidf.shape,te_q2_tfidf.shape))

```

```

Shape of matrix TFIDF Vectorizer on question2
Train data-(283002, 1000),
Test data-(121286, 1000)

```

In [18]:

```

#Removing the the id,qid1,qid2,question1,question2,index once we get TfidfVectorization for ques
tion1 and question2.
data_train.drop(['Unnamed: 0','id','qid1','qid2','question1','question2'],axis=1,inplace=True)
data_test.drop(['Unnamed: 0','id','qid1','qid2','question1','question2'],axis=1,inplace=True)

```

In [19]:

```

#combining all the feature and tfidf features
tr_X_TFIDF=hstack((data_train,tr_q1_tfidf,tr_q2_tfidf))
tr_X_TFIDF=tr_X_TFIDF.toarray()
te_X_TFIDF= hstack((data_test,te_q1_tfidf,te_q2_tfidf))
te_X_TFIDF=te_X_TFIDF.toarray()
print(tr_X_TFIDF.shape)
print(te_X_TFIDF.shape)

```

```

(283002, 2027)
(121286, 2027)

```

## Logistic regression on TFIDF

In [27]:

```

%%time
alpha = [10 ** x for x in range(-2, 2)] # hyperparam for SGD classifier.
import xgboost as xgb
from xgboost.sklearn import XGBClassifier
# read more about SGDClassifier() at http://scikit-
learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters

```

```

# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0
# =0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log',
random_state=2,class_weight='balanced',n_jobs=-1)
    clf.fit(tr_X_TFIDF, y_train_tfidf)
    sig_clf = CalibratedClassifierCV(clf, method='isotonic',cv=3)
    sig_clf.fit(tr_X_TFIDF, y_train_tfidf)
    predict_y = sig_clf.predict_proba(tr_X_TFIDF)
    log_error_array.append(log_loss(y_train_tfidf, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_train_tfidf, predict_y,
labels=clf.classes_, eps=1e-15))

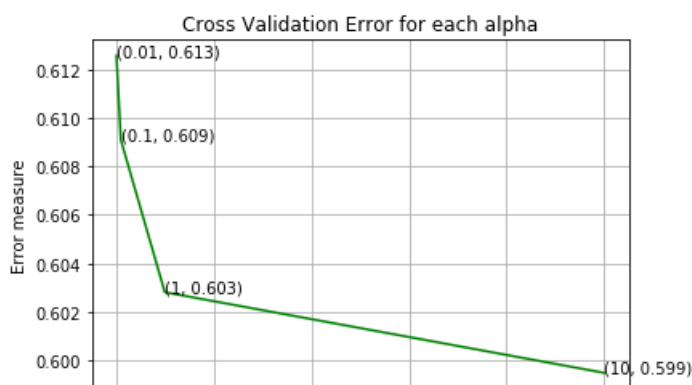
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array,c='g')
for i, txt in enumerate(np.round(log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")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

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=i, penalty='l2', loss='log',
random_state=2,class_weight='balanced',n_jobs=-1)
#clf=XGBClassifier( n_estimators=20,objective= 'binary:logistic',seed=27)
clf.fit(tr_X_TFIDF, y_train_tfidf)
sig_clf = CalibratedClassifierCV(clf, method='isotonic',cv=3)
sig_clf.fit(tr_X_TFIDF, y_train_tfidf)

predict_y = sig_clf.predict_proba(tr_X_TFIDF)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss
is:",log_loss(y_train_tfidf, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(tr_X_TFIDF)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss
is:",log_loss(y_test_tfidf, predict_y, labels=clf.classes_, eps=1e-15))
predicted_y =np.argmax(predict_y,axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test_tfidf, predicted_y)

```

For values of alpha = 0.01 The log loss is: 0.6125677429101645  
 For values of alpha = 0.1 The log loss is: 0.6090948565380095  
 For values of alpha = 1 The log loss is: 0.6028114538603679  
 For values of alpha = 10 The log loss is: 0.5994850612322475



0 2 4 6 8 10  
Alpha i's

For values of best alpha = 10 The train log loss is: 0.6192971909398267  
For values of best alpha = 10 The test log loss is: 0.6224645448213041  
Total number of data points : 121286  
Confusion matrix

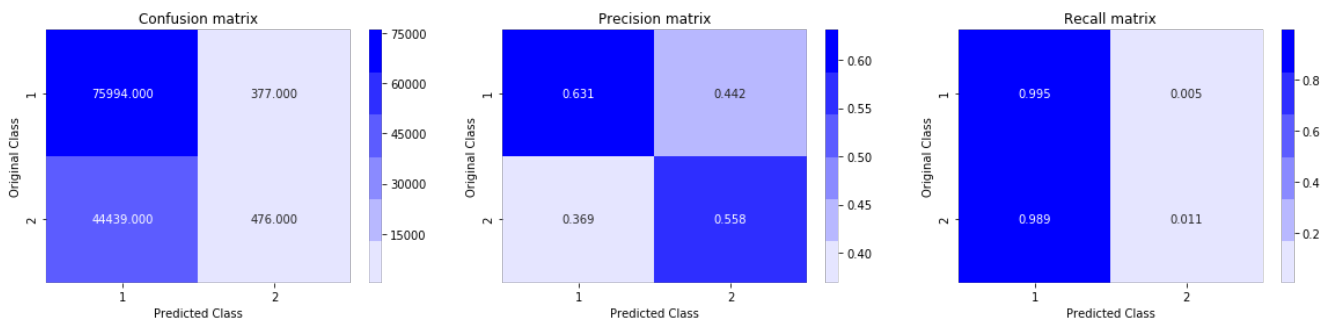
	Predicted NO	Predicted YES
Actual NO	75994	377
Actual YES	44439	476

Precision matrix

	Predicted NO	Predicted YES
Actual NO	0.631006	0.44197
Actual YES	0.368994	0.55803

Recall matrix

	Predicted NO	Predicted YES
Actual NO	0.631006	0.44197
Actual YES	0.368994	0.55803



Wall time: 4h 52min 3s

## Linear-SVM on TFIDF

In [28]:

```
%%time
alpha = [10 ** x for x in range(-3, 2)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(class_weight='balanced',alpha=i, penalty='l1', loss='hinge',
random_state=42,n_jobs=-1)
    clf.fit(tr_X_TFIDF, y_train_tfidf)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(tr_X_TFIDF, y_train_tfidf)
    predict_y = sig_clf.predict_proba(tr_X_TFIDF)
    log_error_array.append(log_loss(y_train_tfidf, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_train_tfidf, predict_y,
labels=clf.classes_, eps=1e-15))

fig, ax = plt.subplots()
```

```

fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], 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()

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42, n_jobs=-1)
clf.fit(tr_X_TFIDF, y_train_tfidf)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(tr_X_TFIDF, y_train_tfidf)

predict_y = sig_clf.predict_proba(tr_X_TFIDF)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train_tfidf, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(te_X_TFIDF)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test_tfidf, predict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y, axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test_tfidf, predicted_y)

```

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning:
The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence
this warning.
    warnings.warn(CV_WARNING, FutureWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)

```

For values of alpha = 0.001 The log loss is: 0.6582699819717955

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning:
The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence
this warning.
    warnings.warn(CV_WARNING, FutureWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)

```

For values of alpha = 0.01 The log loss is: 0.6582699819717955

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning:
The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence
this warning.
    warnings.warn(CV_WARNING, FutureWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
```

For values of alpha = 0.1 The log loss is: 0.6582699819717955

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning:
The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence
this warning.
    warnings.warn(CV_WARNING, FutureWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
```

For values of alpha = 1 The log loss is: 0.5947455281018672

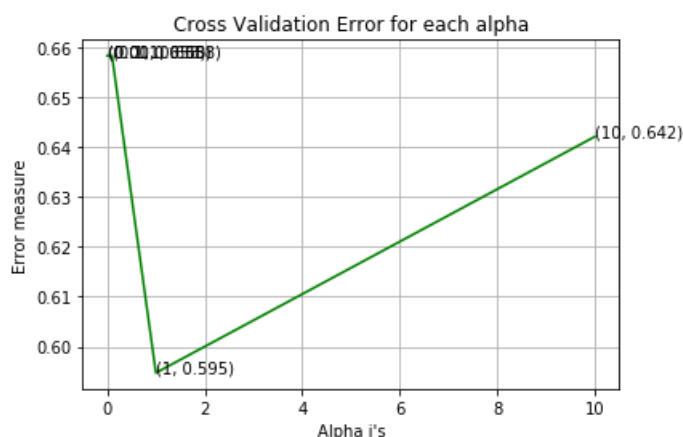
```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\stochastic_gradient.py:561:
ConvergenceWarning: Maximum number of iteration reached before convergence. Consider increasing ma
x_iter to improve the fit.
    ConvergenceWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning:
The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence
this warning.
    warnings.warn(CV_WARNING, FutureWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\stochastic_gradient.py:561:
ConvergenceWarning: Maximum number of iteration reached before convergence. Consider increasing ma
x_iter to improve the fit.
    ConvergenceWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\stochastic_gradient.py:561:
ConvergenceWarning: Maximum number of iteration reached before convergence. Consider increasing ma
x_iter to improve the fit.
    ConvergenceWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
```

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\stochastic_gradient.py:561:
ConvergenceWarning: Maximum number of iteration reached before convergence. Consider increasing ma
x_iter to improve the fit.
    ConvergenceWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)

```

For values of alpha = 10 The log loss is: 0.6420525099642131



```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning:
The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence
this warning.
    warnings.warn(CV_WARNING, FutureWarning)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:453: RuntimeWarning: overflow en
countered in exp
    E = np.exp(AB[0] * F + AB[1])
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\calibration.py:455: RuntimeWarning: invalid
value encountered in multiply
    TEP_minus_T1P = P * (T * E - T1)

```

For values of best alpha = 1 The train log loss is: 0.6122552309000592

For values of best alpha = 1 The test log loss is: 0.6160174504016971

Total number of data points : 121286

Confusion matrix

	Predicted NO	Predicted YES
Actual NO	75833	538
Actual YES	43946	969

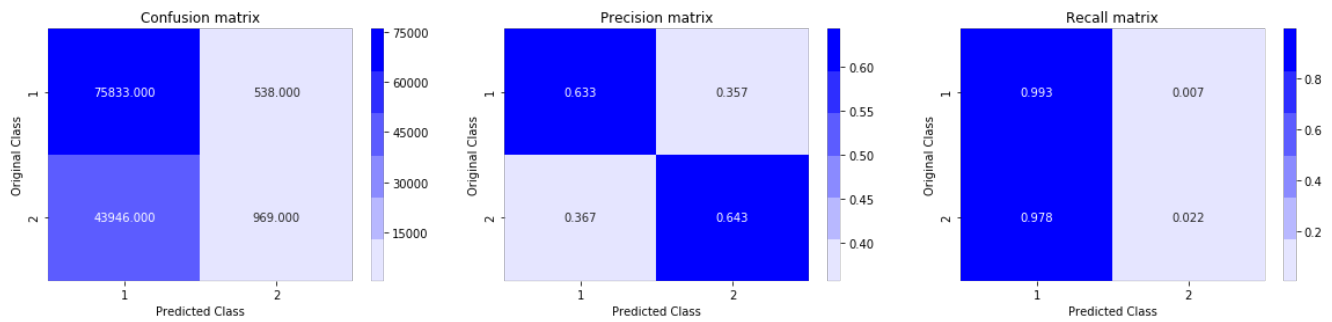
Precision matrix

	Predicted NO	Predicted YES
Actual NO	0.633108	0.357001
Actual YES	0.366892	0.642999

Recall matrix

	Predicted NO	Predicted YES
Actual NO	0.633108	0.357001
Actual YES	0.366892	0.642999





Wall time: 10h 13min 52s

## 2. Hyperparameter tune XgBoost using RandomSearch to reduce the log-loss on TFIDF Weighted W2V.

### 4.6 XGBoost Hyp

In [70]:

```
%%time
import xgboost as xgb
from xgboost.sklearn import XGBClassifier
params={"learning_rate" : [0.05, 0.10, 0.15, 0.20, 0.02] ,
        "max_depth" : [ 4, 10, 12, 15,30],
        "min_child_weight" : [ 1, 3, 5, 7 ],
        "gamma" : [ 0.0, 0.1, 0.2 , 0.3, 0.4 ],
        "colsample_bytree" : [ 0.5 , 0.7 ] }
model = XGBClassifier( n_estimators=20,objective= 'binary:logistic',seed=27,n_jobs=-1)
gs=GridSearchCV(estimator=model,cv=3,n_jobs=-1,scoring
='neg_log_loss',verbose=True,param_grid=params,return_train_score=True)
gs.fit(X_train[:20000],y_train[:20000])
```

Fitting 3 folds for each of 1000 candidates, totalling 3000 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 42 tasks | elapsed: 2.0min
[Parallel(n_jobs=-1)]: Done 192 tasks | elapsed: 9.4min
[Parallel(n_jobs=-1)]: Done 442 tasks | elapsed: 22.0min
[Parallel(n_jobs=-1)]: Done 792 tasks | elapsed: 39.7min
[Parallel(n_jobs=-1)]: Done 1242 tasks | elapsed: 62.4min
[Parallel(n_jobs=-1)]: Done 1792 tasks | elapsed: 94.8min
[Parallel(n_jobs=-1)]: Done 2442 tasks | elapsed: 137.4min
[Parallel(n_jobs=-1)]: Done 3000 out of 3000 | elapsed: 174.4min finished
```

Wall time: 2h 54min 45s

In [118]:

```
results=pd.DataFrame(gs.cv_results_).sort_values(by='rank_test_score').head(56)
results
```

Out[118]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_colsample_bytree	param_gamma	param_learn
966	13.915628	0.111048	0.078003	3.617543e-06	0.7	0.4	0.2
566	14.204998	0.069078	0.097879	3.131560e-03	0.7	0	0.2
333	14.222404	0.175731	0.078003	4.408441e-03	0.7	0.4	0.2

666	14.092431 mean_fit_time	0.175734 std_fit_time	0.078005 mean_score_time	1.189441e-06 std_score_time	0.7 param_colsample_bytree	0.1 param_gamma	0.2 param_learn
766	13.754334	0.074637	0.093603	2.206180e-02	0.7	0.2	0.2
866	13.915625	0.038212	0.083204	7.353333e-03	0.7	0.3	0.2
765	14.445953	0.066188	0.083201	7.353782e-03	0.7	0.2	0.2
667	13.744025	0.287096	0.083202	7.354850e-03	0.7	0.1	0.2
767	13.468327	0.106820	0.078000	1.949921e-06	0.7	0.2	0.2
267	10.358018	0.028967	0.083679	4.741240e-04	0.5	0.2	0.2
565	14.365220	0.181419	0.081538	2.716102e-03	0.7	0	0.2
266	10.555716	0.103093	0.083347	2.056499e-03	0.5	0.2	0.2
167	10.546839	0.121370	0.084013	8.190666e-04	0.5	0.1	0.2
967	13.520418	0.109819	0.083203	7.350467e-03	0.7	0.4	0.2
366	10.690335	0.122424	0.082681	4.701360e-04	0.5	0.3	0.2
367	10.405973	0.138627	0.084681	9.442586e-04	0.5	0.3	0.2
867	13.780426	0.187786	0.083200	7.351984e-03	0.7	0.3	0.2
864	15.214035	0.344709	0.088403	7.353389e-03	0.7	0.3	0.2
464	11.811035	0.244820	0.087679	4.733966e-04	0.5	0.4	0.2
66	10.598064	0.060019	0.084009	1.946680e-07	0.5	0	0.2
271	11.587333	0.022488	0.084347	9.434167e-04	0.5	0.2	0.2
65	11.009235	0.113846	0.084344	4.711457e-04	0.5	0	0.2

865	14.758811 mean_fit_time	0.435611 std_fit_time	0.093600 mean_score_time	1.273576e-02 std_score_time	0.7 param_colsample_bytree	0.3 param_gamma	0.2 param_learn
771	15.205172	0.070155	0.083203	7.354232e-03	0.7	0.2	0.2
64	11.737157	0.157437	0.086008	2.449507e-03	0.5	0	0.2
467	10.443355	0.098621	0.088680	1.246853e-03	0.5	0.4	0.2
567	13.912636	0.294046	0.093268	4.705836e-04	0.7	0	0.2
265	10.924898	0.071748	0.083678	4.724943e-04	0.5	0.2	0.2
665	14.539649	0.183708	0.088403	7.353614e-03	0.7	0.1	0.2
364	11.986251	0.315543	0.084680	9.443145e-04	0.5	0.3	0.2
965	14.117344	0.017645	0.081802	5.373779e-03	0.7	0.4	0.2
571	15.350354	0.027894	0.083748	5.465437e-03	0.7	0	0.2
770	15.719983	0.062834	0.088403	7.354344e-03	0.7	0.2	0.2
264	11.423771	0.052011	0.087013	2.943546e-03	0.5	0.2	0.2
564	15.016865	0.156883	0.084275	3.199846e-03	0.7	0	0.2
371	11.623468	0.096896	0.085679	2.055804e-03	0.5	0.3	0.2
365	11.291876	0.163076	0.084681	1.248573e-03	0.5	0.3	0.2
570	15.985007	0.240899	0.093407	6.775351e-03	0.7	0	0.2
166	10.719058	0.235413	0.085014	1.415347e-03	0.5	0.1	0.2
67	10.552896	0.051683	0.084675	4.713150e-04	0.5	0	0.2
670	15.907290	0.155657	0.078000	9.602742e-07	0.7	0.1	0.2

164	11.502922	0.257823	0.085013	2.060172e-06	0.5	0.1	0.2
	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_colsample_bytree	param_gamma	param_learn
971	15.314472	0.279747	0.093604	2.922181e-06	0.7	0.4	0.2
964	14.973096	0.121324	0.083204	7.352209e-03	0.7	0.4	0.2
466	10.789619	0.078898	0.086146	2.612890e-03	0.5	0.4	0.2
471	12.119001	0.050819	0.086682	2.057842e-03	0.5	0.4	0.2
375	12.827378	0.060354	0.088348	4.713710e-04	0.5	0.3	0.2
664	14.981657	0.181187	0.083201	7.357154e-03	0.7	0.1	0.2
871	15.614089	0.314983	0.088602	7.500877e-03	0.7	0.3	0.2
465	10.991251	0.126774	0.085012	8.167310e-04	0.5	0.4	0.2
547	13.551620	0.191632	0.093603	5.285991e-06	0.7	0	0.15
569	16.909619	0.100443	0.098483	4.376062e-03	0.7	0	0.2
945	14.399237	0.099483	0.083204	7.353164e-03	0.7	0.4	0.15
764	14.877561	0.171998	0.083200	7.354681e-03	0.7	0.2	0.2
647	13.707623	0.246446	0.083202	7.351029e-03	0.7	0.1	0.15
544	14.992062	0.166565	0.083199	7.354625e-03	0.7	0	0.15
747	13.695946	0.137778	0.083203	7.354794e-03	0.7	0.2	0.15

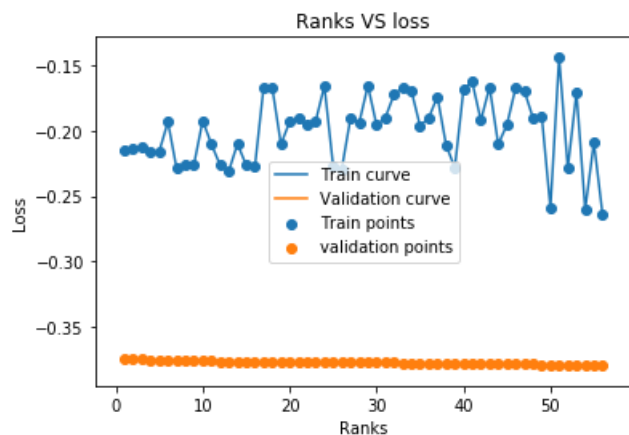
56 rows × 21 columns



In [123]:

```
x=[x for x in range(1,57)]
y_train_loss=results.mean_train_score.tolist()
y_test_loss=results.mean_test_score.tolist()
#function to plot lines
#plt.xscale('log')
plt.plot(x,y_train_loss,label="Train curve")
plt.plot(x,y_test_loss,label="Validation curve")
plt.scatter(x,y_train_loss,label='Train points')
```

```
plt.scatter(x, y_train_loss, label='train points')
plt.scatter(x, y_test_loss, label='validation points')
plt.xlabel("Ranks")
plt.ylabel("Loss")
plt.title("Ranks VS loss")
plt.legend()
plt.show()
```



## Observation

1. Model is overfitting once I choose the best parameters {'colsample\_bytree': 0.7, 'gamma': 0.2, 'learning\_rate': 0.2, 'max\_depth': 10, 'min\_child\_weight': 5} train-logloss:0.002412 valid-logloss:0.374075 lot of difference between train and validation log loss
2. From the above graph we can observe that at 56th ranked point there is no difference between train and validation loss hence i have selected the paramaters of the 56th ranked datapoint in results dataframe.{'colsample\_bytree': 0.7, 'gamma': 0.2, 'learning\_rate': 0.15, 'max\_depth': 10, 'min\_child\_weight': 7} .
3. Index of 56th ranked datapoint is 747.

In [120]:

```
results.params[747]
```

Out[120]:

```
{'colsample_bytree': 0.7,
'gamma': 0.2,
'learning_rate': 0.15,
'max_depth': 10,
'min_child_weight': 7}
```

## 4.6 XGBoost

In [112]:

```
%%time
params = {}
params['objective'] = 'binary:logistic'
params['eval_metric'] = 'logloss'
params['eta'] = 0.02
params['colsample_bytree']= 0.7
params["gamma"]=0.2
params['max_depth']=10
params['min_child_weight']=7

d_train = xgb.DMatrix(X_train, label=y_train)
d_test = xgb.DMatrix(X_test, label=y_test)

watchlist = [(d_train, 'train'), (d_test, 'valid')]

bst1 = xgb.train(params, d_train,
                 , watchlist, early_stopping_rounds=20, verbose_eval=10)
model = bst1.predict(DMatrix(X_test, label=y_test))
```

```
xgmat1 = xgb.DMatrix(x_train,y_train)
predict_y1 = bst1.predict(d_test)
print("The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
[0] train-logloss:0.684462 valid-logloss:0.685346
Multiple eval metrics have been passed: 'valid-logloss' will be used for early stopping.
```

Will train until valid-logloss hasn't improved in 20 rounds.

```
[10] train-logloss:0.599485 valid-logloss:0.613071
[20] train-logloss:0.537971 valid-logloss:0.561511
[30] train-logloss:0.491859 valid-logloss:0.523821
[40] train-logloss:0.454961 valid-logloss:0.494755
[50] train-logloss:0.426053 valid-logloss:0.472528
[60] train-logloss:0.402424 valid-logloss:0.455005
[70] train-logloss:0.383142 valid-logloss:0.441331
[80] train-logloss:0.367218 valid-logloss:0.430515
[90] train-logloss:0.353932 valid-logloss:0.421765
[100] train-logloss:0.342581 valid-logloss:0.414467
[110] train-logloss:0.332798 valid-logloss:0.408675
[120] train-logloss:0.324172 valid-logloss:0.403869
[130] train-logloss:0.316635 valid-logloss:0.399844
[140] train-logloss:0.309682 valid-logloss:0.396455
[150] train-logloss:0.303194 valid-logloss:0.393674
[160] train-logloss:0.297271 valid-logloss:0.391282
[170] train-logloss:0.292126 valid-logloss:0.389246
[180] train-logloss:0.287445 valid-logloss:0.387482
[190] train-logloss:0.282747 valid-logloss:0.386043
[200] train-logloss:0.278738 valid-logloss:0.384653
[210] train-logloss:0.27469 valid-logloss:0.383406
[220] train-logloss:0.271455 valid-logloss:0.382389
[230] train-logloss:0.268307 valid-logloss:0.381438
[240] train-logloss:0.265483 valid-logloss:0.380701
[249] train-logloss:0.262691 valid-logloss:0.379996
The test log loss is: 0.3797524013703534
Wall time: 19min 13s
```

In [113]:

```
predicted_y1 =np.array(predict_y1>0.5,dtype=int)
print("Total number of data points :", len(predicted_y1))
plot_confusion_matrix(y_test, predicted_y1)
```

Total number of data points : 30000

Confusion matrix

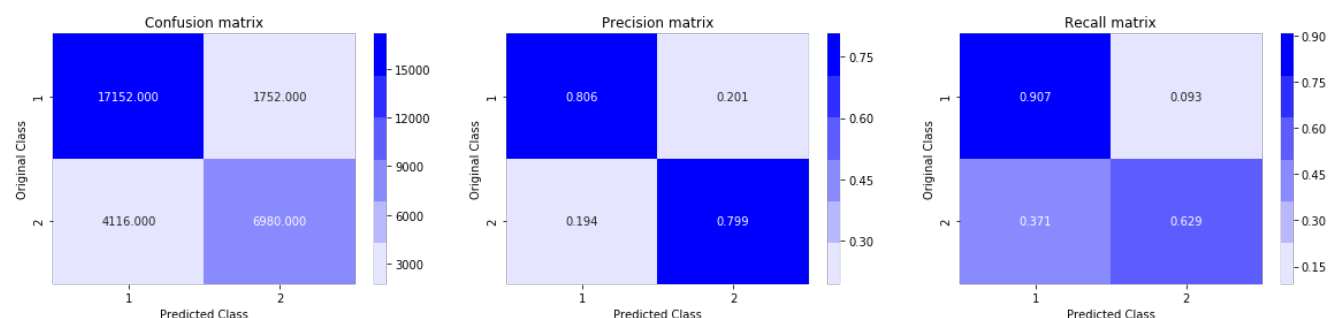
	Predicted NO	Predicted YES
Actual NO	17152	1752
Actual YES	4116	6980

Precision matrix

	Predicted NO	Predicted YES
Actual NO	0.80647	0.200641
Actual YES	0.19353	0.799359

Recall matrix

	Predicted NO	Predicted YES
Actual NO	0.80647	0.200641
Actual YES	0.19353	0.799359



In [146]:

```
from prettytable import PrettyTable

table = PrettyTable()
```

```

table.field_names = ["Vectorizer", "Test loss"]

table.add_row(["Random model","0.89" ])
table.add_row(["logistic regression ","0.59"])
table.add_row(["linear svm","0.59" ])
table.add_row(["logistic regression on TFIDF","0.62"])
table.add_row(["linear svm on TFIDF","0.61" ])
table.add_row(["XG bosst on TFIDF W2V", "0.37" ])

print(table)

```

Vectorizer	Test loss
Random model	0.89
logistic regression	0.59
linear svm	0.59
logistic regression on TFIDF	0.62
linear svm on TFIDF	0.61
XG bosst on TFIDF W2V	0.37