

## Import required libraries

In [112]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import re
4 import time
5 import warnings
6 import sqlite3
7 from sqlalchemy import create_engine # database connection
8 import csv
9 import os
10 warnings.filterwarnings("ignore")
11 import datetime as dt
12 import numpy as np
13 from sklearn.model_selection import StratifiedKFold
14 from sklearn.feature_extraction.text import CountVectorizer
15 import seaborn as sns
16 from prettytable import PrettyTable
17 from sklearn.metrics import confusion_matrix
18 from sklearn.metrics.classification import accuracy_score, log_loss
19 from sklearn.preprocessing import StandardScaler
20 from sklearn.model_selection import RandomizedSearchCV
21 from sklearn.metrics import *
22 from sklearn.feature_extraction.text import TfidfVectorizer
23 from scipy.sparse import hstack
24 from sklearn.calibration import CalibratedClassifierCV
25 from sklearn.model_selection import train_test_split
26 from sklearn.model_selection import GridSearchCV
27 import math
28 import xgboost as xgb
29 #IMPORTING XGBOOST WRAPPER OF SKLEARN
30 from xgboost.sklearn import XGBClassifier
31 import scipy
32
33 from sklearn.model_selection import cross_val_score
34 from sklearn.linear_model import SGDClassifier
35 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
36 from sklearn import model_selection
```

## Save Model For Productionization:

In [113]:

```

1 from sklearn.externals import joblib
2 def saveModelToFile(obj,filename):
3     joblib.dump(obj,open(filename+".pkl","wb"))
4 def openModelFromFile(filename):
5     temp = joblib.load(open(filename+".pkl","rb"))
6     return temp

```

## Function for converting features to numerics :

In [114]:

```

1 # after we read from sql table each entry was read it as a string
2 # we convert all the features into numeric before we apply any model
3 def to_numeric(data, isDF=True):
4     if isDF:
5         cols = list(data.columns)
6         for i in cols:
7             #print('column converted:', end=' ')
8             data[i] = data[i].apply(pd.to_numeric)
9             # print(i)
10    else:
11        # https://stackoverflow.com/questions/7368789/convert-all-strings-in-a-list-to-int
12        data = list(map(int, data.values))
13    return data

```

## Function for stacking matrix and array:

In [115]:

```

1 def hstack_sparse_with_dense(sparse_matrix, array):
2     stacked_mat=scipy.sparse.hstack((sparse_matrix, array))
3     return stacked_mat.tocsr()

```

## Function for Random Model:

```
In [116]: 1 # https://stackoverflow.com/questions/36089083/how-to-use-Log-Loss-in-skLearn-grid-search
2
3 # we need to generate 9 numbers and the sum of numbers should be 1
4 # one solution is to generate 9 numbers and divide each of the numbers by their sum
5 # ref: https://stackoverflow.com/a/18662466/4084039
6 # we create a output array that has exactly same size as the CV data
7 def random_model(y_test,model_summary):
8     predicted_y = np.zeros((len(y_test),2))
9     for i in range(len(y_test)):
10         rand_probs = np.random.rand(1,2)
11         predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
12         loss=log_loss(y_test, predicted_y, eps=1e-15)
13         print("Log loss on Test Data using Random Model",loss)
14
15     predicted_y =np.argmax(predicted_y, axis=1)
16     plot_confusion_matrix(y_test, predicted_y)
17     model_summary.add_row(['Random', "%4f"%loss, "%4f"%loss])
18     return predicted_y
```

**Function for finding optimal value of hyperparameter and plot missclassification error vs hyperparam :**

## Logistic-regression

In [117]:

```
1 ##METHOD-1(CV USING SEARCHING ALGO)
2 def linear_model_lr(x_train, y_train, x_test, y_test, params, CV, search='random'):
3     #INITIALIZE CLASSIFIER
4     clf=SGDClassifier(loss='log', class_weight='balanced', penalty='l2', random_state=32)
5     if search=='random':
6         model= RandomizedSearchCV(clf,\n7             return_train_score=True,\n8             param_distributions=params,\n9             n_iter=len(params['alpha']),\n10            n_jobs=-1,\n11            cv=CV,\n12            scoring='neg_log_loss')
13     model.fit(x_train, y_train)
14
15     train_score= abs(np.array(model.cv_results_['mean_train_score']))
16     cv_score= abs(np.array(model.cv_results_['mean_test_score']))
17
18     for i in range(len(params['alpha'])):
19         print('for alpha: %f,\t\t train log-loss: %.4f,\t cv log-loss: %.4f'%( params['alpha'][i],train_score[i],cv_
20     print('\n\noptimal alpha:%f, log-loss: %.4f'%( model.best_params_['alpha'], abs(model.best_score_)))
21     plt.figure(1,figsize=(7,5))
22     sns.set_style('darkgrid')
23     plt.plot(params['alpha'], train_score, label='Train(Log-Loss)')
24     plt.plot(params['alpha'], cv_score, label='CrossValidation(Log-Loss)')
25     plt.title('Loss-Plot')
26     plt.legend()
27     plt.xlabel('Alpha: Hyperparam')
28     plt.ylabel('Log-Loss')
29     plt.show()
30     return model, abs(min(train_score))#model.best_score_)
```

**Function for finding optimal value of hyperparameter and plot missclassification error vs hyperparam :**

## Linear-SVM

In [118]:

```
1 def linear_model_svm(xtrain, ytrain, xtest, ytest, params):
2     #DISJOINT DATASET
3     loss_tr=[]
4     loss_cv=[]
5     #SPLIT DATA INTO TRAIN AND CV SET
6     xtr, xcv, ytr, ycv= train_test_split(xtrain, ytrain, test_size=.3, shuffle=True)
7     for alpha in params['alpha']:
8
9         clf=SGDClassifier(alpha=alpha, loss='hinge', penalty='l2', random_state=42, class_weight='balanced' , n_jobs=-1)
10
11        #DISJOINT DATASET
12        #xtr, xcv, ytr, ycv= train_test_split(xtrain, ytrain, test_size=.3, shuffle=True)
13
14        #FIT FOR TRAIN DATA
15        clf.fit(xtr, ytr)
16
17        #CALIB MODEL
18        calib_clf=CalibratedClassifierCV(clf, method='sigmoid', cv=3)
19        calib_clf.fit(xtr, ytr)
20
21        #PRED PROBABILITY SCORE AND LOG-LOSS USING CALIB MODEL
22        y_proba_tr= calib_clf.predict_proba(xtr)[:,1]
23        loss_tr.append(log_loss(ytr, y_proba_tr, eps=1e-15))
24
25        y_proba_cv=calib_clf.predict_proba(xcv)[:,1]
26        loss_cv.append(log_loss(ycv, y_proba_cv, eps=1e-15))
27
28    for i in range(len(params['alpha'])):
29        print('for alpha: %f,\t\t train log-loss: %.4f,\t\t cv log-loss: %.4f'%( params['alpha'][i],loss_tr[i],loss_cv[i]))
30
31    #PRED LABELS FOR TEST DATA USING SGDC MODEL
32    optimal_param = params['alpha'][np.argmin(loss_cv)]
33    plt.figure(1,figsize=(7,5))
34    sns.set_style('darkgrid')
35    plt.plot(params['alpha'], loss_tr, label='Train Log-Loss')
36    plt.plot(params['alpha'], loss_cv, label='CV Log-Loss')
37    plt.xlabel('Alpha: Hyperparam')
38    plt.ylabel('Log-Loss')
39    plt.title('Loss-Plot')
40    plt.legend()
41    plt.show()
```

```
42     return optimal_param, min(loss_tr)
43
```

**Function for finding optimal value of hyperparameter :**

**XGBOOST (GBDT)**

In [119]:

```
1 def Ensemble_Classifier(x_train,y_train,CV,params_, tune_param, searchMethod, param_name):
2     #INITIALIZE GBDT CLASSIFIER
3     clf=xgb.XGBClassifier(n_estimators=params_[ 'n_estimators'],\
4                           max_depth=params_[ 'max_depth'],\
5                           eta=.02,\
6                           reg_alpha=params_[ 'reg_alpha'],\
7                           min_child_weight=params_[ 'min_child_weight'],\
8                           gamma=params_[ 'gamma'],\
9                           subsample=params_[ 'subsample'],\
10                          colsample_bytree=params_[ 'colsample_bytree'],\
11                          booster='gbtree',\
12                          objective='binary:logistic')
13     # APPLY RANDOM OR GRID SEARCH FOR HYPERPARAMETER TUNNING
14     if searchMethod=='random':
15         ##len(params[ 'max_depth']),\
16         model=RandomizedSearchCV(clf,\n17                               verbose=1,\n18                               n_jobs=-1,\n19                               cv=CV,\n20                               param_distributions=tune_param,\n21                               n_iter=3,\n22                               return_train_score=True,\n23                               scoring='neg_log_loss')#make_scorer(roc_auc_score,average='weighted'))
24     elif searchMethod=='grid':
25         model= GridSearchCV(estimator=clf,\n26                               verbose=1,\n27                               param_grid=tune_param,\n28                               scoring='neg_log_loss',\n29                               n_jobs=-1,\n30                               cv=CV,\n31                               return_train_score=True\n32 )
33     model.fit(x_train,y_train)
34     #PLOT THE PERFORMANCE OF MODEL ON CROSSVALIDATION DATA FOR EACH HYPERPARAM VALUE
35     train_loss = abs(model.cv_results_[ 'mean_train_score'])
36     cv_loss = abs(model.cv_results_[ 'mean_test_score'])
37     #if len(param_name.split())==1:
38     for i in range(len(train_loss)):
39         print('Train Log-Loss: %.4f'%train_loss[i], 'CV Log-Loss: %.4f'%cv_loss[i])
40     '''plt.figure(1, figsize=(7,5))
41     sns.set_style('darkgrid')
```

```

42     plt.title('Loss-Plot')
43     plt.xlabel('Hyperparam')
44     plt.ylabel('Log-Loss')
45     plt.plot(params_[param_name], train_loss, label='Train')
46     plt.plot(params_[param_name], cv_loss, label='Train')
47     plt.legend()
48     plt.show()'''
49 #elif len(param_name.split())==2:
50 #    hyperparam_matrix(train_loss, cv_loss, param_list=param_name.split())
51
52 return model

```

## Function for measuring performance on test data (linear models):

In [120]:

```

1 def test_performance_linear(param, x_train, y_train, x_test, y_test, model_summary, train_loss, model_use, loss='log'
2     #SGDC CLASSIFIER
3     clf=SGDClassifier(alpha=param,\n4             loss=loss,\n5             penalty='l2',\n6             random_state=42,\n7             n_jobs=-1,\n8             class_weight='balanced')\n9     clf.fit(x_train, y_train)\n10\n11     #CALIB CLASSIFIER\n12     cal_clf=CalibratedClassifierCV(clf, method='sigmoid', cv=3)\n13     cal_clf.fit(x_train, y_train)\n14\n15     #PREDICT CLASS LABEL USING SGDC MODEL\n16     y_pred=clf.predict(x_test)\n17     #PREDICT PROBA USING CALIB MODEL\n18     y_prob=cal_clf.predict_proba(x_test)[:,1]\n19     #PLOT CONFUSION MATRIX\n20     plot_confusion_matrix(y_test, y_pred)\n21     #TEST LOSS\n22     loss=log_loss(y_test, y_prob, eps=1e-15)\n23     print('test log-loss: %.4f'%loss)\n24     model_summary.add_row([model_use, '%.4f'%train_loss, '%.4f'%loss])

```

## Function for measuring performance on test data (XGBOOST model):

```
In [121]: 1 def test_performance(x_train,y_train,x_test,y_test,params_,model_use=None,summary=False):
2     '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL HYPERPARAM'''
3     #INITIALIZE GBDT WITH OPTIMAL VALUE OF HYPERPARAMS
4     clf=xgb.XGBClassifier(n_estimators=params_[ 'n_estimators'],\
5                           max_depth=params_[ 'max_depth'],\
6                           eta=.02,\
7                           reg_alpha=params_[ 'reg_alpha'],\
8                           min_child_weight=params_[ 'min_child_weight'],\
9                           gamma=params_[ 'gamma'],\
10                          subsample=params_[ 'subsample'],\
11                          colsample_bytree=params_[ 'colsample_bytree'],\
12                          booster='gbtree',\
13                          objective='binary:logistic',\
14                          verbose=1)
15
16     clf.fit(x_train, y_train)
17     y_pred=clf.predict(x_test)
18     y_proba=clf.predict_proba(x_test)[:,1]
19     test_loss=log_loss(y_test, y_proba, labels=np.array([0,1]), eps=1e-15)
20     print('test_loss:',test_loss)
21     plot_confusion_matrix(y_test, y_pred)
22     if summary:
23         y_proba_tr=clf.predict_proba(x_train)[:,1]
24         train_loss=log_loss(y_train, y_proba_tr, labels=np.array([0,1]), eps=1e-15)
25         model_summary.add_row([model_use, '%.4f'%train_loss, '%.4f'%test_loss])
```

## Function for filling Nan and None values:

```
In [122]: 1 def fill_none(data):
2     data = data.fillna('')
3     nan_rows = data[data.isnull().any(1)]
4     print(nan_rows)
5     return data
```

## Function for Standardizing data:

In [123]:

```

1 def std_data(train,test,mean):
2     scaler=StandardScaler(with_mean=mean)
3     std_train=scaler.fit_transform(train)
4     std_test=scaler.transform(test)
5     return std_train, std_test

```

## Function for plotting confusion matrix:

In [124]:

```

1 # This function plots the confusion matrices given y_i, y_i_hat.
2 def plot_confusion_matrix(test_y, predict_y):
3     C = confusion_matrix(test_y, predict_y)
4     # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j
5     #precision
6     A =(((C.T)/(C.sum(axis=1))).T)
7     #recall
8     B =(C/C.sum(axis=0))
9
10    plt.figure(figsize=(16,4))
11    labels = [1,2]
12    i=131
13    # representing A in heatmap format
14    cmap=sns.light_palette('blue')
15    for mat in [C,A,B]:
16        plt.subplot(i)
17        sns.heatmap(mat, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
18        plt.xlabel('Predicted Class')
19        plt.ylabel('Original Class')
20        plt.title("Confusion matrix")
21        i+=1
22    plt.show()

```

## 4. Machine Learning Models

In [125]:

```

1 model_summary=PrettyTable()
2 model_summary.field_names = [ 'Model', 'Train(Log-Loss)', 'Test(Log-Loss)' ]

```

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

In [126]:

```
1 #Creating db file from csv
2 if not os.path.isfile('train.db'):
3     disk_engine = create_engine('sqlite:///train.db')
4     start = dt.datetime.now()
5     chunksize = 180000
6     j = 0
7     index_start = 1
8     col_ylabel=[str(i) + '_y' for i in range(0, 96, 1)]
9     col_xlabel=[str(i) + '_x' for i in range(0, 96, 1)]
10    col_names=[ 'Unnamed: 0','id','is_duplicate','cwc_min','cwc_max','csc_min','csc_max','ctc_min',\
11        'ctc_max','last_word_eq','first_word_eq','abs_len_diff','mean_len','token_set_ratio',\
12        'token_sort_ratio','fuzz_ratio','fuzz_partial_ratio','longest_substr_ratio','freq_qid1',\
13        'freq_qid2','q1len','q2len','q1_n_words','q2_n_words','word_Common','word_Total','word_share',\
14        'freq_q1+q2','freq_q1-q2']
15
16    l=[]
17    l.extend(col_names)
18    l.extend(col_xlabel)
19    l.extend(col_ylabel)
20    for df in pd.read_csv('final_features.csv', names=l, chunksize=chunksize, iterator=True, encoding='utf-8'):
21        df.index += index_start
22        j+=1
23        print('{} rows'.format(j*chunksize))
24        df.to_sql('data', disk_engine, if_exists='append')
25        index_start = df.index[-1] + 1
26    j = 0
27    col_names=[ 'Unnamed: 0','id','question1','question2','is_duplicate','cwc_min','cwc_max','csc_min','csc_max','ctc_min',\
28        'ctc_max','last_word_eq','first_word_eq','abs_len_diff','mean_len','token_set_ratio',\
29        'token_sort_ratio','fuzz_ratio','fuzz_partial_ratio','longest_substr_ratio','freq_qid1',\
30        'freq_qid2','q1len','q2len','q1_n_words','q2_n_words','word_Common','word_Total','word_share',\
31        'freq_q1+q2','freq_q1-q2']
32    index_start = 1
33    for df in pd.read_csv('final_features_tf.csv', names=col_names, chunksize=chunksize, iterator=True, encoding='ut
34        df.index += index_start
35        j+=1
36        print('{} rows'.format(j*chunksize))
37        df.to_sql('data_tf', disk_engine, if_exists='append')
38        index_start = df.index[-1] + 1
```

## 4.2 Function for connecting to database and checking for existance of table:

```
In [127]: 1 #http://www.sqlitetutorial.net/sqlite-python/create-tables/
2 def create_connection(db_file):
3     """ create a database connection to the SQLite database specified by db_file
4         :db_file : database file
5         :return : Connection object or None
6     """
7     try:
8         conn = sqlite3.connect(db_file)
9         return conn
10    except Error as e:
11        print(e)
12
13    return None
14
15
16 def checkTableExists(dbcon):
17     cursr = dbcon.cursor()
18     str = "select name from sqlite_master where type='table'"
19     table_names = cursr.execute(str)
20     print("Tables in the databse:")
21     tables =table_names.fetchall()
22     print(tables[0][0])
23     return(len(tables))
```

```
In [128]: 1 read_db = 'train.db'
2 conn_r = create_connection(read_db)
3 checkTableExists(conn_r)
4 conn_r.close()
```

Tables in the databse:  
data

## 4.3 Sample Data:

In [129]:

```

1 %%time
2 # try to sample data according to the computing power you have
3 if os.path.isfile(read_db):
4     conn_r = create_connection(read_db)
5     if conn_r is not None:
6         # for selecting first 1M rows
7         # data = pd.read_sql_query("""SELECT * FROM data LIMIT 100001;""", conn_r)
8
9         # for selecting random points
10        data_tf = pd.read_sql_query("SELECT * From data_tf ORDER BY RANDOM() LIMIT 200001;", conn_r)
11        data=pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 80001;", conn_r)
12        conn_r.commit()
13        conn_r.close()

```

CPU times: user 39.3 s, sys: 12.4 s, total: 51.7 s

Wall time: 51.7 s

In [130]:

```

1 # remove the first row
2 data.drop(data.index[0], inplace=True)
3 data_tf.drop(data_tf.index[0], inplace=True)
4 y_true = data['is_duplicate']
5 y_true_tf = data_tf['is_duplicate']
6 data.drop(['Unnamed: 0', 'id','index','is_duplicate'], axis=1, inplace=True)
7 data_tf.drop(['Unnamed: 0', 'id','index','is_duplicate'], axis=1, inplace=True)

```

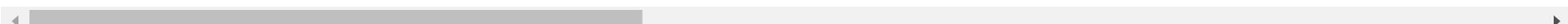
In [131]:

1 data.head(2)

Out[131]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq
1	0.66664444518516	0.399992000159997	0.749981250468738	0.599988000239995	0.714275510349852	0.49999500005	0.0	1.0
2	0.33332222259258	0.199996000079998	0.499991666805553	0.428565306209911	0.39999600004	0.363633057881292	0.0	1.0

2 rows × 218 columns



In [132]: 1 data\_tf.head(2)

Out[132]:

	question1	question2	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last
1	what does donald trump own victory mean for si...	how would a trump presidency affect silicon va...	0.499991666805553	0.499991666805553	0.0	0.0	0.374995312558593	0.29999700003	
2	can you give a ticket to a cop for speeding or...	can a cop give you a ticket for running a stop...	0.428565306209911	0.374995312558593	0.799984000319994	0.444439506227709	0.43749726564209	0.411762283751272	

2 rows × 28 columns



#### 4.4 Fill none and NaN:

In [133]:

```
1 %%time
2 data_tf = fill_none(data_tf)
3 data = fill_none(data)
```

Empty DataFrame

Columns: [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, freq\_qid1, freq\_qid2, q1len, q2len, q1\_n\_words, q2\_n\_words, word\_Common, word\_Total, word\_share, freq\_q1+q2, freq\_q1-q2]  
Index: []

[0 rows x 28 columns]

Empty DataFrame

Columns: [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, ...]

Index: []

[0 rows x 218 columns]

CPU times: user 14 s, sys: 532 ms, total: 14.5 s

Wall time: 14.5 s

## 4.5 convert to numerics

In [134]:

```

1  %%time
2  data = to_numeric(data)
3  y_true = to_numeric(y_true, isDF=False)
4
5  qq_data = data_tf[['question1', 'question2']]
6  cols=list(data_tf.columns)
7  cols.remove('question1')
8  cols.remove('question2')
9  numeric_data = to_numeric(data_tf[cols])
10 data_tf = pd.concat([qq_data, numeric_data], axis=1)
11 y_true_tf = to_numeric(y_true_tf, isDF=False)

```

CPU times: user 11min 1s, sys: 15.5 s, total: 11min 16s  
Wall time: 11min 16s

## 4.6 Split data to train and test set:

In [135]:

```

1 X_train, X_test, y_train, y_test = train_test_split(data, y_true, test_size=0.3)
2 X_train_tf, X_test_tf, y_train_tf, y_test_tf = train_test_split(data_tf, y_true_tf, test_size=0.3 )

```

In [136]:

```

1 x_train_qq = X_train_tf['question1']+ ' '+X_train_tf['question2']
2 x_test_qq = X_test_tf['question1']+ ' '+X_test_tf['question2']

```

## 4.7 Function for featurizing data :

In [137]:

```

1 def tfidfVector(X_train,X_test,max_features=None):
2     tf_idf_vect = TfidfVectorizer(ngram_range=(1,1),min_df=10,max_features=max_features)
3     X_train_tfidf = tf_idf_vect.fit_transform(X_train)
4     print("the type of count vectorizer: ",type(X_train_tfidf))
5     print("the shape of out text TFIDF vectorizer: ",X_train_tfidf.get_shape())
6     print("the number of unique words including both unigrams and bigrams: ", X_train_tfidf.get_shape()[1])
7
8     #processing of test data(convert test data into numerical vectors)
9     X_test_tfidf = tf_idf_vect.transform(X_test)
10    print("the shape of out text BOW vectorizer: ",X_test_tfidf.get_shape())
11    return tf_idf_vect, X_train_tfidf, X_test_tfidf

```

In [138]:

```

1 # Tfifd vector with all features which we use for brute force implementation
2 %time tf_idf_vect_q1, X_train_tfidf_q1, X_test_tfidf_q1 = tfidfVector(x_train_qq.values,\n
3                                         x_test_qq.values,\n
4                                         max_features=None)
5 print()
6 #%%time tf_idf_vect_q2, X_train_tfidf_q2, X_test_tfidf_q2 = tfidfVector(x_train_tf_q2.values,\n
7 #                                         x_test_tf_q2.values,\n
8 #                                         max_features=None)

```

the type of count vectorizer: <class 'scipy.sparse.csr.csr\_matrix'>  
 the shape of out text TFIDF vectorizer: (140000, 10576)  
 the number of unique words including both unigrams and bigrams: 10576  
 the shape of out text BOW vectorizer: (60000, 10576)  
 CPU times: user 4.54 s, sys: 4 ms, total: 4.54 s  
 Wall time: 4.55 s

## 4.8 Stack TFIDF matrix with other features

In [139]:

```

1 stacked_tf_data_train=hstack_sparse_with_dense(X_train_tfidf_q1, np.array(X_train_tf.drop(['question1','question2']),\n2 print(stacked_tf_data_train.shape)

```

(140000, 10602)

In [140]:

```

1 stacked_tf_data_test=hstack_sparse_with_dense(X_test_tfidf_q1, np.array(X_test_tf.drop(['question1','question2']),axis=0)\n2 print(stacked_tf_data_test.shape)

```

(60000, 10602)

## Model:

### Initialize common objects :

In [141]:

```

1 #HYPERPARAM RANGE
2 params={'alpha':[10**x for x in range(-5,5)]}
3 #STRATIFIED CV SPLIT
4 CV = StratifiedKFold(n_splits=10)
5 #TYPES OF MODELS
6 model_use=['Logistic-Regression', 'Linear-SVM', 'XGBoost']

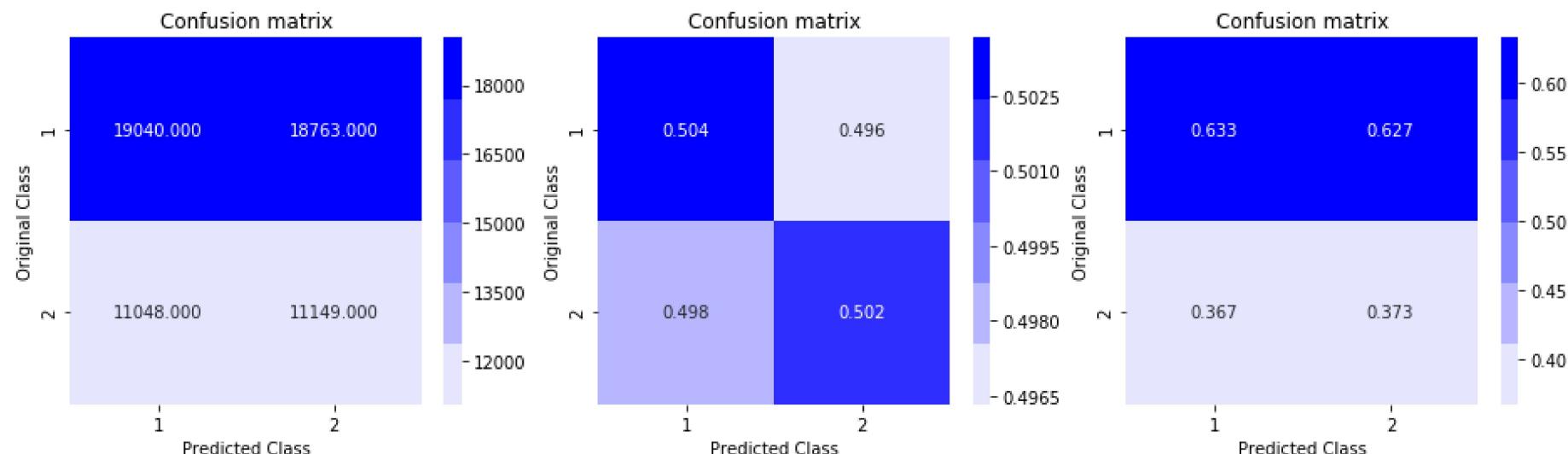
```

## 1.Random Model:

In [142]:

```
1 predicted_y=random_model(y_test_tf,model_summary)
```

Log loss on Test Data using Random Model 0.8847709468609612



In [143]:

```
1 print(model_summary)
```

Model	Train(Log-Loss)	Test(Log-Loss)
Random	0.8848	0.8848

## Logistic Regression:

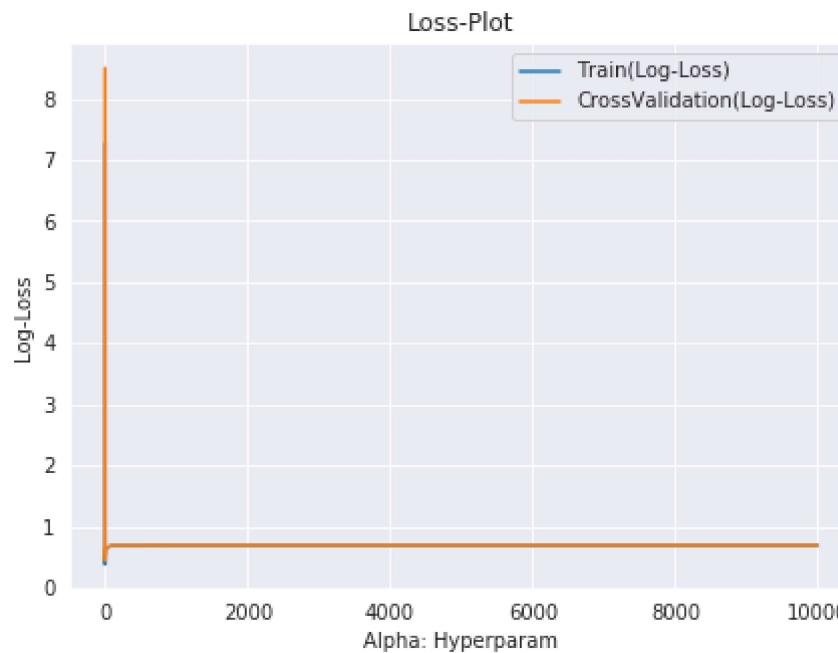
1.Hyperparam tunning and loss plot :

In [144]:

```
1 #STANDARDIZING DATA
2 train, test = std_data(train=stacked_tf_data_train,test=stacked_tf_data_test,mean=False)
3 #TUNNING HYPERPARAM
4 model, opt_train_loss = linear_model_lr(train, y_train_tf, test, y_test_tf, params, CV, search='random')
```

```
for alpha: 0.000010,          train log-loss: 7.2698,      cv log-loss: 8.4978
for alpha: 0.000100,          train log-loss: 6.1089,      cv log-loss: 7.3266
for alpha: 0.001000,          train log-loss: 1.6636,      cv log-loss: 2.2396
for alpha: 0.010000,          train log-loss: 0.3822,      cv log-loss: 0.4763
for alpha: 0.100000,          train log-loss: 0.4101,      cv log-loss: 0.4508
for alpha: 1.000000,          train log-loss: 0.5310,      cv log-loss: 0.5421
for alpha: 10.000000,         train log-loss: 0.6556,      cv log-loss: 0.6574
for alpha: 100.000000,        train log-loss: 0.6905,      cv log-loss: 0.6907
for alpha: 1000.000000,       train log-loss: 0.6929,      cv log-loss: 0.6929
for alpha: 10000.000000,      train log-loss: 0.6931,      cv log-loss: 0.6931
```

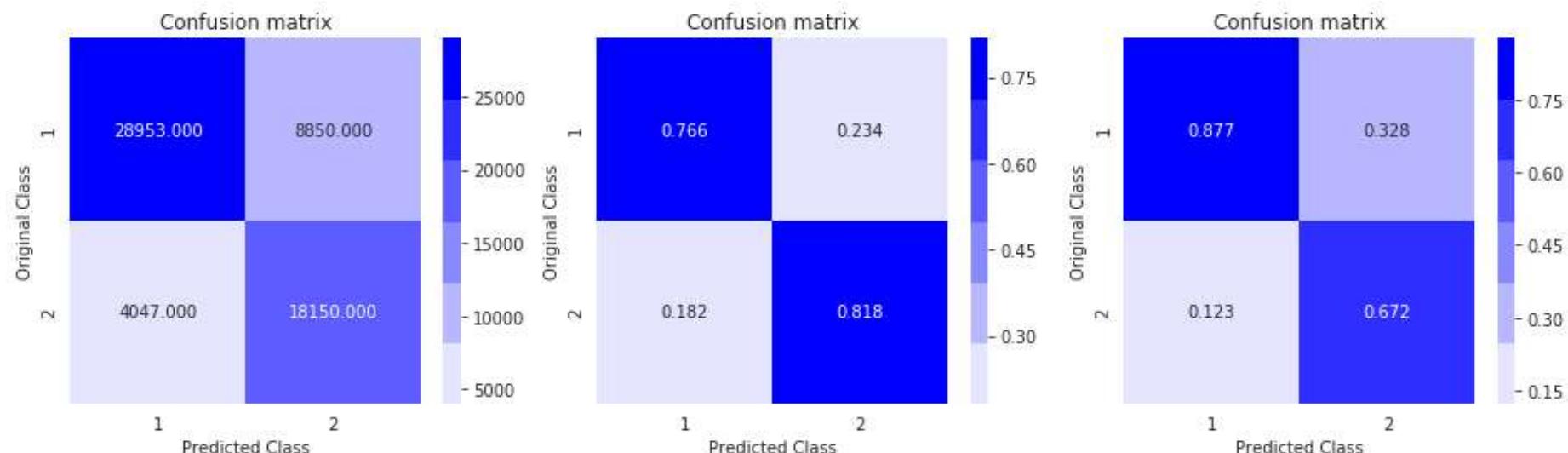
optimal alpha:0.100000, log-loss: 0.4508



## 2.Performance on test data:

In [145]:

```
1 opt_param=model.best_params_[ 'alpha' ]
2 test_performance_linear(opt_param, train, y_train_tf, test, y_test_tf, model_summary, opt_train_loss, model_use[0],)
```



test log-loss: 0.4176

### 3. Summary:

In [146]:

```
1 print(model_summary)
```

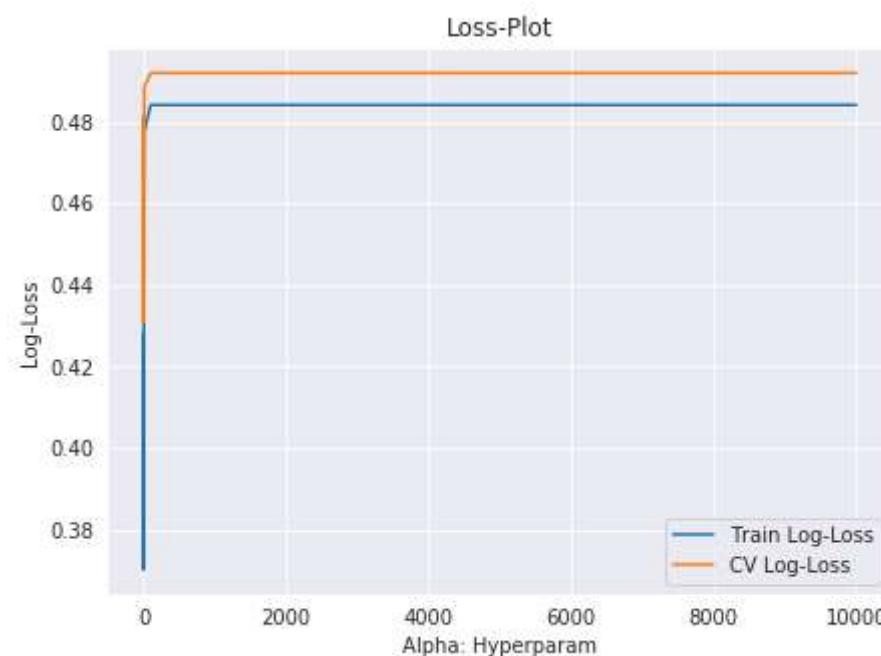
Model	Train(Log-Loss)	Test(Log-Loss)
Random	0.8848	0.8848
Logistic-Regression	0.3822	0.4176

## Linear SVM

### 1.Hyperparam tunning and loss plot :

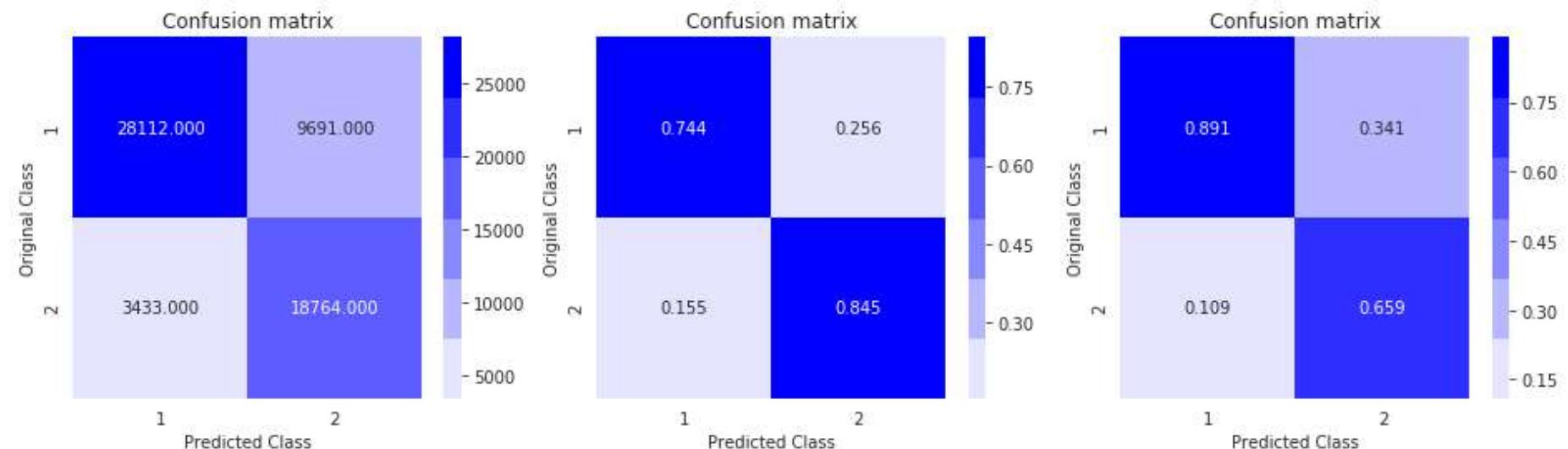
```
In [147]: 1 opt_param, opt_train_loss = linear_model_svm(train, y_train_tf, test, y_test_tf, params)
```

```
for alpha: 0.000010,           train log-loss: 0.4280,           cv log-loss: 0.4813
for alpha: 0.000100,           train log-loss: 0.4275,           cv log-loss: 0.4808
for alpha: 0.001000,           train log-loss: 0.4245,           cv log-loss: 0.4787
for alpha: 0.010000,           train log-loss: 0.3996,           cv log-loss: 0.4576
for alpha: 0.100000,           train log-loss: 0.3703,           cv log-loss: 0.4308
for alpha: 1.000000,           train log-loss: 0.4031,           cv log-loss: 0.4402
for alpha: 10.000000,          train log-loss: 0.4776,           cv log-loss: 0.4886
for alpha: 100.000000,          train log-loss: 0.4841,           cv log-loss: 0.4919
for alpha: 1000.000000,         train log-loss: 0.4841,           cv log-loss: 0.4919
for alpha: 10000.000000,        train log-loss: 0.4841,           cv log-loss: 0.4919
```



## 2.Performance on test data:

In [148]: 1 test\_performance\_linear(opt\_param, train, y\_train\_tf, test, y\_test\_tf, model\_summary, opt\_train\_loss, model\_use[1],



test log-loss: 0.4196

## 2. summary:

In [149]: 1 print(model\_summary)

Model	Train(Log-Loss)	Test(Log-Loss)
Random	0.8848	0.8848
Logistic-Regression	0.3822	0.4176
Linear-SVM	0.3703	0.4196

## XGBOOST:

### Initialize common objects:

In [37]:

```

1 #DICT OF PARAMETERS TO BE TUNED, INITIALLY SET REASONABLE VALUES FOR PARAMETER
2 params_={
3     'n_estimators':128,
4     'max_depth':5,
5     'min_child_weight':1,
6     'gamma':0,
7     'subsample':.8,
8     'colsample_bytree':.8,
9     'reg_alpha':.1
10 }
11
12 #SEARCH METHOD
13 searchMethod='random'
14 #TRAIN AND TEST DATA
15 train=X_train;test=X_test;

```

**[a.] Tune n\_estimators(no. of base learners):**

In [42]:

```

1 %%time
2 #BASE ESTIMATOR RANGE
3 tune_param={}
4 tune_param['n_estimators']=[64,128,256]
5
6 #TUNNING HYPERPARAM
7 model=Ensemble_Classifier(train,y_train,CV,params_, tune_param,searchMethod,param_name=1)
8
9 #UPDATE OPTIMAL VALUE OF PARAMETER
10 params_['n_estimators'] = model.best_params_['n_estimators']

```

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

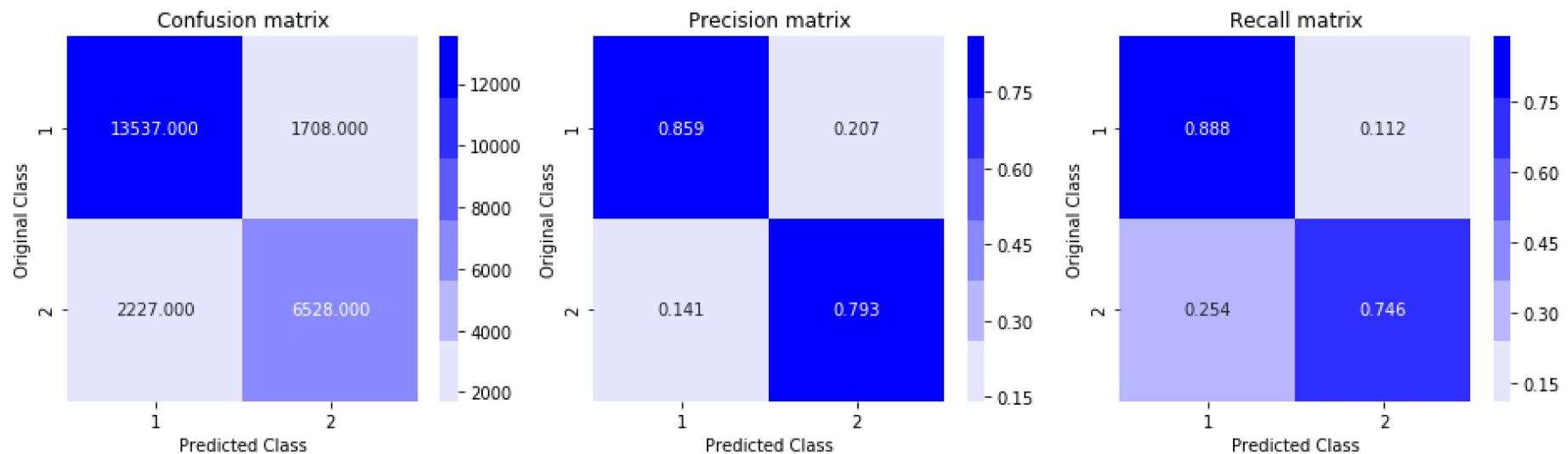
[Parallel(n\_jobs=-1)]: Using backend LokyBackend with 24 concurrent workers.  
[Parallel(n\_jobs=-1)]: Done 14 out of 30 | elapsed: 3.1min remaining: 3.6min  
[Parallel(n\_jobs=-1)]: Done 30 out of 30 | elapsed: 6.0min finished

Train Log-Loss: 0.3344 CV Log-Loss: 0.3532  
Train Log-Loss: 0.2962 CV Log-Loss: 0.3417  
Train Log-Loss: 0.2441 CV Log-Loss: 0.3363

**Test Performance:**

```
In [46]: 1 %time test_performance(train, y_train, test, y_test, params_ )
```

test\_loss: 0.33175645515301455



[b.] Tune `max_depth` and `min_child_weight`:

In [48]:

```
1 %%time
2 #HYPERPARAM RANGE
3 tune_param={}
4 tune_param['max_depth']=[5,7,9]
5 tune_param['min_child_weight']=[1,3,5]
6
7 #TUNNING HYPERPARAM
8 model=Ensemble_Classifier(train, y_train, CV, params_, tune_param, searchMethod, param_name=2)
9
10 #UPDATE OTIMAL VALUE OF PARAMETER
11 params_['max_depth'] = model.best_params_['max_depth']
12 params_['min_child_weight'] = model.best_params_['min_child_weight']
```

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

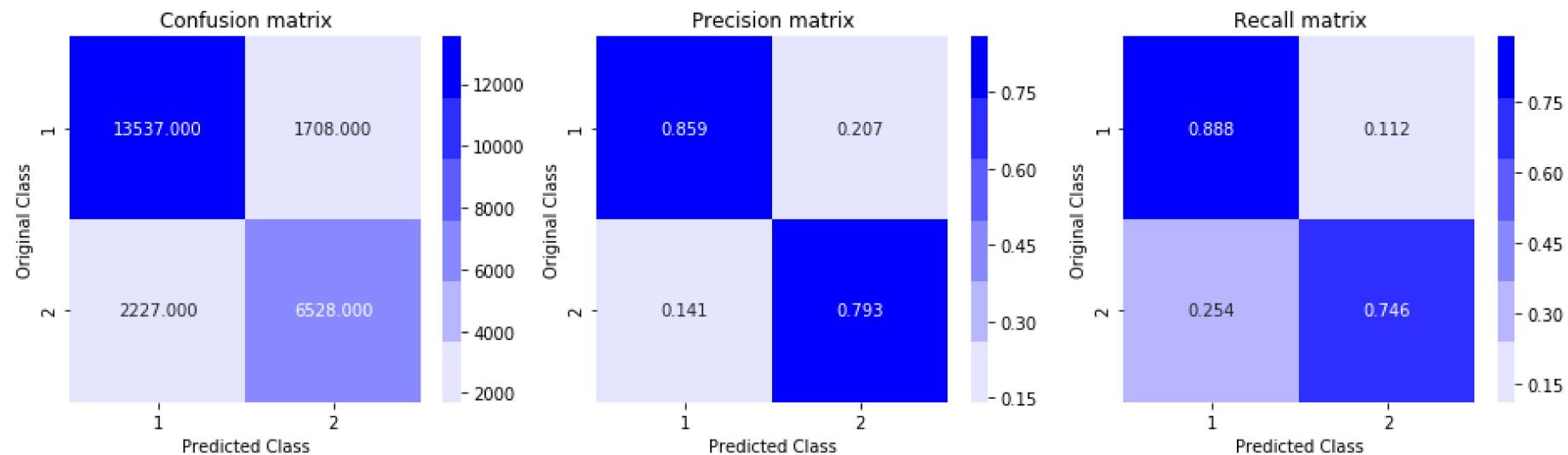
```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 24 concurrent workers.
[Parallel(n_jobs=-1)]: Done 14 out of 30 | elapsed: 6.6min remaining: 7.5min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 13.5min finished
```

```
Train Log-Loss: 0.2441 CV Log-Loss: 0.3363
Train Log-Loss: 0.2482 CV Log-Loss: 0.3367
Train Log-Loss: 0.0645 CV Log-Loss: 0.3394
```

### Test Performance:

```
In [49]: 1 %time test_performance(train, y_train, test, y_test,params_ )
```

test\_loss: 0.33175645515301455



CPU times: user 4min 15s, sys: 2.35 s, total: 4min 17s

Wall time: 4min 14s

**[c.] Tune gamma:**

```
In [52]: 1 #HYPERPARAM RANGE
2 tune_param={}
3 tune_param['gamma']=[i/10.0 for i in range(0,5)]
4
5 #TUNNING HYPERPARAM
6 model=Ensemble_Classifier(train,y_train,CV,params_,tune_param,searchMethod, param_name=1)
7
8 #UPDATE OTIMAL VALUE OF PARAMETER
9 params_['gamma'] = model.best_params_['gamma']
```

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

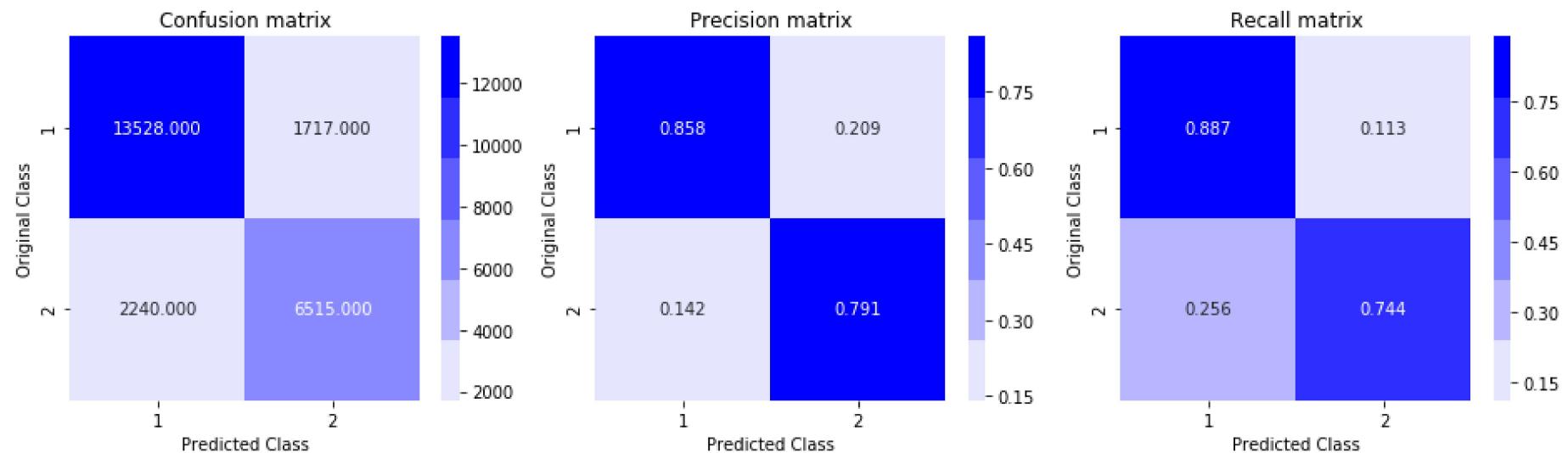
```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 24 concurrent workers.
[Parallel(n_jobs=-1)]: Done 14 out of 30 | elapsed: 6.6min remaining: 7.6min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 10.3min finished
```

```
Train Log-Loss: 0.2445 CV Log-Loss: 0.3362
Train Log-Loss: 0.2441 CV Log-Loss: 0.3363
Train Log-Loss: 0.2440 CV Log-Loss: 0.3362
```

### Test Performance:

```
In [53]: 1 test_performance(train, y_train, test, y_test,params_ )
```

test\_loss: 0.33178602452932426



[d.] Tune subsample and colsample\_bytree:

In [57]:

```
1 %%time
2 #HYPERPARAM RANGE
3 tune_param={}
4 tune_param['subsample']=[.7,.8,.9]#[i/10.0 for i in range(6,10)],
5 tune_param['colsample_bytree']=[.7,.8,.9]#[i/10.0 for i in range(6,10)]
6
7 #TUNNING HYPERPARAM
8 model=Ensemble_Classifier(train,y_train,CV,params_,tune_param,searchMethod, param_name=2)
9
10 #UPDATE OTIMAL VALUE OF PARAMETER
11 params_['subsample'] = model.best_params_['subsample']
12 params_['colsample_bytree'] = model.best_params_['colsample_bytree']
```

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

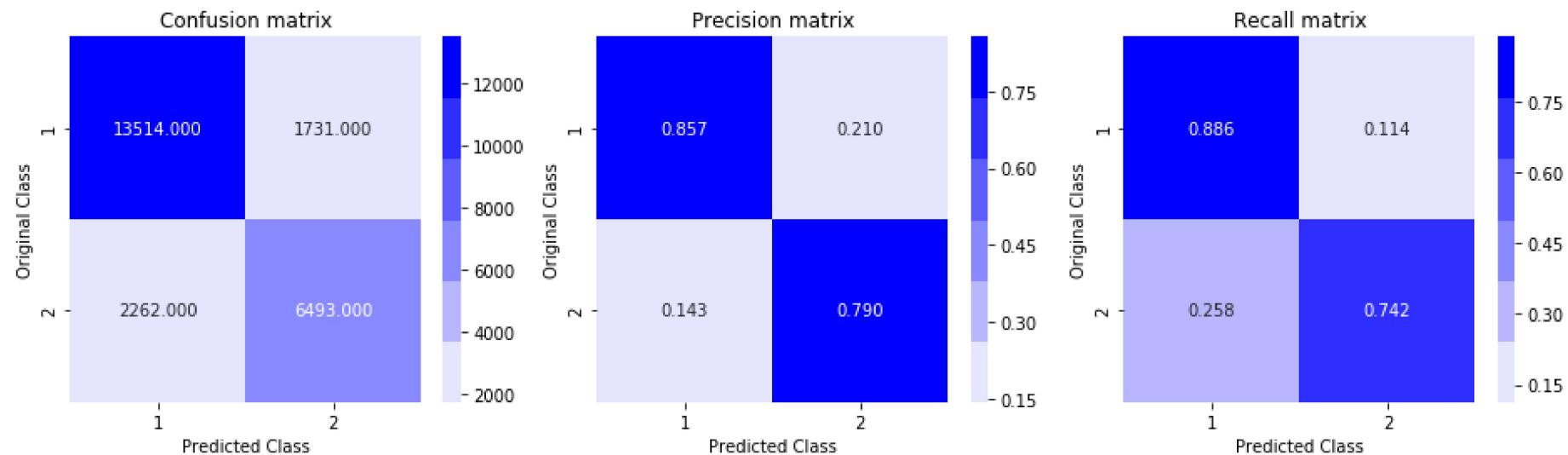
```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 24 concurrent workers.
[Parallel(n_jobs=-1)]: Done 14 out of 30 | elapsed: 7.5min remaining: 8.6min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 11.6min finished
```

```
Train Log-Loss: 0.2450 CV Log-Loss: 0.3360
Train Log-Loss: 0.2434 CV Log-Loss: 0.3367
Train Log-Loss: 0.2434 CV Log-Loss: 0.3366
CPU times: user 4min 43s, sys: 1.39 s, total: 4min 45s
Wall time: 16min 16s
```

### Test Performance:

```
In [58]: 1 test_performance(train, y_train, test, y_test,params_ )
```

test\_loss: 0.33225498789842156



[e.] Tune reg\_alpha:

In [59]:

```
1 #HYPERPARAM RANGE
2 tune_param={}
3 tune_param['reg_alpha']=[0, 0.001, 0.005, 0.01, 0.05]
4
5 #TUNNING HYPERPARAM
6 model=Ensemble_Classifier(train,y_train,CV,params_,tune_param,searchMethod, param_name=1)
7
8 #UPDATE OTIMAL VALUE OF PARAMETER
9 params_['reg_alpha'] = model.best_params_['reg_alpha']
```

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

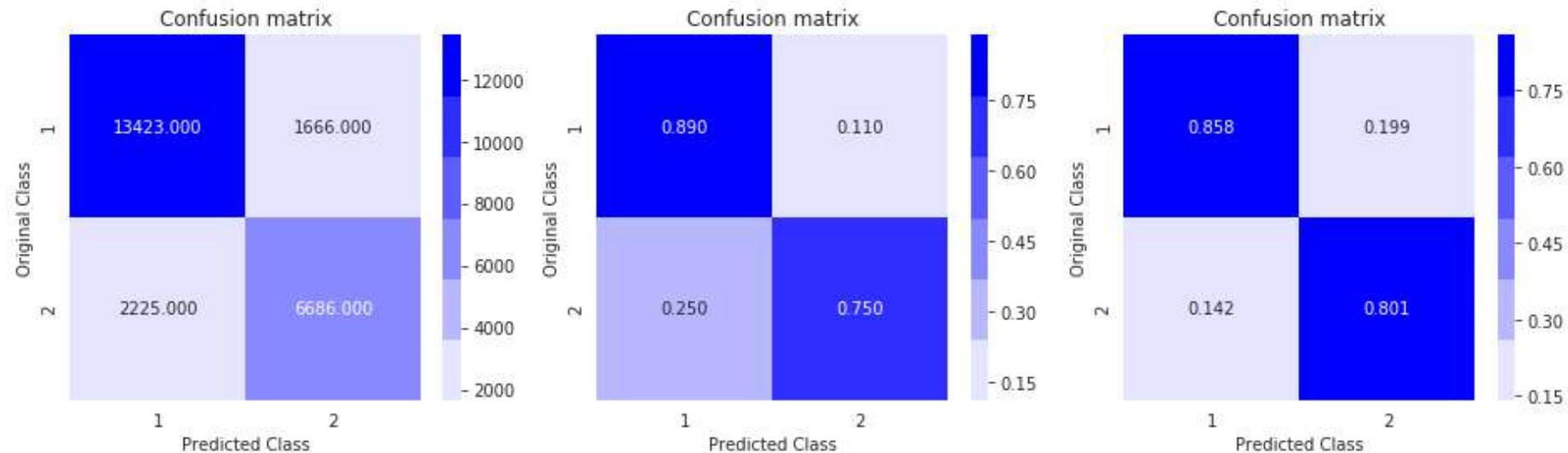
```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 24 concurrent workers.
[Parallel(n_jobs=-1)]: Done 14 out of 30 | elapsed: 7.7min remaining: 8.8min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 11.9min finished
```

```
Train Log-Loss: 0.2455 CV Log-Loss: 0.3359
Train Log-Loss: 0.2447 CV Log-Loss: 0.3360
Train Log-Loss: 0.2452 CV Log-Loss: 0.3359
```

### Test Performance:

In [154]: 1 test\_performance(train, y\_train, test, y\_test, params\_, model\_use[2], True)

test\_loss: 0.33007172532686385



### Summary:

In [155]: 1 print(model\_summary)

Model	Train(Log-Loss)	Test(Log-Loss)
Random	0.8848	0.8848
Logistic-Regression	0.3822	0.4176
Linear-SVM	0.3703	0.4196
XGBoost	0.2561	0.3301

### Save final optimal model:

```
In [151]: 1 saveModelToFile(model, 'XGBOOST_TUNED_MODEL')
2 model=openModelFromFile('XGBOOST_TUNED_MODEL')
3 train=X_train;test=X_test;
```

## Conclusion:

```
In [156]: 1 print(model_summary)
```

Model	Train(Log-Loss)	Test(Log-Loss)
Random	0.8848	0.8848
Logistic-Regression	0.3822	0.4176
Linear-SVM	0.3703	0.4196
XGBoost	0.2561	0.3301

**XGBOOST performs very well with log-loss:**

- a. train log-loss = .2561
- b. test log-loss = .3301

## Procedure:

1. We have to solve problem, given a pair of questions these question are similar or not.
2. We proceed with basic details of dataset like :
  - a. How many data points?
  - b. Is data is balanced or imbalanced?
  - c. How many null values are there?
3. We do Exploratory Data Analysis and Feature Extraction :

- a. We extract Frequency based features like length of question, frequency of question, number of words in question, word share etc
- b. We extract more advanced features using fuzzywuzzy (ex fuzz-ratio, fuzz-partial-ratio, token-sort-ratio etc)
- c. We did EDA on Extracted fetures (PDF, Violin plot, Histogram etc).

4. We featurize textual data:

- a. TF-IDF Weighted Word2VEC (using spacy pretrained GLO-VEC(Global Vectors))
- b. TF-IDF Vector (using sklearn)

5. We build different ML Models on extracted features (Textual + Freq Based + FuzzyWuzzy):

- a. Logistic Regression(TF-IDF)
- b. Linear SVM(TF-IDF)
- c. XGBOOST(TFIDF-W2V)

Reference Links:

1. <https://www.appliedaicourse.com/> (<https://www.appliedaicourse.com/>)
2. <https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/> (<https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/>)
3. <https://spacy.io/usage/vectors-similarity> (<https://spacy.io/usage/vectors-similarity>)
4. <https://towardsdatascience.com/fine-tuning-xgboost-in-python-like-a-boss-b4543ed8b1e> (<https://towardsdatascience.com/fine-tuning-xgboost-in-python-like-a-boss-b4543ed8b1e>)
5. <https://machinelearningmastery.com/tune-number-size-decision-trees-xgboost-python/> (<https://machinelearningmastery.com/tune-number-size-decision-trees-xgboost-python/>)