



```
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 nltk.corpus import stopwords
14 from sklearn.decomposition import TruncatedSVD
15 from sklearn.preprocessing import normalize
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.manifold import TSNE
18 import seaborn as sns
19 from sklearn.neighbors import KNeighborsClassifier
20 from sklearn.metrics import confusion_matrix
21 from sklearn.metrics.classification import accuracy_score, log_loss
22 from sklearn.feature_extraction.text import TfidfVectorizer
23 from collections import Counter
24 from scipy.sparse import hstack
25 from sklearn.multiclass import OneVsRestClassifier
26 from sklearn.svm import SVC
27 from sklearn.model_selection import StratifiedKFold
28 from collections import Counter, defaultdict
29 from sklearn.calibration import CalibratedClassifierCV
30 from sklearn.naive_bayes import MultinomialNB
31 from sklearn.naive_bayes import GaussianNB
32 from sklearn.model_selection import train_test_split
33 from sklearn.model_selection import GridSearchCV
34 import math
```

```
35 from sklearn.metrics import normalized_mutual_info_score
36 from sklearn.ensemble import RandomForestClassifier
37
38
39
40 from sklearn.model_selection import cross_val_score
41 from sklearn.linear_model import SGDClassifier
42 from mlxtend.classifier import StackingClassifier
43
44 from sklearn import model_selection
45 from sklearn.linear_model import LogisticRegression
46 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

```
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=1
6     j = 0
7     index_start = 1
8     for df in pd.read_csv('final_features.csv', names=['Unnamed: 0', 'id', 'is_duplicate', 'cwc_min', 'cwc_i
9         chunksize=chunksize, iterator=True, encoding='utf-8', nrows=10000):
10         df.index += index_start
11         j+=1
12         print('{} rows'.format(j*chunksize))
13         df.to_sql('data', disk_engine, if_exists='append')
14         index_start = df.index[-1] + 1
```

1 rows  
2 rows  
3 rows  
4 rows  
5 rows  
6 rows  
7 rows  
8 rows  
9 rows  
10 rows  
11 rows  
12 rows  
13 rows  
14 rows  
15 rows  
16 rows  
17 rows  
18 rows  
19 rows  
20 rows  
21 rows

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

```
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

```

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

```

```

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

```

```
1  data.shape
```

```
(99999, 794)
```

```
1  data.head(5)
```

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max
1	0.999980000399992	0.833319444675922	0.999983333611106	0.999983333611106	0.916659027841435	0.78570867350947
2	0.799984000319994	0.39999600004	0.749981250468738	0.599988000239995	0.699993000069999	0.46666355557629
3	0.399992000159997	0.333327777870369	0.399992000159997	0.249996875039062	0.39999600004	0.28571224491253
4	0.0	0.0	0.0	0.0	0.0	0.0
5	0.399992000159997	0.19999800002	0.999950002499875	0.66664444518516	0.571420408279882	0.30768994084660

```
5 rows x 794 columns
```

## 4.2 Converting strings to numerics

```
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 cols = list(data.columns)
4 for i in cols:
5     data[i] = data[i].apply(pd.to_numeric)
6     print(i)
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
```

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

## 4.3 Random train test split( 80:20)

```
1 X_train,X_test, y_train, y_test = train_test_split(data, y_true, stratify=y_true, test_size=0.2)
```

```
1 print("Number of data points in train data :",X_train.shape)
2 print("Number of data points in test data :",X_test.shape)
```

Number of data points in train data : (79999, 794)

Number of data points in test data : (20000, 794)

```
1 print("-"*10, "Distribution of output variable in train data", "-"*10)
2 train_distr = Counter(y_train)
3 train_len = len(y_train)
4 print("Class 0: ",int(train_distr[0])/train_len,"Class 1: ", int(train_distr[1])/train_len)
5 print("-"*10, "Distribution of output variable in train data", "-"*10)
6 test_distr = Counter(y_test)
7 test_len = len(y_test)
8 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: 0.6274578432230403 Class 1: 0.37254215677695973

----- Distribution of output variable in train data -----

Class 0: 0.37255 Class 1: 0.37255



```

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
6      A = (((C.T)/(C.sum(axis=1))).T)
7      #divid each element of the confusion matrix with the sum of elements in that column
8
9      # C = [[1, 2],
10     #      [3, 4]]
11     # C.T = [[1, 3],
12     #        [2, 4]]
13     # C.sum(axis = 1)  axis=0 corresonds to columns and axis=1 corresponds to rows in two dimensional
14     # C.sum(axix =1) = [[3, 7]]
15     # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
16     #                             [2/3, 4/7]]
17
18     # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
19     #                             [3/7, 4/7]]
20     # sum of row elements = 1
21
22     B = (C/C.sum(axis=0))
23     #divid each element of the confusion matrix with the sum of elements in that row
24     # C = [[1, 2],
25     #      [3, 4]]
26     # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to rows in two dimensional
27     # C.sum(axix =0) = [[4, 6]]
28     # (C/C.sum(axis=0)) = [[1/4, 2/6],
29     #                       [3/4, 4/6]]
30     plt.figure(figsize=(20,4))
31
32     labels = [1,2]
33     # representing A in heatmap format
34     cmap=sns.light_palette("blue")

```

```
35 plt.subplot(1, 3, 1)
36 sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
37 plt.xlabel('Predicted Class')
38 plt.ylabel('Original Class')
39 plt.title("Confusion matrix")
40
41 plt.subplot(1, 3, 2)
42 sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
43 plt.xlabel('Predicted Class')
44 plt.ylabel('Original Class')
45 plt.title("Precision matrix")
46
47 plt.subplot(1, 3, 3)
48 # representing B in heatmap format
49 sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
50 plt.xlabel('Predicted Class')
51 plt.ylabel('Original Class')
52 plt.title("Recall matrix")
53
54 plt.show()
```

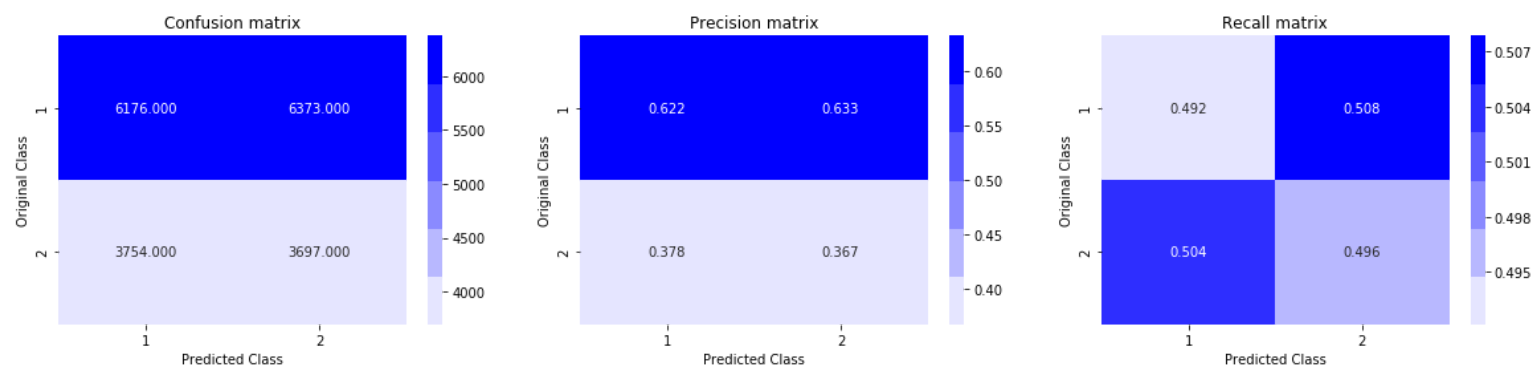
## 4.4 Building a random model (Finding worst-case log-loss)

```

1  # we need to generate 9 numbers and the sum of numbers should be 1
2  # one solution is to generate 9 numbers and divide each of the numbers by their sum
3  # ref: https://stackoverflow.com/a/18662466/4084039
4  # we create a output array that has exactly same size as the CV data
5  predicted_y = np.zeros((test_len,2))
6  for i in range(test_len):
7      rand_probs = np.random.rand(1,2)
8      predicted_y[i] = ((rand_probs/sum(sum(rand_probs))))[0])
9  print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y, eps=1e-15))
10
11 predicted_y =np.argmax(predicted_y, axis=1)
12 plot_confusion_matrix(y_test, predicted_y)

```

Log loss on Test Data using Random Model 0.9007478073282229



## Logistic Regression with hyperparameter tuning

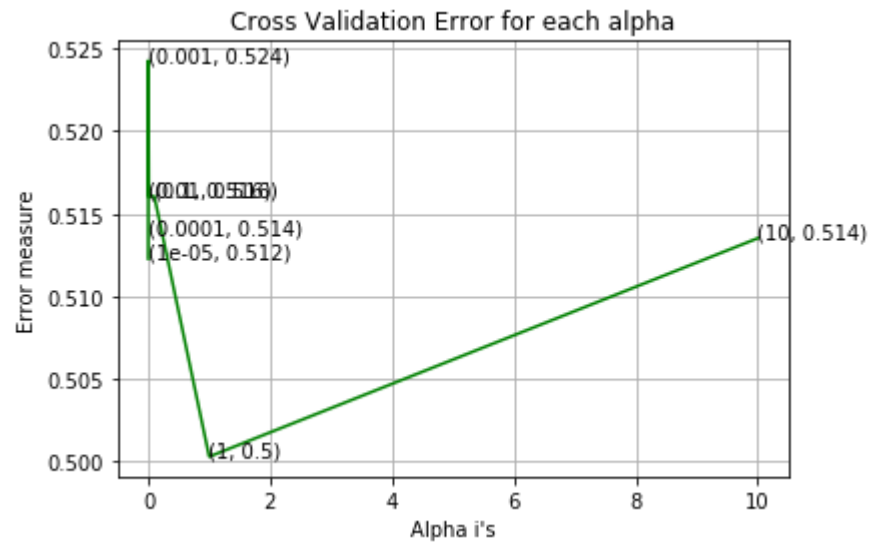
```

1  alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=1000,
7  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
8  # class_weight=None, warm_start=False, average=False, n_iter=None)
9
10 # some of methods
11 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
12 # predict(X) Predict class labels for samples in X.
13
14 #-----
15 # video link:
16 #-----
17
18
19 log_error_array=[]
20 for i in alpha:
21     clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
22     clf.fit(X_train, y_train)
23     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
24     sig_clf.fit(X_train, y_train)
25     predict_y = sig_clf.predict_proba(X_test)
26     log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
27     print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
28
29 fig, ax = plt.subplots()
30 ax.plot(alpha, log_error_array, c='g')
31 for i, txt in enumerate(np.round(log_error_array, 3)):
32     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
33 plt.grid()
34 plt.title("Cross Validation Error for each alpha")

```

```
35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
42 clf.fit(X_train, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(X_train, y_train)
45
46 predict_y = sig_clf.predict_proba(X_train)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
48 predict_y = sig_clf.predict_proba(X_test)
49 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
50 predicted_y = np.argmax(predict_y, axis=1)
51 print("Total number of data points :", len(predicted_y))
52 plot_confusion_matrix(y_test, predicted_y)
```

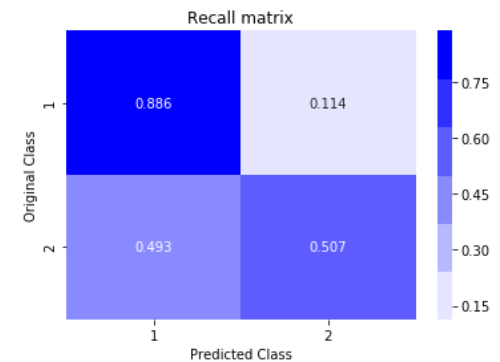
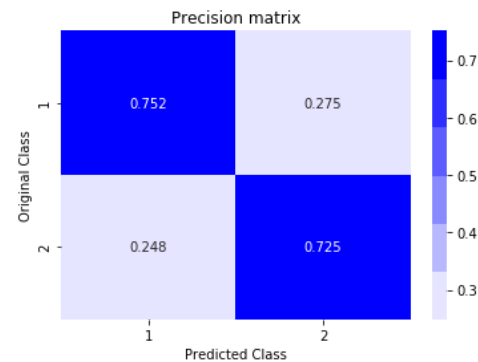
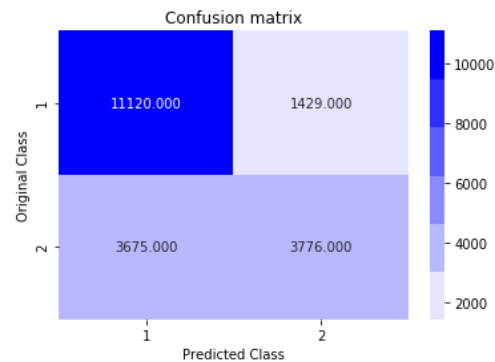
```
For values of alpha = 1e-05 The log loss is: 0.5122785366740716
For values of alpha = 0.0001 The log loss is: 0.5138446543654717
For values of alpha = 0.001 The log loss is: 0.5242688922957729
For values of alpha = 0.01 The log loss is: 0.5160512998933138
For values of alpha = 0.1 The log loss is: 0.5160503173962908
For values of alpha = 1 The log loss is: 0.5002981377384517
For values of alpha = 10 The log loss is: 0.513515646484234
```



For values of best alpha = 1 The train log loss is: 0.49672963589730734

For values of best alpha = 1 The test log loss is: 0.5002981377384517

Total number of data points : 20000



## Linear SVM with hyperparameter tuning

```

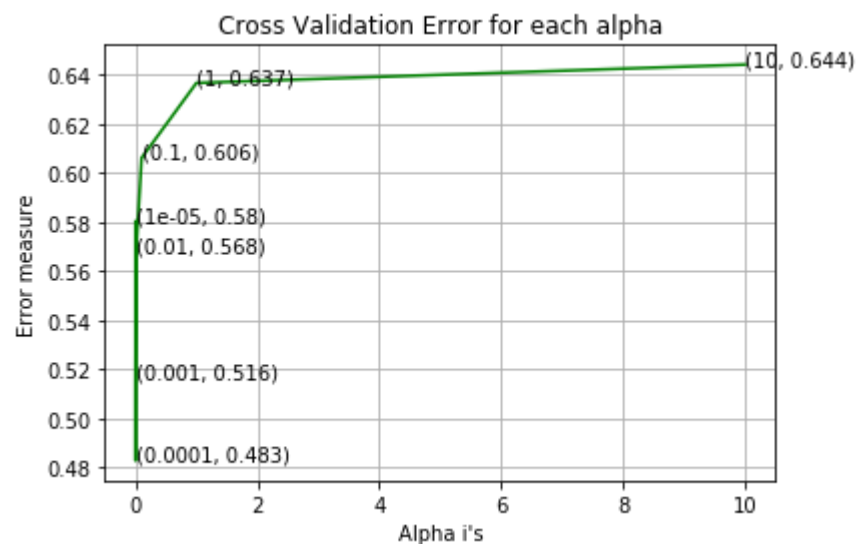
1  alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=1000,
7  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
8  # class_weight=None, warm_start=False, average=False, n_iter=None)
9
10 # some of methods
11 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
12 # predict(X) Predict class labels for samples in X.
13
14 #-----
15 # video link:
16 #-----
17
18
19 log_error_array=[]
20 for i in alpha:
21     clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
22     clf.fit(X_train, y_train)
23     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
24     sig_clf.fit(X_train, y_train)
25     predict_y = sig_clf.predict_proba(X_test)
26     log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
27     print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
28
29 fig, ax = plt.subplots()
30 ax.plot(alpha, log_error_array, c='g')
31 for i, txt in enumerate(np.round(log_error_array, 3)):
32     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
33 plt.grid()
34 plt.title("Cross Validation Error for each alpha")

```

```
35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l1', loss='hinge', random_state=42)
42 clf.fit(X_train, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(X_train, y_train)
45
46 predict_y = sig_clf.predict_proba(X_train)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
48 predict_y = sig_clf.predict_proba(X_test)
49 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
50 predicted_y = np.argmax(predict_y, axis=1)
51 print("Total number of data points :", len(predicted_y))
52 plot_confusion_matrix(y_test, predicted_y)
```

```
For values of alpha = 1e-05 The log loss is: 0.5802471570896839
For values of alpha = 0.0001 The log loss is: 0.4828094006631191
For values of alpha = 0.001 The log loss is: 0.516201285754687
For values of alpha = 0.01 The log loss is: 0.5678119444737264
For values of alpha = 0.1 The log loss is: 0.6060723428406178
For values of alpha = 1 The log loss is: 0.6366103451105977
For values of alpha = 10 The log loss is: 0.644078585611343
```

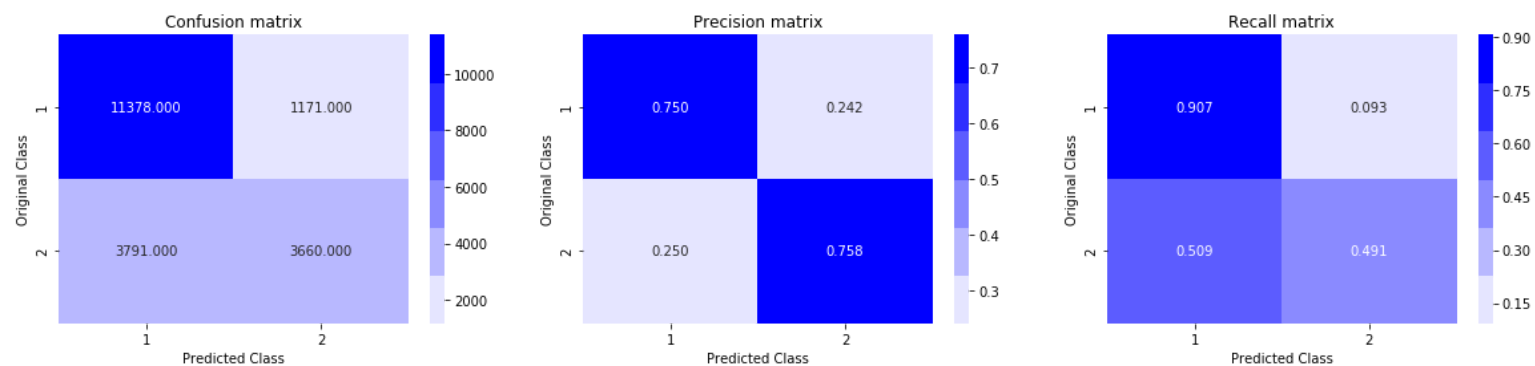




For values of best alpha = 0.0001 The train log loss is: 0.4777288603749968

For values of best alpha = 0.0001 The test log loss is: 0.4828094006631191

Total number of data points : 20000



## Hyper parameter tuning for XGBOOST

```
1 from sklearn.model_selection import RandomizedSearchCV
2 from sklearn.metrics import log_loss
3 from xgboost import XGBClassifier
```

```

1  n_estimators=list(range(100,500,100))
2  learning_rate=[0.01,0.02,0.03,0.04,0.05,0.06,0.07,0.08,0.09,0.1,0.2,0.3]
3  max_depth=list(range(3,20,1))
4  gamma=[i/10.0 for i in range(0,5)]
5  min_child_weight =list(range(1,20,1))
6  subsample=[0.5,0.6,0.7,0.8,0.9,1.0]
7  colsample_bytree=[0.5,0.6,0.7,0.8,0.9,1.0]
8  scale_pos_weight=list(range(0,5,1))
9  reg_alpha= [1e-5,1e-4,1e-3 ,1e-2, 0.1, 1, 100]
10
11 param_distributions = dict(n_estimators=n_estimators,max_depth=max_depth,learning_rate=learning_rate,g
12                             min_child_weight=min_child_weight,subsample=subsample,colsample_bytree=colsa
13                             scale_pos_weight=scale_pos_weight,reg_alpha=reg_alpha)
14 print(param_distributions)
15
16 # instantiate and fit the grid
17 grid = RandomizedSearchCV(XGBClassifier(), param_distributions, cv=3, scoring='neg_log_loss', return_tr

```

```

{'n_estimators': [100, 200, 300, 400], 'max_depth': [3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19], 'learnin
g_rate': [0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3], 'gamma': [0.0, 0.1, 0.2, 0.3, 0.4], 'min_ch
ild_weight': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19], 'subsample': [0.5, 0.6, 0.7, 0.8, 0.9, 1.
0], 'colsample_bytree': [0.5, 0.6, 0.7, 0.8, 0.9, 1.0], 'scale_pos_weight': [0, 1, 2, 3, 4], 'reg_alpha': [1e-05, 0.0001,
0.001, 0.01, 0.1, 1, 100]}

```

```

1  grid.fit(X_train,y_train)
2
3  # examine the best model
4  print(grid.best_score_)
5  print(grid.best_params_)

```

```
-0.34538638494778556
```

```

{'subsample': 0.7, 'scale_pos_weight': 2, 'reg_alpha': 0.01, 'n_estimators': 100, 'min_child_weight': 16, 'max_depth': 17,
'learning_rate': 0.08, 'gamma': 0.3, 'colsample_bytree': 0.8}

```

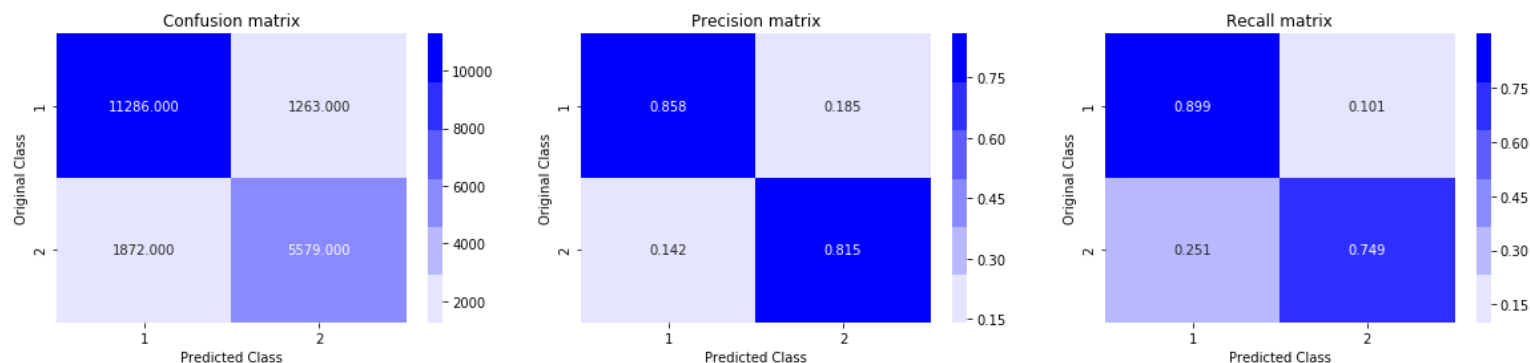
## XGBOOST

```
1 import xgboost as xgb
2 params = {}
3 params['n_estimators'] = 100
4 params['min_child_weight'] = 16
5 params['subsample'] = 0.7
6 params['colsample_bytree'] = 0.8
7 params['gamma'] = 0.3
8 params['reg_alpha'] = 0.01
9 params['objective'] = 'binary:logistic'
10 params['eval_metric'] = 'logloss'
11 params['eta'] = 0.01
12 params['max_depth'] = 17
13
14 d_train = xgb.DMatrix(X_train, label=y_train)
15 d_test = xgb.DMatrix(X_test, label=y_test)
16
17 watchlist = [(d_train, 'train'), (d_test, 'valid')]
18
19 bst = xgb.train(params, d_train, 400, watchlist, early_stopping_rounds=20, verbose_eval=10)
20
21 xgdmatrix = xgb.DMatrix(X_train, y_train)
22 predict_y = bst.predict(d_test)
23 print("The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
[10:04:59] C:\Users\Administrator\Desktop\xgboost\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 586 extra nodes, 2 pruned nodes, max_depth=17
[10:05:09] C:\Users\Administrator\Desktop\xgboost\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 530 extra nodes, 2 pruned nodes, max_depth=17
```

```
1 predicted_y = np.array(predict_y>0.5, dtype=int)
2 print("Total number of data points :", len(predicted_y))
3 plot_confusion_matrix(y_test, predicted_y)
```

Total number of data points : 20000



## Importing Simple Tf-idf features file

Lets try Logistic regression and Linear SVM on it

```
1 from scipy import sparse
2 final_features_tfidf = sparse.load_npz("final_features_tfidf.npz")
```

```
1 final_features_tfidf
```

```
<100000x163096 sparse matrix of type '<class 'numpy.float64''>'
  with 4505604 stored elements in Compressed Sparse Row format>
```

```
1 y_true=pd.read_csv('y_true.csv',header=None)
2 y_true=pd.DataFrame(y_true)
```

```
1 y_true.head(5)
```

	0	1
0	121958	1
1	146867	0
2	131932	1
3	365838	1
4	259178	0

```
1 y_true.drop([0],axis=1,inplace=True)
```

```
1 y_true.head(5)
```

	1
0	1
1	0
2	1
3	1
4	0

```
1 X_train,X_test, y_train, y_test = train_test_split(final_features_tfidf, y_true, stratify=y_true, test_:
```

## Logistic Regression with hyperparameter tuning

```

1  alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=1000,
7  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
8  # class_weight=None, warm_start=False, average=False, n_iter=None)
9
10 # some of methods
11 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
12 # predict(X) Predict class labels for samples in X.
13
14 #-----
15 # video link:
16 #-----
17
18
19 log_error_array=[]
20 for i in alpha:
21     clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
22     clf.fit(X_train, y_train)
23     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
24     sig_clf.fit(X_train, y_train)
25     predict_y = sig_clf.predict_proba(X_test)
26     log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
27     print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
28
29 fig, ax = plt.subplots()
30 ax.plot(alpha, log_error_array, c='g')
31 for i, txt in enumerate(np.round(log_error_array, 3)):
32     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
33 plt.grid()
34 plt.title("Cross Validation Error for each alpha")

```

```

35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
42 clf.fit(X_train, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(X_train, y_train)
45
46 predict_y = sig_clf.predict_proba(X_train)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
48 predict_y = sig_clf.predict_proba(X_test)
49 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
50 predicted_y = np.argmax(predict_y, axis=1)
51 print("Total number of data points :", len(predicted_y))
52 plot_confusion_matrix(y_test, predicted_y)

```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

For values of alpha = 1e-05 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

For values of alpha = 0.0001 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

assed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 0.001 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 0.01 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 0.1 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 1 The log loss is: 0.6449071484847831

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

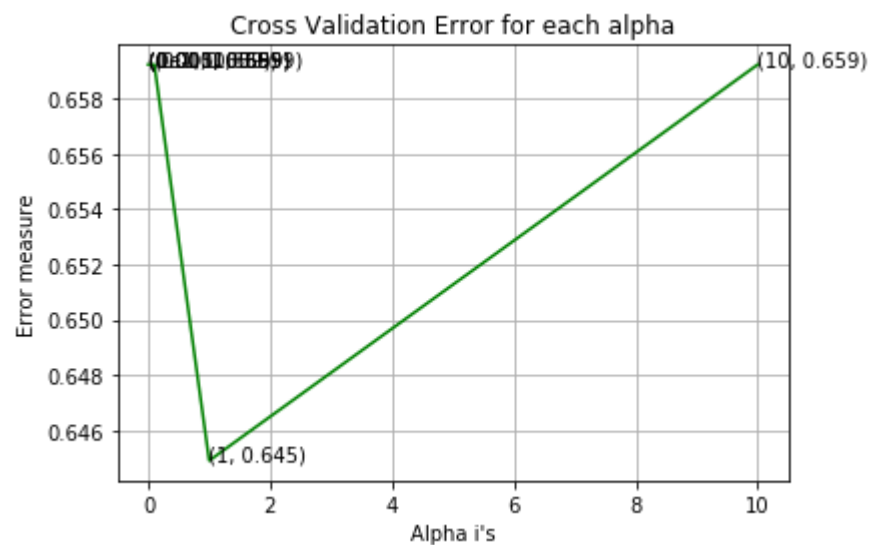
```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```



For values of alpha = 10 The log loss is: 0.6592212539965192



C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

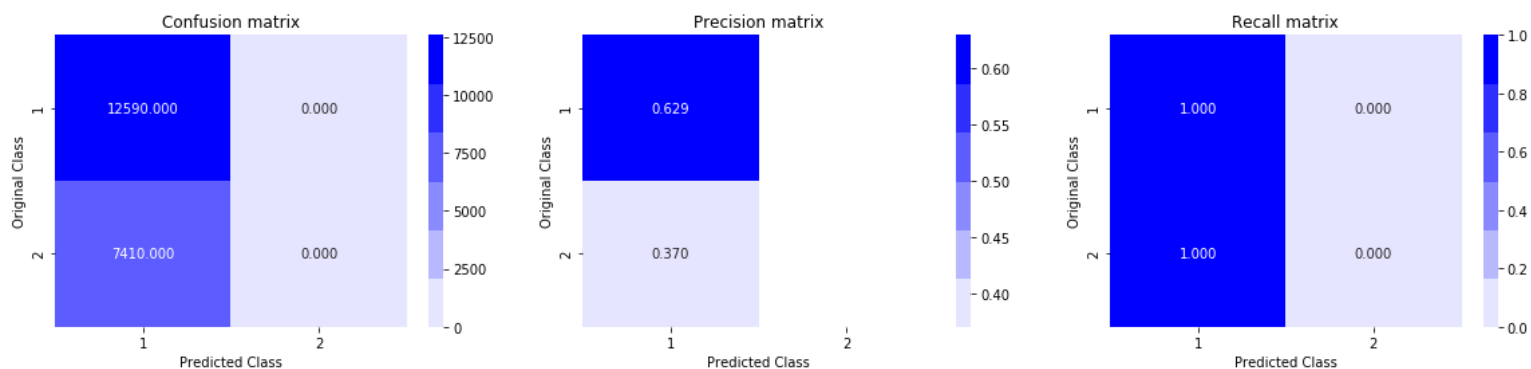
C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of best alpha = 1 The train log loss is: 0.6448625946071465

For values of best alpha = 1 The test log loss is: 0.6449071484847831

Total number of data points : 20000



## Linear SVM with hyperparameter tuning

```

1  alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=1000,
7  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
8  # class_weight=None, warm_start=False, average=False, n_iter=None)
9
10 # some of methods
11 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
12 # predict(X) Predict class labels for samples in X.
13
14 #-----
15 # video link:
16 #-----
17
18
19 log_error_array=[]
20 for i in alpha:
21     clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
22     clf.fit(X_train, y_train)
23     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
24     sig_clf.fit(X_train, y_train)
25     predict_y = sig_clf.predict_proba(X_test)
26     log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
27     print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
28
29 fig, ax = plt.subplots()
30 ax.plot(alpha, log_error_array, c='g')
31 for i, txt in enumerate(np.round(log_error_array, 3)):
32     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
33 plt.grid()
34 plt.title("Cross Validation Error for each alpha")

```

```

35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l1', loss='hinge', random_state=42)
42 clf.fit(X_train, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(X_train, y_train)
45
46 predict_y = sig_clf.predict_proba(X_train)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
48 predict_y = sig_clf.predict_proba(X_test)
49 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
50 predicted_y = np.argmax(predict_y, axis=1)
51 print("Total number of data points :", len(predicted_y))
52 plot_confusion_matrix(y_test, predicted_y)

```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

For values of alpha = 1e-05 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

y = column\_or\_1d(y, warn=True)

For values of alpha = 0.0001 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

assed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 0.001 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 0.01 The log loss is: 0.6592212539965192

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 0.1 The log loss is: 0.620823682677829

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 1 The log loss is: 0.6368395560836245

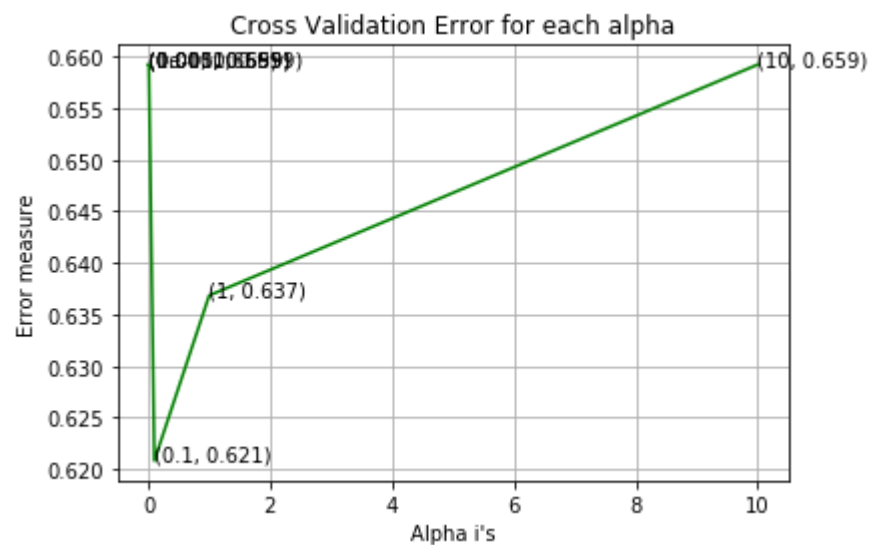
C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of alpha = 10 The log loss is: 0.6592212539965192



C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

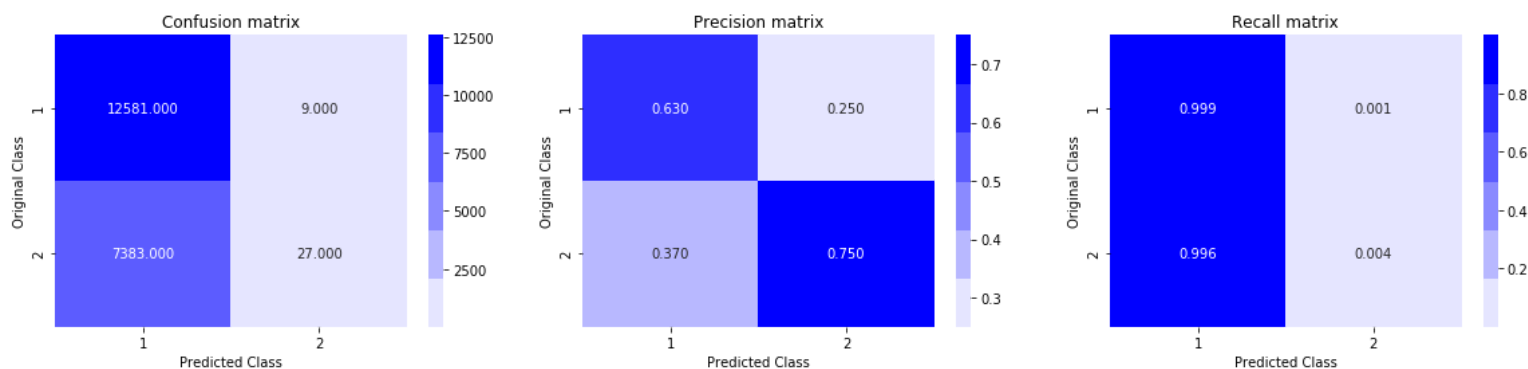
C:\Users\deepak\Anaconda3\lib\site-packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

For values of best alpha = 0.1 The train log loss is: 0.6196195882146978

For values of best alpha = 0.1 The test log loss is: 0.6208236826777829

Total number of data points : 20000



## Summary

```

1  from prettytable import PrettyTable
2  x=PrettyTable()
3  x.field_names = ["Model","Featurization", "Train log-loss", "Test log-loss"]
4  x.add_row(["Logistic Regression","TFIDF-W2V",0.4967,0.5002])
5  x.add_row(["Linear SVM","TFIDF-W2V",0.4777,0.4828])
6  x.add_row(["XGboost","TFIDF-W2V",0.2200,0.3282])
7  x.add_row(["Logistic Regression","TFIDF",0.6448,0.6449])
8  x.add_row(["Linear SVM","TFIDF",0.6196,0.6208])
9
10 print(x)

```

Model	Featurization	Train log-loss	Test log-loss
Logistic Regression	TFIDF-W2V	0.4967	0.5002
Linear SVM	TFIDF-W2V	0.4777	0.4828
XGboost	TFIDF-W2V	0.22	0.3282
Logistic Regression	TFIDF	0.6448	0.6449
Linear SVM	TFIDF	0.6196	0.6208