

Personalized cancer diagnosis

Task 1

1. Business Problem

1.1. Description

Source: <https://www.kaggle.com/c/msk-redefining-cancer-treatment/>

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

Context:

Source: <https://www.kaggle.com/c/msk-redefining-cancer-treatment/discussion/35336#198462>


Problem statement :

Classify the given genetic variations/mutations based on evidence from text-based clinical literature.

1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

1. <https://www.forbes.com/sites/matthewherper/2017/06/03/a-new-cancer-drug-helped-almost-everyone-who-took-it-almost-heres-what-it-teaches-us/#2a44ee2f6b25> (<https://www.forbes.com/sites/matthewherper/2017/06/03/a-new-cancer-drug-helped-almost-everyone-who-took-it-almost-heres-what-it-teaches-us/#2a44ee2f6b25>)

- (http://savefrom.net/?
url=https%3A%2F%2Fwww.youtube.com
2. <https://www.youtube.com/watch?v=UwbuW7oK8rk> (<https://www.youtube.com/watch?v=UwbuW7oK8rk>) chrome&utm_medium=extensions&utm_c
(http://savefrom.net/?
url=https%3A%2F%2Fwww.youtube.com%
3. <https://www.youtube.com/watch?v=qxXRKVompl8> (<https://www.youtube.com/watch?v=qxXRKVompl8>) chrome&utm_medium=extensions&utm_c
- 

1.3. Real-world/Business objectives and constraints.

- No low-latency requirement.
- Interpretability is important.
- Errors can be very costly.
- Probability of a data-point belonging to each class is needed.

2. Machine Learning Problem Formulation

2.1. Data

2.1.1. Data Overview

- Source: <https://www.kaggle.com/c/msk-redefining-cancer-treatment/data> (<https://www.kaggle.com/c/msk-redefining-cancer-treatment/data>)
- We have two data files: one contains the information about the genetic mutations and the other contains the clinical evidence (text) that human experts/pathologists use to classify the genetic mutations.
- Both these data files have a common column called ID
- Data file's information:
 - training_variants (ID , Gene, Variations, Class)
 - training_text (ID, Text)

2.1.2. Example Data Point

training_variants

ID, Gene, Variation, Class

0, FAM58A, Truncating Mutations, 1

1, CBL, W802*, 2

2, CBL, Q249E, 2

...

training_text

ID, Text

0||Cyclin-dependent kinases (CDKs) regulate a variety of fundamental cellular processes. CDK10 stands out as one of the last orphan CDKs for which no activating cyclin has been identified and no kinase activity revealed. Previous work has shown that CDK10 silencing increases ETS2 (v-ets erythroblastosis virus E26 oncogene homolog 2)-driven activation of the MAPK pathway, which confers tamoxifen resistance to breast cancer cells. The precise mechanisms by which CDK10 modulates ETS2 activity, and more generally the functions of CDK10, remain elusive. Here we demonstrate that CDK10 is a cyclin-dependent kinase by identifying cyclin M as an activating cyclin. Cyclin M, an orphan cyclin, is the product of FAM58A, whose mutations cause STAR syndrome, a human developmental anomaly whose features include toe syndactyly, telecanthus, and anogenital and renal malformations. We show that STAR syndrome-associated cyclin M mutants are unable to interact with CDK10. Cyclin M silencing phenocopies CDK10 silencing in increasing c-Raf and in conferring tamoxifen resistance to breast cancer cells. CDK10/cyclin M phosphorylates ETS2 in vitro, and in cells it positively controls ETS2 degradation by the proteasome. ETS2 protein levels are increased in cells derived from a STAR patient, and this increase is attributable to decreased cyclin M levels. Altogether, our results reveal an additional regulatory mechanism for ETS2, which plays key roles in cancer and development. They also shed light on the molecular mechanisms underlying STAR syndrome. Cyclin-dependent kinases (CDKs) play a pivotal role in the control of a number of fundamental cellular processes (1). The human genome contains 21 genes encoding proteins that can be considered as members of the CDK family owing to their sequence similarity with bona fide CDKs, those known to be activated by cyclins (2). Although discovered almost 20 y ago (3, 4), CDK10 remains one of the two CDKs without an identified cyclin partner. This knowledge gap has largely impeded the exploration of its biological functions. CDK10 can act as a positive cell cycle regulator in some cells (5, 6) or as a tumor suppressor in others (7, 8). CDK10 interacts with the ETS2 (v-ets erythroblastosis virus E26 oncogene homolog 2) transcription factor and inhibits its transcriptional activity through an unknown mechanism (9). CDK10 knockdown derepresses ETS2, which increases the expression of the c-Raf protein kinase, activates the MAPK pathway, and induces resistance of MCF7 cells to tamoxifen (6). ...

2.2. Mapping the real-world problem to an ML problem

2.2.1. Type of Machine Learning Problem

There are nine different classes a genetic mutation can be classified into => Multi class classification problem

2.2.2. Performance Metric

Source: <https://www.kaggle.com/c/msk-redefining-cancer-treatment#evaluation> (<https://www.kaggle.com/c/msk-redefining-cancer-treatment#evaluation>)

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

Objective: Predict the probability of each data-point belonging to each of the nine classes.

Constraints:

* Interpretability * Class probabilities are needed. * Penalize the errors in class probabilities => Metric is Log-loss. * No Latency constraints.

2.3. Train, CV and Test Datasets

Split the dataset randomly into three parts train, cross validation and test with 64%, 16%, 20% of data respectively

3. Exploratory Data Analysis

```
In [7]: import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
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 sklearn.linear_model import SGDClassifier
from imblearn.over_sampling import SMOTE
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 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
warnings.filterwarnings("ignore")

from mlxtend.classifier import StackingClassifier

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
```

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

```
In [8]: data = pd.read_csv('training_variants.csv')
print('Number of data points : ', data.shape[0])
print('Number of features : ', data.shape[1])
print('Features : ', data.columns.values)
data.head()
```

```
Number of data points : 3321
Number of features : 4
Features : ['ID' 'Gene' 'Variation' 'Class']
```

Out[8]:

	ID	Gene	Variation	Class
0	0	FAM58A	Truncating Mutations	1
1	1	CBL	W802*	2
2	2	CBL	Q249E	2
3	3	CBL	N454D	3
4	4	CBL	L399V	4

training/training_variants is a comma separated file containing the description of the genetic mutations used for training. Fields are

- **ID** : the id of the row used to link the mutation to the clinical evidence
- **Gene** : the gene where this genetic mutation is located
- **Variation** : the aminoacid change for this mutations
- **Class** : 1-9 the class this genetic mutation has been classified on

3.1.2. Reading Text Data

```
In [9]: # note the separator in this file
data_text = pd.read_csv("training_text.csv", sep="\\|\\|", engine="python", names=["ID", "TEXT"], skiprows=1)
print('Number of data points : ', data_text.shape[0])
print('Number of features : ', data_text.shape[1])
print('Features : ', data_text.columns.values)
data_text.head()
```

```
Number of data points : 3321
Number of features : 2
Features : ['ID' 'TEXT']
```

Out[9]:

	ID	TEXT
0	0	Cyclin-dependent kinases (CDKs) regulate a var...
1	1	Abstract Background Non-small cell lung canc...
2	2	Abstract Background Non-small cell lung canc...
3	3	Recent evidence has demonstrated that acquired...
4	4	Oncogenic mutations in the monomeric Casitas B...

```
In [10]: import nltk
nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to C:\Users\Himanshu
[nltk_data]   Pc\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

Out[10]: True

3.1.3. Preprocessing of text

```

In [11]: # Loading stop words from nltk library
stop_words = set(stopwords.words('english'))

def nlp_preprocessing(total_text, index, column):
    if type(total_text) is not int:
        string = ""
        # replace every special char with space
        total_text = re.sub('[^a-zA-Z0-9\n]', ' ', total_text)
        # replace multiple spaces with single space
        total_text = re.sub('\s+', ' ', total_text)
        # converting all the chars into lower-case.
        total_text = total_text.lower()

        for word in total_text.split():
            # if the word is a not a stop word then retain that word from the data
            if not word in stop_words:
                string += word + " "

        data_text[column][index] = string

```

```

In [12]: #text processing stage.
start_time = time.clock()
for index, row in data_text.iterrows():
    if type(row['TEXT']) is str:
        nlp_preprocessing(row['TEXT'], index, 'TEXT')
    else:
        print("there is no text description for id:",index)
print('Time took for preprocessing the text :',time.clock() - start_time, "seconds")

```

```

there is no text description for id: 1109
there is no text description for id: 1277
there is no text description for id: 1407
there is no text description for id: 1639
there is no text description for id: 2755
Time took for preprocessing the text : 131.04027170000006 seconds

```



```
In [13]: #merging both gene_variations and text data based on ID
result = pd.merge(data, data_text,on='ID', how='left')
result.head()
```

Out[13]:

	ID	Gene	Variation	Class	TEXT
0	0	FAM58A	Truncating Mutations	1	cyclin dependent kinases cdks regulate variety...
1	1	CBL	W802*	2	abstract background non small cell lung cancer...
2	2	CBL	Q249E	2	abstract background non small cell lung cancer...
3	3	CBL	N454D	3	recent evidence demonstrated acquired uniparen...
4	4	CBL	L399V	4	oncogenic mutations monomeric casitas b lineag...

```
In [14]: result[result.isnull().any(axis=1)]
```

Out[14]:

	ID	Gene	Variation	Class	TEXT
1109	1109	FANCA	S1088F	1	NaN
1277	1277	ARID5B	Truncating Mutations	1	NaN
1407	1407	FGFR3	K508M	6	NaN
1639	1639	FLT1	Amplification	6	NaN
2755	2755	BRAF	G596C	7	NaN

```
In [15]: result.loc[result['TEXT'].isnull(),'TEXT'] = result['Gene'] + ' '+result['Variation']
```

```
In [16]: result[result['ID']==1109]
```

Out[16]:

	ID	Gene	Variation	Class	TEXT
1109	1109	FANCA	S1088F	1	FANCA S1088F

3.1.4. Test, Train and Cross Validation Split

3.1.4.1. Splitting data into train, test and cross validation (64:20:16)

```
In [17]: y_true = result['Class'].values
result.Gene      = result.Gene.str.replace('\s+', '_')
result.Variation = result.Variation.str.replace('\s+', '_')

# split the data into test and train by maintaining same distribution of output variable 'y_true' [stratify=y_true]
X_train, test_df, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_size=0.2)
# split the train data into train and cross validation by maintaining same distribution of output variable 'y_train' [st
train_df, cv_df, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)
```

We split the data into train, test and cross validation data sets, preserving the ratio of class distribution in the original data set

```
In [18]: print('Number of data points in train data:', train_df.shape[0])
print('Number of data points in test data:', test_df.shape[0])
print('Number of data points in cross validation data:', cv_df.shape[0])
```

```
Number of data points in train data: 2124
Number of data points in test data: 665
Number of data points in cross validation data: 532
```

3.1.4.2. Distribution of y_i's in Train, Test and Cross Validation datasets

```

In [19]: # it returns a dict, keys as class labels and values as the number of data points in that class
train_class_distribution = train_df['Class'].value_counts().sortlevel()
test_class_distribution = test_df['Class'].value_counts().sortlevel()
cv_class_distribution = cv_df['Class'].value_counts().sortlevel()

my_colors = 'rgbkymc'
train_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in train data')
plt.grid()
plt.show()

# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':', train_class_distribution.values[i], '(', np.round((train_class_dist

print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

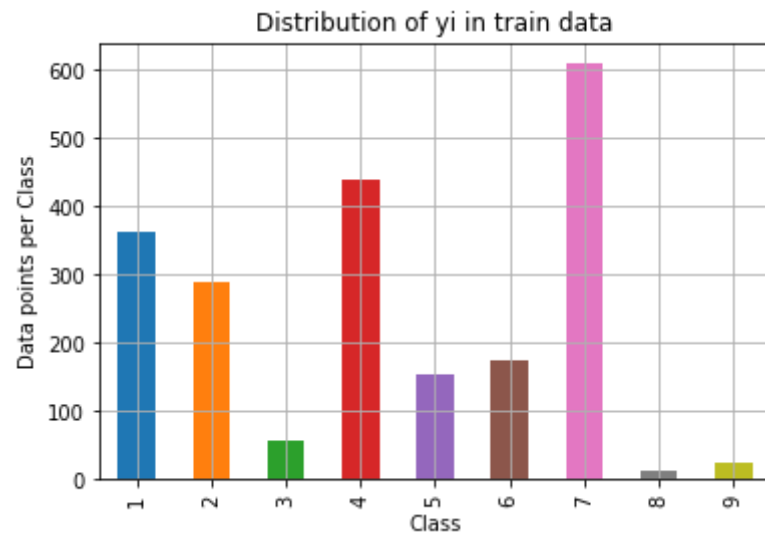
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':', test_class_distribution.values[i], '(', np.round((test_class_distri

print('-'*80)
my_colors = 'rgbkymc'
cv_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in cross validation data')

```

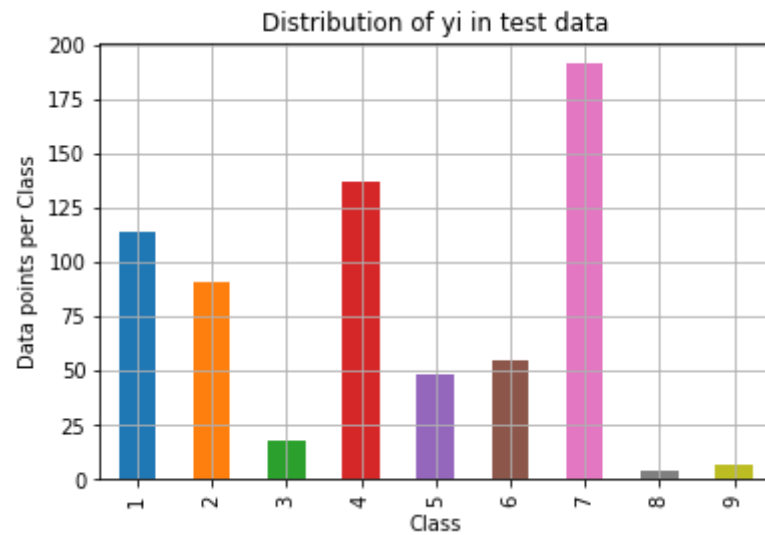
```
plt.grid()
plt.show()

# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ': ', cv_class_distribution.values[i], ' (', np.round((cv_class_distributi
```

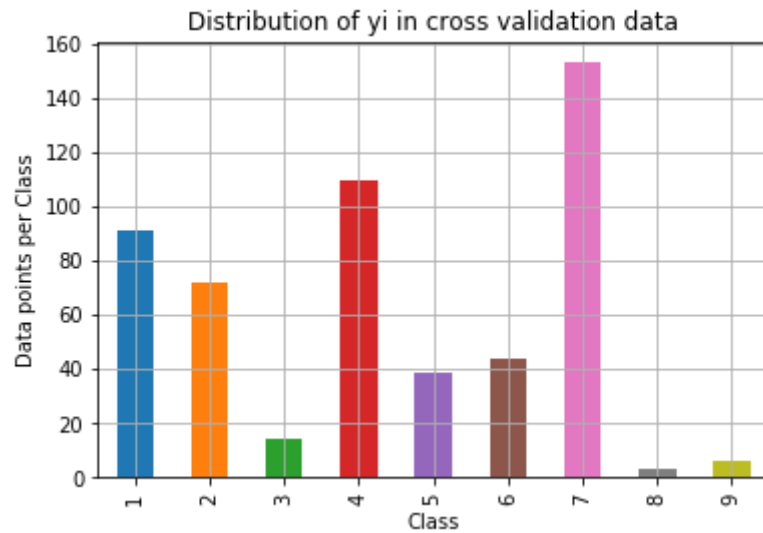


Number of data points in class 7 : 609 (28.672 %)
 Number of data points in class 4 : 439 (20.669 %)

Number of data points in class 1 : 363 (17.09 %)
Number of data points in class 2 : 289 (13.606 %)
Number of data points in class 6 : 176 (8.286 %)
Number of data points in class 5 : 155 (7.298 %)
Number of data points in class 3 : 57 (2.684 %)
Number of data points in class 9 : 24 (1.13 %)
Number of data points in class 8 : 12 (0.565 %)



Number of data points in class 7 : 191 (28.722 %)
Number of data points in class 4 : 137 (20.602 %)
Number of data points in class 1 : 114 (17.143 %)
Number of data points in class 2 : 91 (13.684 %)
Number of data points in class 6 : 55 (8.271 %)
Number of data points in class 5 : 48 (7.218 %)
Number of data points in class 3 : 18 (2.707 %)
Number of data points in class 9 : 7 (1.053 %)
Number of data points in class 8 : 4 (0.602 %)



Number of data points in class 7 : 153 (28.759 %)
Number of data points in class 4 : 110 (20.677 %)
Number of data points in class 1 : 91 (17.105 %)
Number of data points in class 2 : 72 (13.534 %)
Number of data points in class 6 : 44 (8.271 %)
Number of data points in class 5 : 39 (7.331 %)
Number of data points in class 3 : 14 (2.632 %)
Number of data points in class 9 : 6 (1.128 %)
Number of data points in class 8 : 3 (0.564 %)

3.2 Prediction using a 'Random' Model

In a 'Random' Model, we generate the NINE class probabilities randomly such that they sum to 1.


```

In [20]: # 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]]

    labels = [1,2,3,4,5,6,7,8,9]
    # representing A in heatmap format
    print("-"*20, "Confusion matrix", "-"*20)
    plt.figure(figsize=(20,7))
    sns.heatmap(C, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()

    print("-"*20, "Precision matrix (Columm Sum=1)", "-"*20)
    plt.figure(figsize=(20,7))

```



```
sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()

# representing B in heatmap format
print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
plt.figure(figsize=(20,7))
sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()
```

```

In [21]: # 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
test_data_len = test_df.shape[0]
cv_data_len = cv_df.shape[0]

# we create a output array that has exactly same size as the CV data
cv_predicted_y = np.zeros((cv_data_len,9))
for i in range(cv_data_len):
    rand_probs = np.random.rand(1,9)
    cv_predicted_y[i] = ((rand_probs/sum(sum(rand_probs))))[0])
print("Log loss on Cross Validation Data using Random Model",log_loss(y_cv,cv_predicted_y, eps=1e-15))

# Test-Set error.
#we create a output array that has exactly same as the test data
test_predicted_y = np.zeros((test_data_len,9))
for i in range(test_data_len):
    rand_probs = np.random.rand(1,9)
    test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs))))[0])
print("Log loss on Test Data using Random Model",log_loss(y_test,test_predicted_y, eps=1e-15))

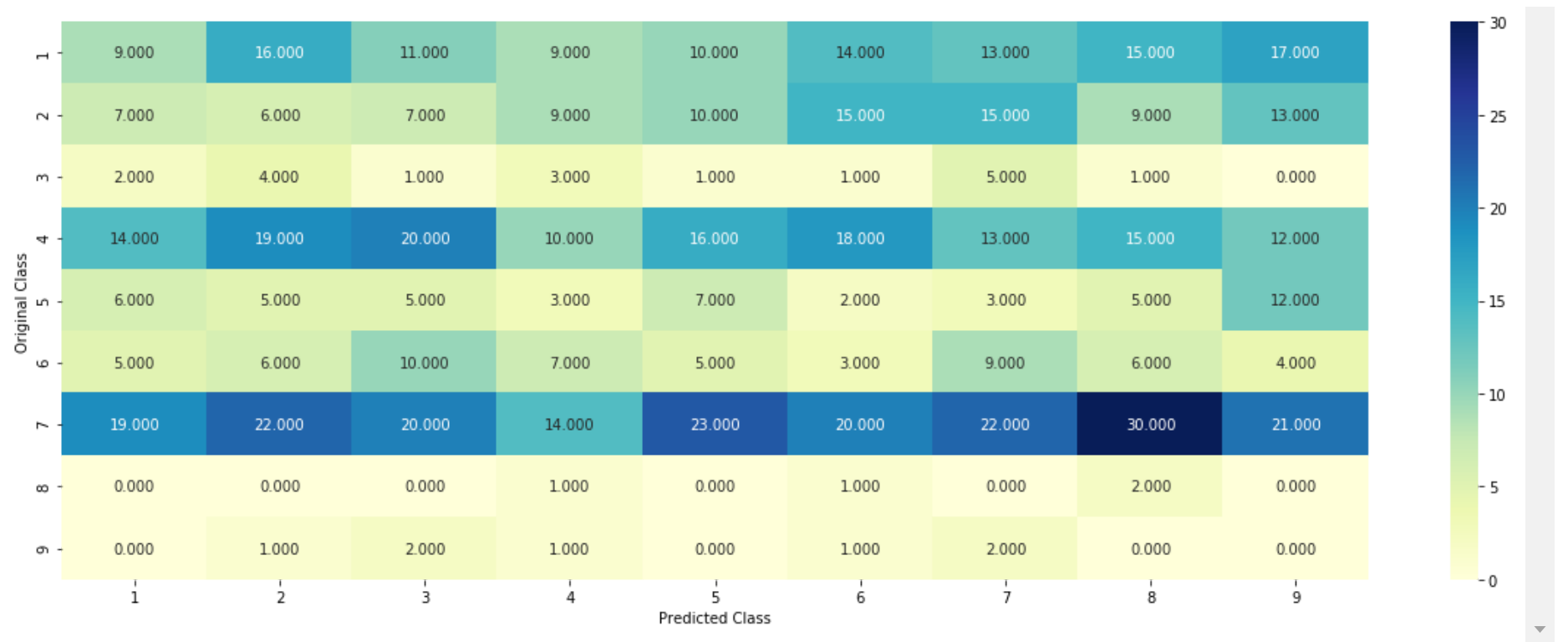
predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

```

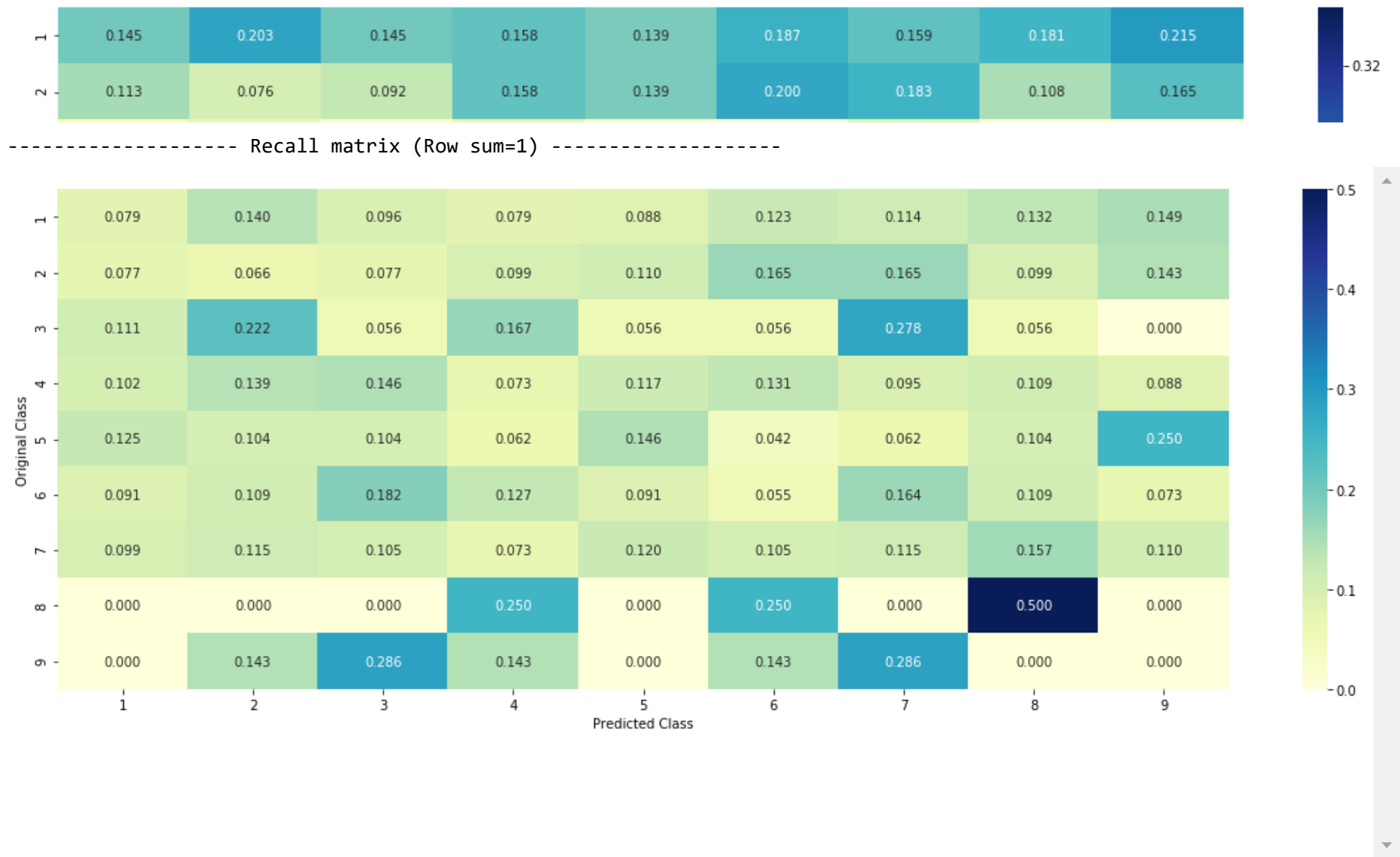
Log loss on Cross Validation Data using Random Model 2.4471631691399027

Log loss on Test Data using Random Model 2.4936180352741486

----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----



3.3 Univariate Analysis

```

In [22]: # code for response coding with Laplace smoothing.
# alpha : used for Laplace smoothing
# feature: ['gene', 'variation']
# df: ['train_df', 'test_df', 'cv_df']
# algorithm
# -----
# Consider all unique values and the number of occurrences of given feature in train data dataframe
# build a vector (1*9) , the first element = (number of times it occurred in class1 + 10*alpha / number of times it occurred)
# gv_dict is like a look up table, for every gene it stores a (1*9) representation of it
# for a value of feature in df:
# if it is in train data:
# we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
# if it is not there is train:
# we add [1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
# return 'gv_fea'
# -----

# get_gv_fea_dict: Get Gene variation Feature Dict
def get_gv_fea_dict(alpha, feature, df):
    # value_count: it contains a dict like
    # print(train_df['Gene'].value_counts())
    # output:
    #      {BRCA1      174
    #       TP53      106
    #       EGFR       86
    #       BRCA2       75
    #       PTEN       69
    #       KIT        61
    #       BRAF        60
    #       ERBB2       47
    #       PDGFRA      46
    #       ...}
    # print(train_df['Variation'].value_counts())
    # output:
    # {
    #   Truncating_Mutations      63
    #   Deletion                   43
    #   Amplification              43
    #   Fusions                    22
    #   Overexpression             3
    #   E17K                      3

```

```

# Q61L                                3
# S222D                                2
# P130S                                2
# ...
# }
value_count = train_df[feature].value_counts()

# gv_dict : Gene Variation Dict, which contains the probability array for each gene/variation
gv_dict = dict()

# denominator will contain the number of time that particular feature occurred in whole data
for i, denominator in value_count.items():
    # vec will contain (p(yi==1/Gi) probability of gene/variation belongs to particular class
    # vec is 9 dimensional vector
    vec = []
    for k in range(1,10):
        # print(train_df.loc[(train_df['Class']==1) & (train_df['Gene']=='BRCA1')])
        #
        #      ID      Gene      Variation      Class
        # 2470  2470  BRCA1      S1715C      1
        # 2486  2486  BRCA1      S1841R      1
        # 2614  2614  BRCA1      M1R      1
        # 2432  2432  BRCA1      L1657P      1
        # 2567  2567  BRCA1      T1685A      1
        # 2583  2583  BRCA1      E1660G      1
        # 2634  2634  BRCA1      W1718L      1
        # cls_cnt.shape[0] will return the number of rows

        cls_cnt = train_df.loc[(train_df['Class']==k) & (train_df[feature]==i)]

        # cls_cnt.shape[0](numerator) will contain the number of time that particular feature occurred in whole data
        vec.append((cls_cnt.shape[0] + alpha*10)/ (denominator + 90*alpha))

    # we are adding the gene/variation to the dict as key and vec as value
    gv_dict[i]=vec
return gv_dict

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    #
    #      {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.068181818181818177, 0.13636363636363635, 0.25, 0.193181
    #      'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.27040816326530615, 0.061224489795
    #      'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.068181818181818177, 0.068181818181818177, 0.062

```

```

#      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.078787878787878782, 0.1393939393
#      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.46540880503144655, 0.07547169811
#      'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295, 0.066225165562
#      'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333333333333334, 0.07333333333333334, 0.09333333333
#      ...
#      }
gv_dict = get_gv_fea_dict(alpha, feature, df)
# value_count is similar in get_gv_fea_dict
value_count = train_df[feature].value_counts()

# gv_fea: Gene_variation feature, it will contain the feature for each feature value in the data
gv_fea = []
# for every feature values in the given data frame we will check if it is there in the train data then we will add t
# if not we will add [1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9] to gv_fea
for index, row in df.iterrows():
    if row[feature] in dict(value_count).keys():
        gv_fea.append(gv_dict[row[feature]])
    else:
        gv_fea.append([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
#      gv_fea.append([-1,-1,-1,-1,-1,-1,-1,-1,-1])
return gv_fea

```

when we calculate the probability of a feature belongs to any particular class, we apply laplace smoothing

- $(\text{numerator} + 10 \cdot \alpha) / (\text{denominator} + 90 \cdot \alpha)$

3.2.1 Univariate Analysis on Gene Feature

Q1. Gene, What type of feature it is ?

Ans. Gene is a categorical variable

Q2. How many categories are there and How they are distributed?

```
In [23]: unique_genes = train_df['Gene'].value_counts()
print('Number of Unique Genes :', unique_genes.shape[0])
# the top 10 genes that occurred most
print(unique_genes.head(10))
```

Number of Unique Genes : 233

BRCA1 182

TP53 100

EGFR 88

BRCA2 85

PTEN 77

BRAF 60

KIT 54

ALK 48

ERBB2 42

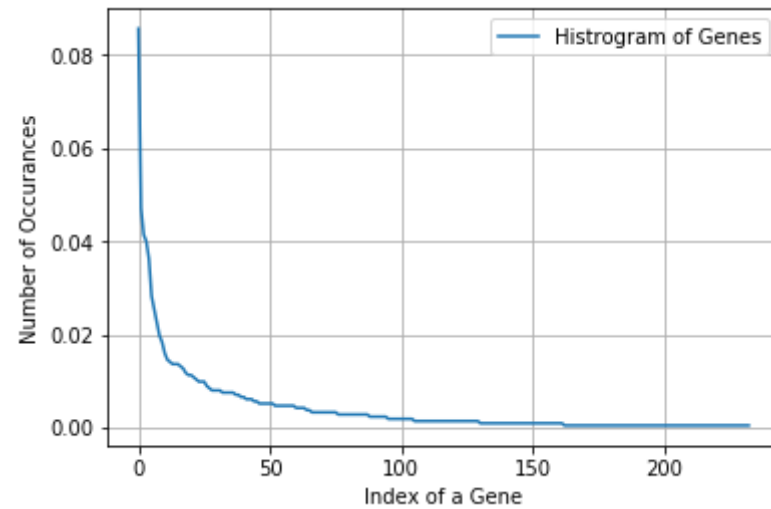
PDGFRA 39

Name: Gene, dtype: int64

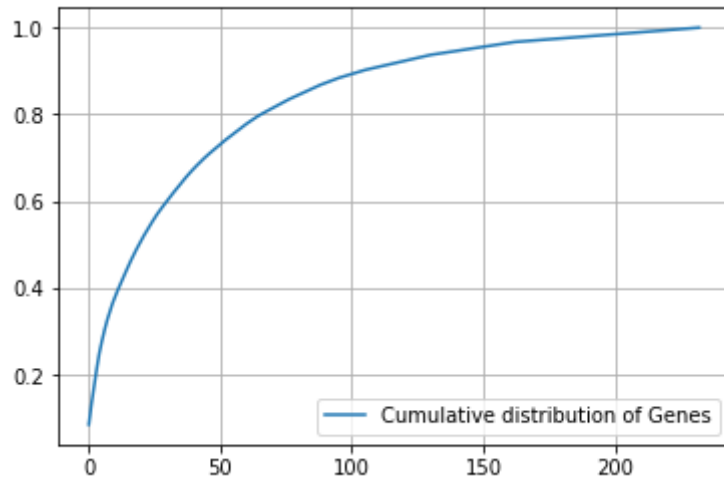
```
In [24]: print("Ans: There are", unique_genes.shape[0] , "different categories of genes in the train data, and they are distributed
```

Ans: There are 233 different categories of genes in the train data, and they are distributed as follows


```
In [25]: s = sum(unique_genes.values);  
h = unique_genes.values/s;  
plt.plot(h, label="Histogram of Genes")  
plt.xlabel('Index of a Gene')  
plt.ylabel('Number of Occurances')  
plt.legend()  
plt.grid()  
plt.show()
```



```
In [26]: c = np.cumsum(h)
plt.plot(c,label='Cumulative distribution of Genes')
plt.grid()
plt.legend()
plt.show()
```



Q3. How to featurize this Gene feature ?

Ans. there are two ways we can featurize this variable check out this video: <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

1. One hot Encoding
2. Response coding

We will choose the appropriate featurization based on the ML model we use. For this problem of multi-class classification with categorical features, one-hot encoding is better for Logistic regression while response coding is better for Random Forests.

```
In [27]: #response-coding of the Gene feature  
# alpha is used for Laplace smoothing  
alpha = 1  
# train gene feature  
train_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", train_df))  
# test gene feature  
test_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", test_df))  
# cross validation gene feature  
cv_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", cv_df))
```

```
In [28]: print("train_gene_feature_responseCoding is converted feature using response coding method. The shape of gene feature:",  
train_gene_feature_responseCoding is converted feature using response coding method. The shape of gene feature: (2124,  
9)
```

```
In [29]: # one-hot encoding of Gene feature.  
from sklearn.feature_extraction.text import TfidfVectorizer  
gene_vectorizer = TfidfVectorizer()  
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])  
test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])  
cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

```
In [30]: train_df['Gene'].head()
```

```
Out[30]: 1533      ALK  
195      EGFR  
1099     BAP1  
1046     TSC2  
62      PTPRT  
Name: Gene, dtype: object
```

```
In [31]: gene_vectorizer.get_feature_names()
```

```
Out[31]: ['abl1',  
          'acvr1',  
          'ago2',  
          'akt1',  
          'akt2',  
          'akt3',  
          'alk',  
          'apc',  
          'ar',  
          'araf',  
          'arid1a',  
          'arid1b',  
          'arid2',  
          'arid5b',  
          'asx11',  
          'asx12',  
          'atm',  
          'atrx',  
          'aurka',  
          ...]
```

```
In [32]: print("train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature:",
```

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 233)
```

Q4. How good is this gene feature in predicting y_i ?

There are many ways to estimate how good a feature is, in predicting y_i . One of the good methods is to build a proper ML model using just this feature. In this case, we will build a logistic regression model using only Gene feature (one hot encoded) to predict y_i .

```

In [33]: alpha = [10 ** x for x in range(-5, 1)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier
# -----
# 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:
#-----

cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_gene_feature_onehotCoding, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_gene_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)

```

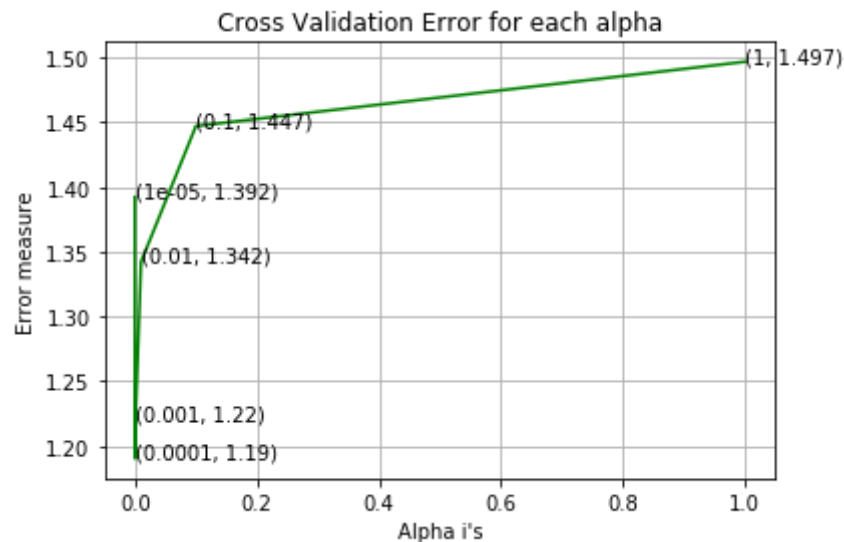
```

clf.fit(train_gene_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_gene_feature_onehotCoding, y_train)

predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_))

```

For values of alpha = 1e-05 The log loss is: 1.3923613423027779
 For values of alpha = 0.0001 The log loss is: 1.1903838279158911
 For values of alpha = 0.001 The log loss is: 1.2199799972987682
 For values of alpha = 0.01 The log loss is: 1.3421815296588195
 For values of alpha = 0.1 The log loss is: 1.44721331408222
 For values of alpha = 1 The log loss is: 1.4969154954899642



For values of best alpha = 0.0001 The train log loss is: 1.0456368485897742
For values of best alpha = 0.0001 The cross validation log loss is: 1.1903838279158911
For values of best alpha = 0.0001 The test log loss is: 1.224581817827015

Q5. Is the Gene feature stable across all the data sets (Test, Train, Cross validation)?

Ans. Yes, it is. Otherwise, the CV and Test errors would be significantly more than train error.

```
In [34]: print("Q6. How many data points in Test and CV datasets are covered by the ", unique_genes.shape[0], " genes in train da

test_coverage=test_df[test_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]
cv_coverage=cv_df[cv_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]

print('Ans\n1. In test data',test_coverage, 'out of',test_df.shape[0], ":",(test_coverage/test_df.shape[0])*100)
print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0]," :", (cv_coverage/cv_df.shape[0])*100)
```

Q6. How many data points in Test and CV datasets are covered by the 233 genes in train dataset?

Ans

1. In test data 644 out of 665 : 96.84210526315789
2. In cross validation data 514 out of 532 : 96.61654135338345

3.2.2 Univariate Analysis on Variation Feature

Q7. Variation, What type of feature is it ?

Ans. Variation is a categorical variable

Q8. How many categories are there?

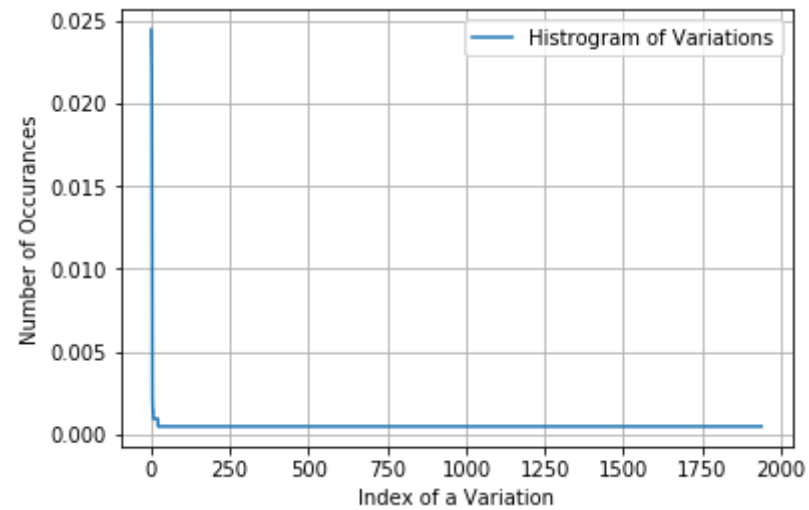
```
In [35]: unique_variations = train_df['Variation'].value_counts()
print('Number of Unique Variations :', unique_variations.shape[0])
# the top 10 variations that occurred most
print(unique_variations.head(10))
```

```
Number of Unique Variations : 1939
Truncating_Mutations      52
Amplification              48
Deletion                  46
Fusions                   20
Overexpression             5
E17K                      3
Q61R                      3
Y64A                      2
ETV6-NTRK3_Fusion         2
Q61L                      2
Name: Variation, dtype: int64
```

```
In [36]: print("Ans: There are", unique_variations.shape[0] , "different categories of variations in the train data, and they are
```

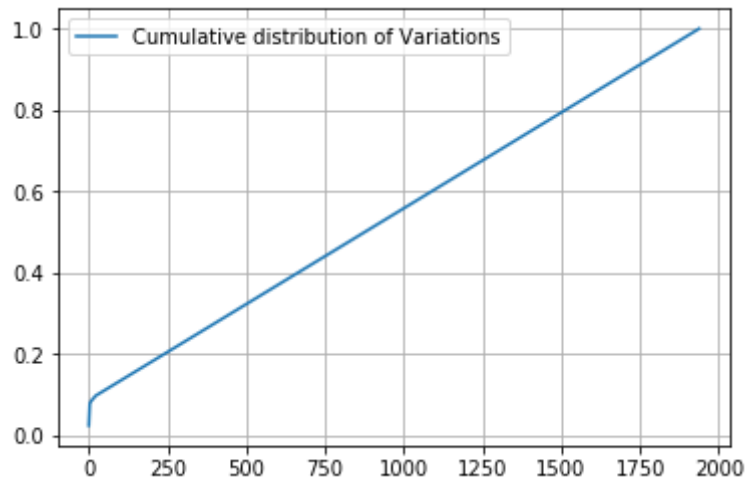
Ans: There are 1939 different categories of variations in the train data, and they are distributed as follows


```
In [37]: s = sum(unique_variations.values);  
h = unique_variations.values/s;  
plt.plot(h, label="Histogram of Variations")  
plt.xlabel('Index of a Variation')  
plt.ylabel('Number of Occurances')  
plt.legend()  
plt.grid()  
plt.show()
```



```
In [38]: c = np.cumsum(h)
print(c)
plt.plot(c,label='Cumulative distribution of Variations')
plt.grid()
plt.legend()
plt.show()

[0.02448211 0.04708098 0.06873823 ... 0.99905838 0.99952919 1.          ]
```



Q9. How to featurize this Variation feature ?

Ans. There are two ways we can featurize this variable check out this video: <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

1. One hot Encoding
2. Response coding

We will be using both these methods to featurize the Variation Feature

```
In [39]: # alpha is used for laplace smoothing
alpha = 1
# train gene feature
train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", train_df))
# test gene feature
test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", test_df))
# cross validation gene feature
cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", cv_df))
```

```
In [40]: print("train_variation_feature_responseCoding is a converted feature using the response coding method. The shape of Vari
```

train_variation_feature_responseCoding is a converted feature using the response coding method. The shape of Variation feature: (2124, 9)

```
In [41]: # one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer()
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(train_df['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(test_df['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(cv_df['Variation'])
```

```
In [42]: print("train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape of Varia
```

train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape of Variation feature: (2124, 1967)

Q10. How good is this Variation feature in predicting y_i ?

Let's build a model just like the earlier!

```
In [43]: alpha = [10 ** x for x in range(-5, 1)]
```

```
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier
# -----
# 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:
#-----

cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_variation_feature_onehotCoding, y_train)

    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_variation_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)

    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

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

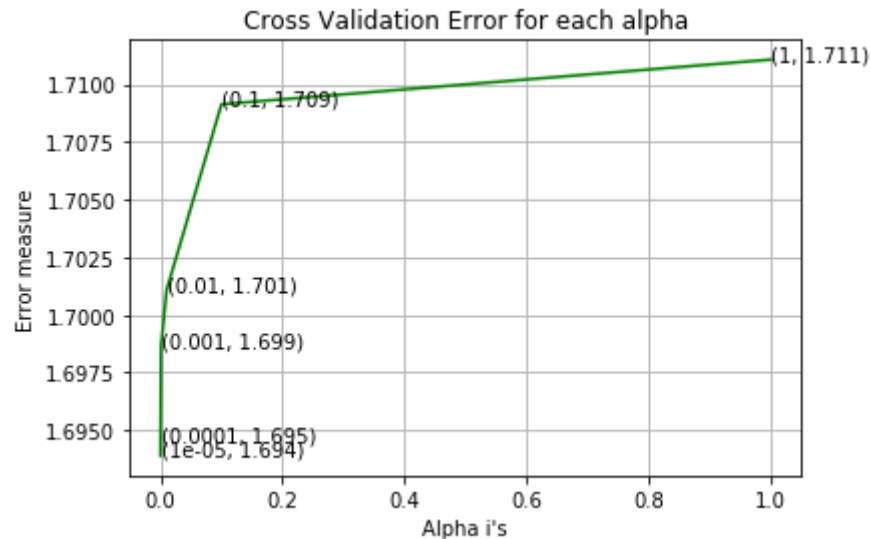
```

best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_variation_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_variation_feature_onehotCoding, y_train)

predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, la
predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.c

```

For values of alpha = 1e-05 The log loss is: 1.6939017334679705
 For values of alpha = 0.0001 The log loss is: 1.6945466116539385
 For values of alpha = 0.001 The log loss is: 1.6986274793368845
 For values of alpha = 0.01 The log loss is: 1.7010801015821841
 For values of alpha = 0.1 The log loss is: 1.7091356951616932
 For values of alpha = 1 The log loss is: 1.7110732925280796



For values of best alpha = 1e-05 The train log loss is: 0.6885760806303544
For values of best alpha = 1e-05 The cross validation log loss is: 1.6939017334679705
For values of best alpha = 1e-05 The test log loss is: 1.692557742855498

Q11. Is the Variation feature stable across all the data sets (Test, Train, Cross validation)?

Ans. Not sure! But lets be very sure using the below analysis.

```
In [44]: print("Q12. How many data points are covered by total ", unique_variations.shape[0], " genes in test and cross validation data sets?  
test_coverage=test_df[test_df['Variation'].isin(list(set(train_df['Variation'])))].shape[0]  
cv_coverage=cv_df[cv_df['Variation'].isin(list(set(train_df['Variation'])))].shape[0]  
print('Ans\n1. In test data',test_coverage, 'out of',test_df.shape[0], ":",(test_coverage/test_df.shape[0])*100)  
print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0],":", (cv_coverage/cv_df.shape[0])*100)
```

Q12. How many data points are covered by total 1939 genes in test and cross validation data sets?

Ans

1. In test data 81 out of 665 : 12.180451127819548
2. In cross validation data 55 out of 532 : 10.338345864661653

3.2.3 Univariate Analysis on Text Feature

1. How many unique words are present in train data?
2. How are word frequencies distributed?
3. How to featurize text field?
4. Is the text feature useful in predicting y_i ?
5. Is the text feature stable across train, test and CV datasets?

```
In [45]: # cls_text is a data frame
# for every row in data fram consider the 'TEXT'
# split the words by space
# make a dict with those words
# increment its count whenever we see that word
```

```
def extract_dictionary_paddle(cls_text):
    dictionary = defaultdict(int)
    for index, row in cls_text.iterrows():
        for word in row['TEXT'].split():
            dictionary[word] +=1
    return dictionary
```

```
In [46]: import math
#https://stackoverflow.com/a/1602964
def get_text_responsecoding(df):
    text_feature_responseCoding = np.zeros((df.shape[0],9))
    for i in range(0,9):
        row_index = 0
        for index, row in df.iterrows():
            sum_prob = 0
            for word in row['TEXT'].split():
                sum_prob += math.log(((dict_list[i].get(word,0)+10 )/(total_dict.get(word,0)+90)))
            text_feature_responseCoding[row_index][i] = math.exp(sum_prob/len(row['TEXT'].split()))
            row_index += 1
    return text_feature_responseCoding
```

```

In [47]: # building a TfidfVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer(min_df=3)
train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])
# getting all the feature names (words)
train_text_features = text_vectorizer.get_feature_names()

# train_text_feature_onehotCoding.sum(axis=0).A1 will sum every row and returns (1*number of features) vector
train_text_fea_counts = train_text_feature_onehotCoding.sum(axis=0).A1

# zip(list(text_features),text_fea_counts) will zip a word with its number of times it occurred
text_fea_dict = dict(zip(list(train_text_features),train_text_fea_counts))

print("Total number of unique words in train data :", len(train_text_features))

```

Total number of unique words in train data : 53251

```

In [48]: dict_list = []
# dict_list =[] contains 9 dictionaries each corresponds to a class
for i in range(1,10):
    cls_text = train_df[train_df['Class']==i]
    # build a word dict based on the words in that class
    dict_list.append(extract_dictionary_paddle(cls_text))
    # append it to dict_list

# dict_list[i] is build on i'th class text data
# total_dict is build on whole training text data
total_dict = extract_dictionary_paddle(train_df)

confuse_array = []
for i in train_text_features:
    ratios = []
    max_val = -1
    for j in range(0,9):
        ratios.append((dict_list[j][i]+10)/(total_dict[i]+90))
    confuse_array.append(ratios)
confuse_array = np.array(confuse_array)

```



```
In [49]: #response coding of text features
train_text_feature_responseCoding = get_text_responsecoding(train_df)
test_text_feature_responseCoding = get_text_responsecoding(test_df)
cv_text_feature_responseCoding = get_text_responsecoding(cv_df)
```

```
In [50]: # https://stackoverflow.com/a/16202486
# we convert each row values such that they sum to 1
train_text_feature_responseCoding = (train_text_feature_responseCoding.T/train_text_feature_responseCoding.sum(axis=1)).T
test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.sum(axis=1)).T
cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=1)).T
```

```
In [51]: # don't forget to normalize every feature
train_text_feature_onehotCoding = normalize(train_text_feature_onehotCoding, axis=0)

# we use the same vectorizer that was trained on train data
test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
# don't forget to normalize every feature
test_text_feature_onehotCoding = normalize(test_text_feature_onehotCoding, axis=0)

# we use the same vectorizer that was trained on train data
cv_text_feature_onehotCoding = text_vectorizer.transform(cv_df['TEXT'])
# don't forget to normalize every feature
cv_text_feature_onehotCoding = normalize(cv_text_feature_onehotCoding, axis=0)
```

```
In [52]: #https://stackoverflow.com/a/2258273/4084039
sorted_text_fea_dict = dict(sorted(text_fea_dict.items(), key=lambda x: x[1] , reverse=True))
sorted_text_occur = np.array(list(sorted_text_fea_dict.values()))
```

```
# Number of words for a given frequency.
print(Counter(sorted_text_occur))
```

Counter({0.006778898601341044: 323, 0.020247838411974747: 315, 0.02599478516689304: 230, 0.08075811055028549: 198, 0.027280172470686686: 187, 0.3334301102596617: 154, 0.02640842829903654: 139, 0.2630218486255754: 136, 0.04960262222527532: 124, 0.019515328038108123: 116, 0.0412806823829328: 114, 0.06048422297980386: 107, 0.04378463267343804: 104, 0.045588954419229194: 102, 0.02750139999689016: 102, 0.08095835642677447: 101, 0.03674411640037657: 96, 0.02851548578881293: 93, 0.018521470977529174: 93, 0.02398856737357453: 91, 0.11162546445913982: 87, 0.02611027015939002: 83, 0.024553061201490384: 83, 0.02028362117072389: 81, 0.025277394329811062: 80, 0.034063944263414886: 74, 0.057655801558905165: 69, 0.024105894119439652: 69, 0.03616623462092275: 68, 0.04625838386168363: 66, 0.011538728412169793: 66, 0.045117349314266406: 65, 0.060132077212428296: 64, 0.015954345803845527: 64, 0.015278394299085958: 64, 0.010666328316126863: 61, 0.012165072426700747: 60, 0.14126522875485378: 59, 0.019000116954134977: 59, 0.03172043730718365: 58, 0.010813205234015433: 58, 0.05458298374528034: 57, 0.048609927946438154: 57, 0.01688600518636504: 57, 0.11956931986914933: 56, 0.011925728985283077: 56, 0.01787903546766058: 55, 0.019771879025761604: 54, 0.02840434294523626: 53, 0.01743294755039674: 53, 0.2613747186542437: 51, 0.02575847276845654: 51, 0.021753937238536733: 51, 0.0841132378555108: 50, 0.0297773168391361: 50, 0.014642724061217699: 50, 0.014016052963012986: 49, 0.013829805270869485: 49, 0.05784485731202661: 48, 0.05612446563881204: 48, 0.025315712502879048: 48, 0.02091353584940624: 48, 0.008867545707589833: 48, 0.15714026344475393: 47, 0.04433849180233074: 47, 0.035556729811749455: 47, 0.021637696637199593: 47, 0.01889794924424844: 47, 0.009007577359912054: 47, 0.05181780289785548: 46, 0.08619869539787903: 45, 0.0361033006844828: 45, 0.015570921062912246: 45, 0.10580036735157682: 44, 0.034298702064023134: 44, 0.03814334729230173: 43, 0.03427750446994655: 43, 0.025468930236770937: 43, 0.019047278413806226: 43, 0.0643273054192798: 42, 0.05283188634163555: 42, 0.07346754590970934: 41, 0.03029818838477491: 41, 0.05281685659807308: 40, 0.036304555340950624: 39, 0.03402355822715871: 39, 0.03036840350769085

```

In [54]: # Train a Logistic regression+Calibration model using text features which are on-hot encoded
alpha = [10 ** x for x in range(-5, 1)]

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier
# -----
# 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:
#-----

cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_text_feature_onehotCoding, y_train)

    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_text_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

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

```

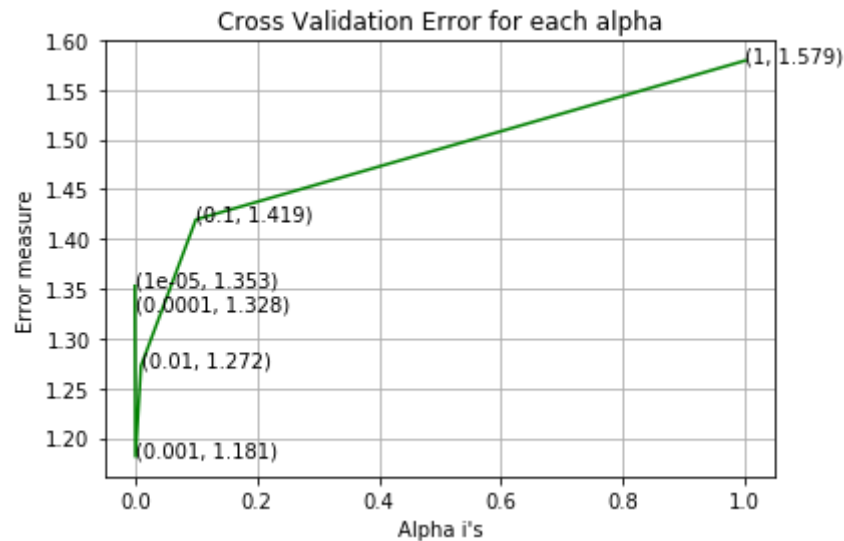
```

best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_text_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_text_feature_onehotCoding, y_train)

predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, la
predict_y = sig_clf.predict_proba(test_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.c

```

For values of alpha = 1e-05 The log loss is: 1.352687725086972
 For values of alpha = 0.0001 The log loss is: 1.328358163551927
 For values of alpha = 0.001 The log loss is: 1.1810889220250345
 For values of alpha = 0.01 The log loss is: 1.272310193592132
 For values of alpha = 0.1 The log loss is: 1.4194282352345788
 For values of alpha = 1 The log loss is: 1.5792073015809702



For values of best alpha = 0.001 The train log loss is: 0.7026859588234048
 For values of best alpha = 0.001 The cross validation log loss is: 1.1810889220250345
 For values of best alpha = 0.001 The test log loss is: 1.1825601373352446

Q. Is the Text feature stable across all the data sets (Test, Train, Cross validation)?

Ans. Yes, it seems like!

```
In [55]: def get_intersec_text(df):
          df_text_vec = TfidfVectorizer(min_df=3)
          df_text_fea = df_text_vec.fit_transform(df['TEXT'])
          df_text_features = df_text_vec.get_feature_names()

          df_text_fea_counts = df_text_fea.sum(axis=0).A1
          df_text_fea_dict = dict(zip(list(df_text_features), df_text_fea_counts))
          len1 = len(set(df_text_features))
          len2 = len(set(train_text_features) & set(df_text_features))
          return len1, len2
```

```
In [56]: len1, len2 = get_intersec_text(test_df)
          print(np.round((len2/len1)*100, 3), "% of word of test data appeared in train data")
          len1, len2 = get_intersec_text(cv_df)
          print(np.round((len2/len1)*100, 3), "% of word of Cross Validation appeared in train data")
```

```
97.571 % of word of test data appeared in train data
98.244 % of word of Cross Validation appeared in train data
```

4. Machine Learning Models

In [57]: *#Data preparation for ML models.*

#Misc. functions for ML models

```
def predict_and_plot_confusion_matrix(train_x, train_y, test_x, test_y, clf):  
    clf.fit(train_x, train_y)  
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")  
    sig_clf.fit(train_x, train_y)  
    pred_y = sig_clf.predict(test_x)  
  
    # for calculating log_loss we will provide the array of probabilities belongs to each class  
    print("Log loss :", log_loss(test_y, sig_clf.predict_proba(test_x)))  
    # calculating the number of data points that are misclassified  
    print("Number of mis-classified points :", np.count_nonzero((pred_y - test_y))/test_y.shape[0])  
    plot_confusion_matrix(test_y, pred_y)
```

In [58]: **def** report_log_loss(train_x, train_y, test_x, test_y, clf):

```
    clf.fit(train_x, train_y)  
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")  
    sig_clf.fit(train_x, train_y)  
    sig_clf_probs = sig_clf.predict_proba(test_x)  
    return log_loss(test_y, sig_clf_probs, eps=1e-15)
```

```

In [59]: # this function will be used just for naive bayes
# for the given indices, we will print the name of the features
# and we will check whether the feature present in the test point text or not
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = TfidfVectorizer()
    var_count_vec = TfidfVectorizer()
    text_count_vec = TfidfVectorizer(min_df=3)

    gene_vec = gene_count_vec.fit(train_df['Gene'])
    var_vec = var_count_vec.fit(train_df['Variation'])
    text_vec = text_count_vec.fit(train_df['TEXT'])

    fea1_len = len(gene_vec.get_feature_names())
    fea2_len = len(var_count_vec.get_feature_names())

    word_present = 0
    for i,v in enumerate(indices):
        if (v < fea1_len):
            word = gene_vec.get_feature_names()[v]
            yes_no = True if word == gene else False
            if yes_no:
                word_present += 1
                print(i, "Gene feature [{}] present in test data point [{}]" .format(word,yes_no))
        elif (v < fea1_len+fea2_len):
            word = var_vec.get_feature_names()[v-(fea1_len)]
            yes_no = True if word == var else False
            if yes_no:
                word_present += 1
                print(i, "variation feature [{}] present in test data point [{}]" .format(word,yes_no))
        else:
            word = text_vec.get_feature_names()[v-(fea1_len+fea2_len)]
            yes_no = True if word in text.split() else False
            if yes_no:
                word_present += 1
                print(i, "Text feature [{}] present in test data point [{}]" .format(word,yes_no))

    print("Out of the top ",no_features," features ", word_present, "are present in query point")

```

Stacking the three types of features

```
In [60]: # merging gene, variance and text features

# building train, test and cross validation data sets
# a = [[1, 2],
#      [3, 4]]
# b = [[4, 5],
#      [6, 7]]
# hstack(a, b) = [[1, 2, 4, 5],
#                [ 3, 4, 6, 7]]

train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding, train_variation_feature_onehotCoding))
test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding, test_variation_feature_onehotCoding))
cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding, cv_variation_feature_onehotCoding))

train_x_onehotCoding = hstack((train_gene_var_onehotCoding, train_text_feature_onehotCoding)).tocsr()
train_y = np.array(list(train_df['Class']))

test_x_onehotCoding = hstack((test_gene_var_onehotCoding, test_text_feature_onehotCoding)).tocsr()
test_y = np.array(list(test_df['Class']))

cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(cv_df['Class']))

train_gene_var_responseCoding = np.hstack((train_gene_feature_responseCoding, train_variation_feature_responseCoding))
test_gene_var_responseCoding = np.hstack((test_gene_feature_responseCoding, test_variation_feature_responseCoding))
cv_gene_var_responseCoding = np.hstack((cv_gene_feature_responseCoding, cv_variation_feature_responseCoding))

train_x_responseCoding = np.hstack((train_gene_var_responseCoding, train_text_feature_responseCoding))
test_x_responseCoding = np.hstack((test_gene_var_responseCoding, test_text_feature_responseCoding))
cv_x_responseCoding = np.hstack((cv_gene_var_responseCoding, cv_text_feature_responseCoding))
```



```
In [61]: print("One hot encoding features :")
print("(number of data points * number of features) in train data = ", train_x_onehotCoding.shape)
print("(number of data points * number of features) in test data = ", test_x_onehotCoding.shape)
print("(number of data points * number of features) in cross validation data =", cv_x_onehotCoding.shape)
```

```
One hot encoding features :
(number of data points * number of features) in train data = (2124, 55451)
(number of data points * number of features) in test data = (665, 55451)
(number of data points * number of features) in cross validation data = (532, 55451)
```

```
In [62]: print(" Response encoding features :")
print("(number of data points * number of features) in train data = ", train_x_responseCoding.shape)
print("(number of data points * number of features) in test data = ", test_x_responseCoding.shape)
print("(number of data points * number of features) in cross validation data =", cv_x_responseCoding.shape)
```

```
Response encoding features :
(number of data points * number of features) in train data = (2124, 27)
(number of data points * number of features) in test data = (665, 27)
(number of data points * number of features) in cross validation data = (532, 27)
```

4.1. Base Line Model

4.1.1. Naive Bayes

4.1.1.1. Hyper parameter tuning

```

In [63]: # find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/sklearn.naive\_bayes
# -----
# default paramters
# sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)

# some of methods of MultinomialNB()
# fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
# predict(X) Perform classification on an array of test vectors X.
# predict_log_proba(X) Return log-probability estimates for the test vector X.
# -----
# video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.Ca
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
# -----
# video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

alpha = [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100, 1000]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = MultinomialNB(alpha=i)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabillites we use log-probability estimates
    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

```

```

fig, ax = plt.subplots()
ax.plot(np.log10(alpha), cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],str(txt)), (np.log10(alpha[i]),cv_log_error_array[i]))
plt.grid()
plt.xticks(np.log10(alpha))
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(cv_log_error_array)
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, la
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.c

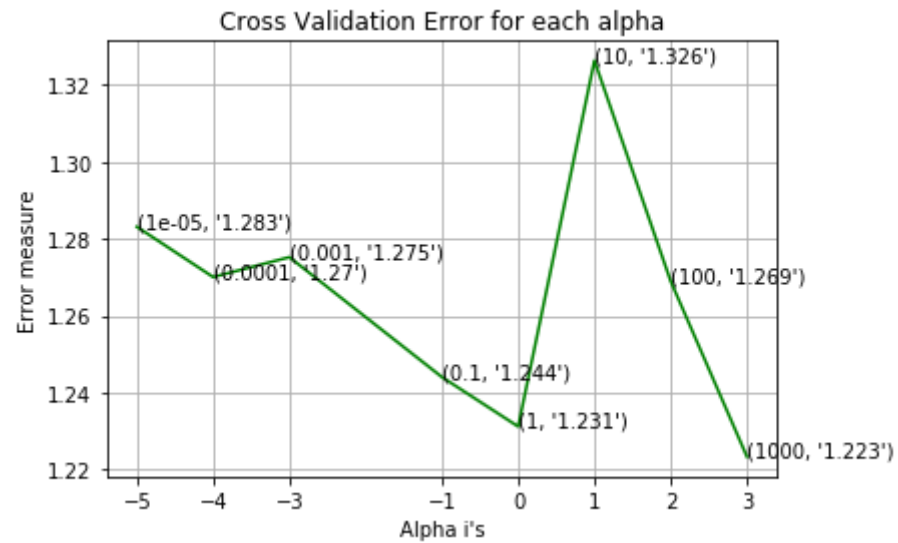
```

```

for alpha = 1e-05
Log Loss : 1.283137886452169
for alpha = 0.0001
Log Loss : 1.2700567355551982
for alpha = 0.001
Log Loss : 1.275229415079927
for alpha = 0.1
Log Loss : 1.244032669530854
for alpha = 1
Log Loss : 1.231145148707748
for alpha = 10
Log Loss : 1.3263855461367018
for alpha = 100

```

Log Loss : 1.2689838379474125
for alpha = 1000
Log Loss : 1.2233119279860467



For values of best alpha = 1000 The train log loss is: 0.9442994720988458
For values of best alpha = 1000 The cross validation log loss is: 1.2233119279860467
For values of best alpha = 1000 The test log loss is: 1.233771684693973

4.1.1.2. Testing the model with best hyper paramters

```

In [64]: # find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/sklearn.naive\_bayes.MultinomialNB.html
# -----
# default paramters
# sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)

# some of methods of MultinomialNB()
# fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
# predict(X) Perform classification on an array of test vectors X.
# predict_log_proba(X) Return log-probability estimates for the test vector X.
# -----
# video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
# -----

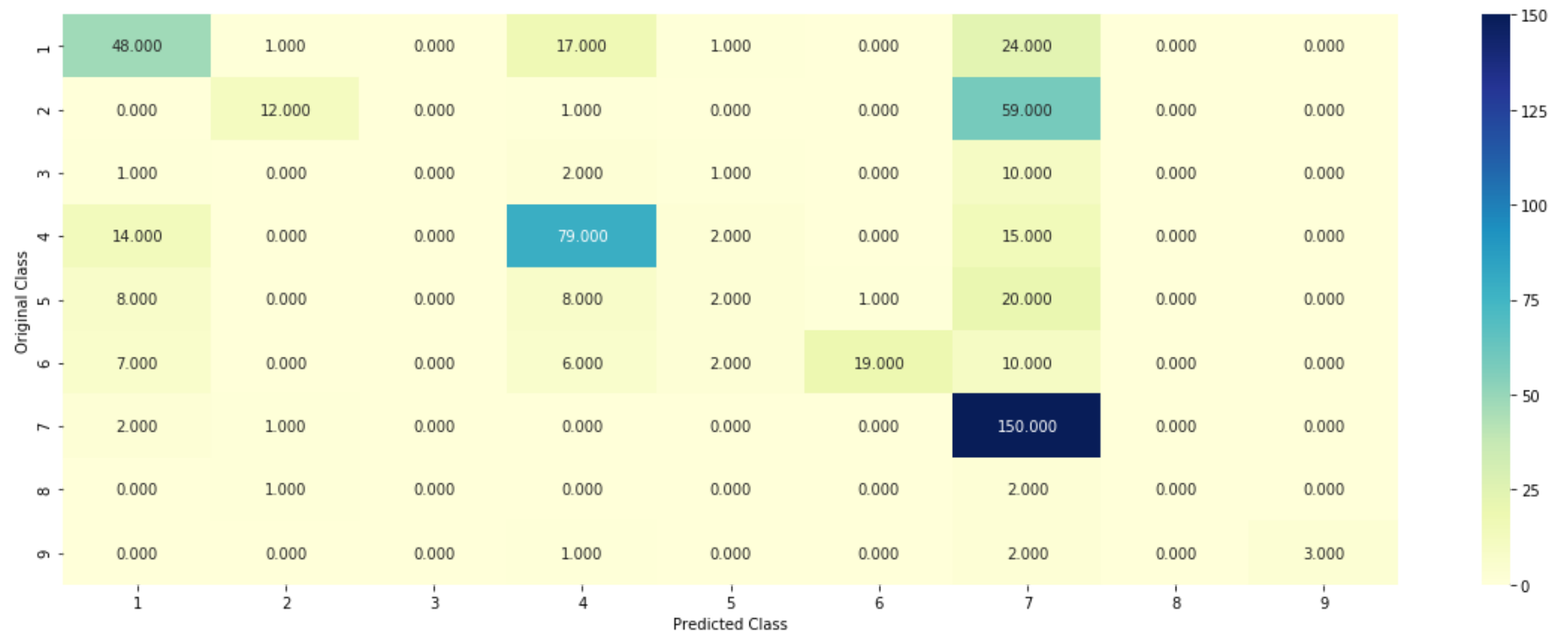
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
# to avoid rounding error while multiplying probabilities we use log-probability estimates
print("Log Loss :", log_loss(cv_y, sig_clf_probs))
print("Number of missclassified point :", np.count_nonzero((sig_clf.predict(cv_x_onehotCoding) - cv_y))/cv_y.shape[0])
plot_confusion_matrix(cv_y, sig_clf.predict(cv_x_onehotCoding.toarray()))

```

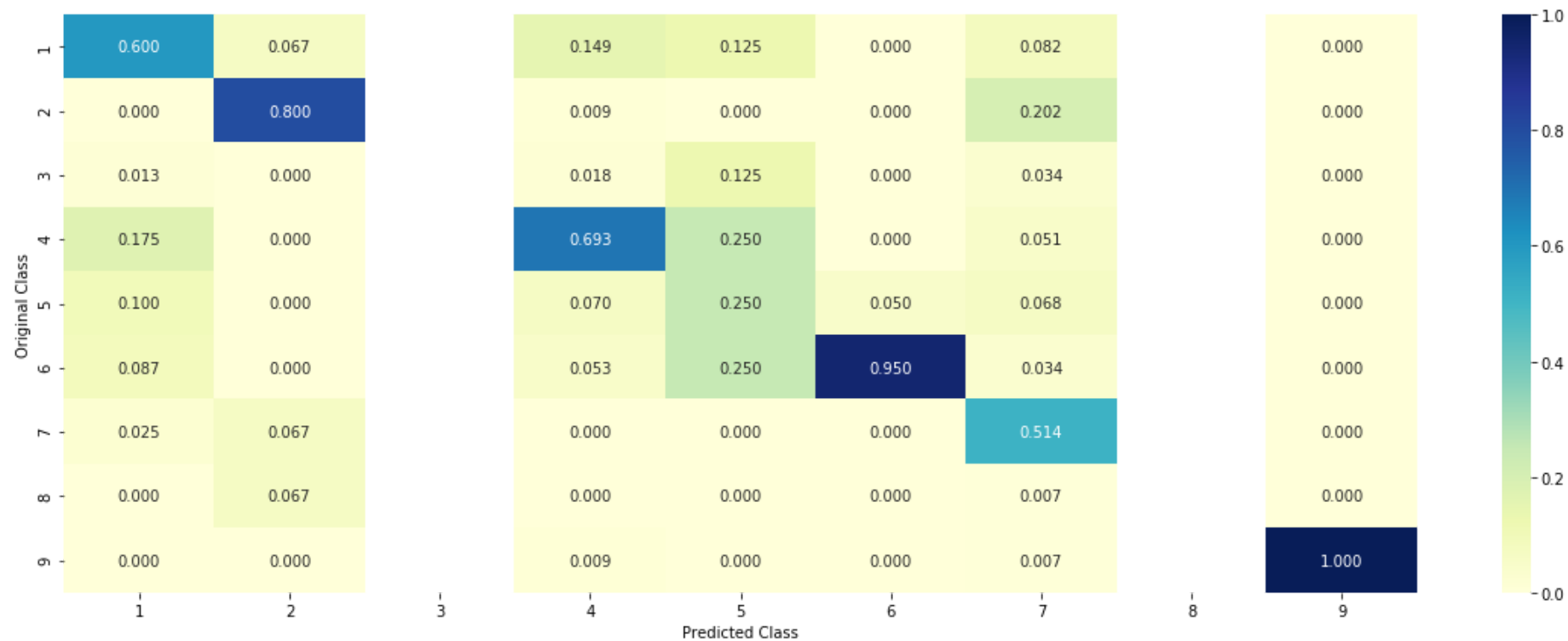
Log Loss : 1.2233119279860467

Number of missclassified point : 0.4116541353383459

----- Confusion matrix -----

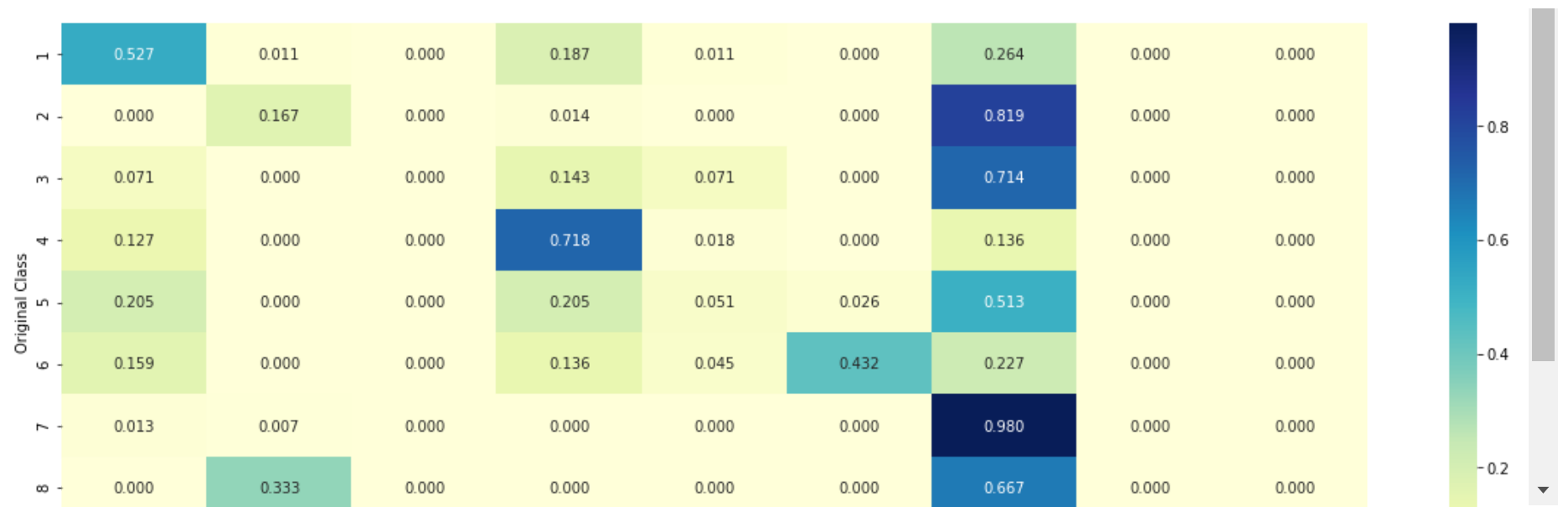


----- Precision matrix (Column Sum=1) -----



----- Recall matrix (Row sum=1) -----





4.1.1.3. Feature Importance, Correctly classified point


```
In [65]: test_point_index = 1
no_feature = 100
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices=np.argsort(abs(-clf.coef_))[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['
```

Predicted Class : 7

Predicted Class Probabilities: [[4.710e-02 3.012e-01 4.500e-03 5.040e-02 2.460e-02 2.240e-02 5.446e-01
4.800e-03 4.000e-04]]

Actual Class : 7

```
-----
15 Text feature [cells] present in test data point [True]
17 Text feature [cell] present in test data point [True]
18 Text feature [activated] present in test data point [True]
19 Text feature [kinase] present in test data point [True]
20 Text feature [downstream] present in test data point [True]
21 Text feature [contrast] present in test data point [True]
22 Text feature [activation] present in test data point [True]
23 Text feature [presence] present in test data point [True]
24 Text feature [factor] present in test data point [True]
25 Text feature [expressing] present in test data point [True]
26 Text feature [growth] present in test data point [True]
27 Text feature [inhibitor] present in test data point [True]
28 Text feature [phosphorylation] present in test data point [True]
29 Text feature [shown] present in test data point [True]
30 Text feature [10] present in test data point [True]
33 Text feature [also] present in test data point [True]
34 Text feature [suggest] present in test data point [True]
36 Text feature [however] present in test data point [True]
37 Text feature [treated] present in test data point [True]
38 Text feature [addition] present in test data point [True]
39 Text feature [compared] present in test data point [True]
40 Text feature [found] present in test data point [True]
41 Text feature [recently] present in test data point [True]
42 Text feature [independent] present in test data point [True]
44 Text feature [treatment] present in test data point [True]
```

45 Text feature [well] present in test data point [True]
46 Text feature [showed] present in test data point [True]
47 Text feature [mutations] present in test data point [True]
48 Text feature [similar] present in test data point [True]
49 Text feature [interestingly] present in test data point [True]
50 Text feature [increased] present in test data point [True]
51 Text feature [tyrosine] present in test data point [True]
52 Text feature [figure] present in test data point [True]
53 Text feature [potential] present in test data point [True]
54 Text feature [described] present in test data point [True]
55 Text feature [followed] present in test data point [True]
57 Text feature [higher] present in test data point [True]
58 Text feature [sensitive] present in test data point [True]
59 Text feature [mechanism] present in test data point [True]
60 Text feature [enhanced] present in test data point [True]
62 Text feature [using] present in test data point [True]
63 Text feature [constitutive] present in test data point [True]
64 Text feature [mutant] present in test data point [True]
66 Text feature [inhibition] present in test data point [True]
67 Text feature [consistent] present in test data point [True]
68 Text feature [demonstrated] present in test data point [True]
70 Text feature [constitutively] present in test data point [True]
71 Text feature [observed] present in test data point [True]
72 Text feature [may] present in test data point [True]
73 Text feature [inhibited] present in test data point [True]
74 Text feature [activating] present in test data point [True]
75 Text feature [serum] present in test data point [True]
77 Text feature [inhibitors] present in test data point [True]
78 Text feature [without] present in test data point [True]
79 Text feature [total] present in test data point [True]
80 Text feature [furthermore] present in test data point [True]
83 Text feature [pathways] present in test data point [True]
84 Text feature [reported] present in test data point [True]
85 Text feature [various] present in test data point [True]
87 Text feature [expression] present in test data point [True]
88 Text feature [mutation] present in test data point [True]
89 Text feature [leading] present in test data point [True]
90 Text feature [recent] present in test data point [True]
91 Text feature [including] present in test data point [True]
93 Text feature [increase] present in test data point [True]
94 Text feature [proliferation] present in test data point [True]
96 Text feature [confirmed] present in test data point [True]

97 Text feature [performed] present in test data point [True]
98 Text feature [absence] present in test data point [True]
99 Text feature [two] present in test data point [True]
Out of the top 100 features 70 are present in query point

4.1.1.4. Feature Importance, Incorrectly classified point

```
In [66]: test_point_index = 100
no_feature = 100
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['
```

Predicted Class : 1

Predicted Class Probabilities: [[0.3822 0.0239 0.001 0.3099 0.0152 0.0109 0.2535 0.0034 0.]]

Actual Class : 4

```
-----
10 Text feature [type] present in test data point [True]
11 Text feature [protein] present in test data point [True]
12 Text feature [wild] present in test data point [True]
15 Text feature [one] present in test data point [True]
16 Text feature [two] present in test data point [True]
17 Text feature [containing] present in test data point [True]
18 Text feature [dna] present in test data point [True]
19 Text feature [therefore] present in test data point [True]
20 Text feature [results] present in test data point [True]
22 Text feature [also] present in test data point [True]
23 Text feature [function] present in test data point [True]
24 Text feature [binding] present in test data point [True]
26 Text feature [region] present in test data point [True]
27 Text feature [effect] present in test data point [True]
28 Text feature [control] present in test data point [True]
29 Text feature [three] present in test data point [True]
30 Text feature [table] present in test data point [True]
32 Text feature [four] present in test data point [True]
33 Text feature [using] present in test data point [True]
34 Text feature [determined] present in test data point [True]
35 Text feature [specific] present in test data point [True]
38 Text feature [indicate] present in test data point [True]
41 Text feature [shown] present in test data point [True]
43 Text feature [complex] present in test data point [True]
45 Text feature [loss] present in test data point [True]
46 Text feature [similar] present in test data point [True]
```

47 Text feature [ability] present in test data point [True]
49 Text feature [human] present in test data point [True]
50 Text feature [either] present in test data point [True]
51 Text feature [several] present in test data point [True]
54 Text feature [important] present in test data point [True]
59 Text feature [amino] present in test data point [True]
60 Text feature [result] present in test data point [True]
61 Text feature [different] present in test data point [True]
62 Text feature [addition] present in test data point [True]
63 Text feature [page] present in test data point [True]
67 Text feature [observed] present in test data point [True]
68 Text feature [expression] present in test data point [True]
70 Text feature [fig] present in test data point [True]
72 Text feature [present] present in test data point [True]
74 Text feature [described] present in test data point [True]
75 Text feature [many] present in test data point [True]
76 Text feature [cancer] present in test data point [True]
78 Text feature [analysis] present in test data point [True]
79 Text feature [whether] present in test data point [True]
80 Text feature [gene] present in test data point [True]
82 Text feature [respectively] present in test data point [True]
83 Text feature [structure] present in test data point [True]
84 Text feature [even] present in test data point [True]
86 Text feature [analyzed] present in test data point [True]
87 Text feature [transcription] present in test data point [True]
88 Text feature [well] present in test data point [True]
90 Text feature [example] present in test data point [True]
91 Text feature [mutant] present in test data point [True]
93 Text feature [including] present in test data point [True]
95 Text feature [used] present in test data point [True]
97 Text feature [compared] present in test data point [True]
98 Text feature [thus] present in test data point [True]
99 Text feature [domain] present in test data point [True]
Out of the top 100 features 59 are present in query point

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

```

In [67]: # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

# methods of
# fit(X, y) : Fit the model using X as training data and y as target values
# predict(X):Predict the class labels for the provided data
# predict_proba(X):Return probability estimates for the test data X.
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-geometric-intuition
#-----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV
# -----
# default parameters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
#-----
# video link:
#-----

alpha = [5, 11, 15, 21, 31, 41, 51, 99]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = KNeighborsClassifier(n_neighbors=i)
    clf.fit(train_x_responseCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_responseCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilities we use log-probability estimates

```

```

    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

best_alpha = np.argmin(cv_log_error_array)
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

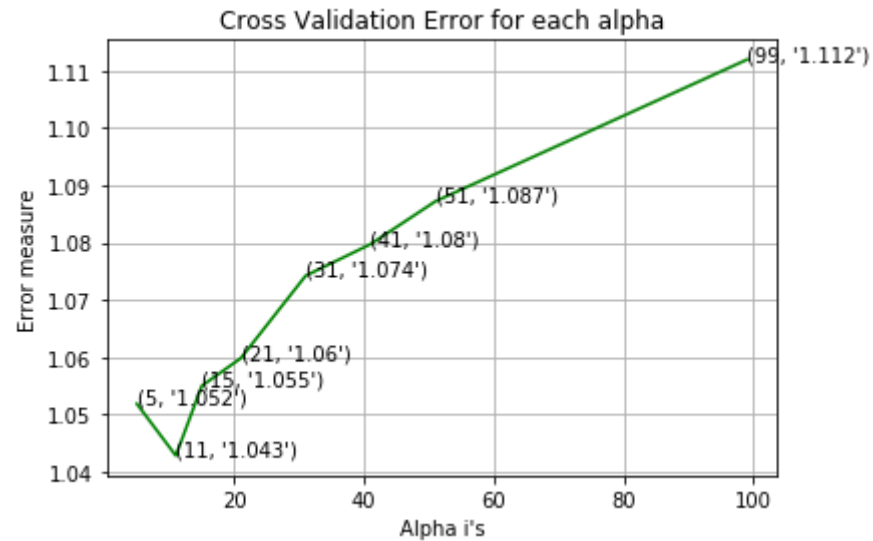
predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, la
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.c

for alpha = 5
Log Loss : 1.0518944444555418
for alpha = 11
Log Loss : 1.042853704031021
for alpha = 15
Log Loss : 1.0550279749751241
for alpha = 21
Log Loss : 1.0598376325018803
for alpha = 31
Log Loss : 1.0742354684838895
for alpha = 41
Log Loss : 1.079743560091223
for alpha = 51
Log Loss : 1.0872246594236858

```

for alpha = 99

Log Loss : 1.1119448514615644



For values of best alpha = 11 The train log loss is: 0.6472674048734882

For values of best alpha = 11 The cross validation log loss is: 1.042853704031021

For values of best alpha = 11 The test log loss is: 1.0788862450064254

4.2.2. Testing the model with best hyper paramters

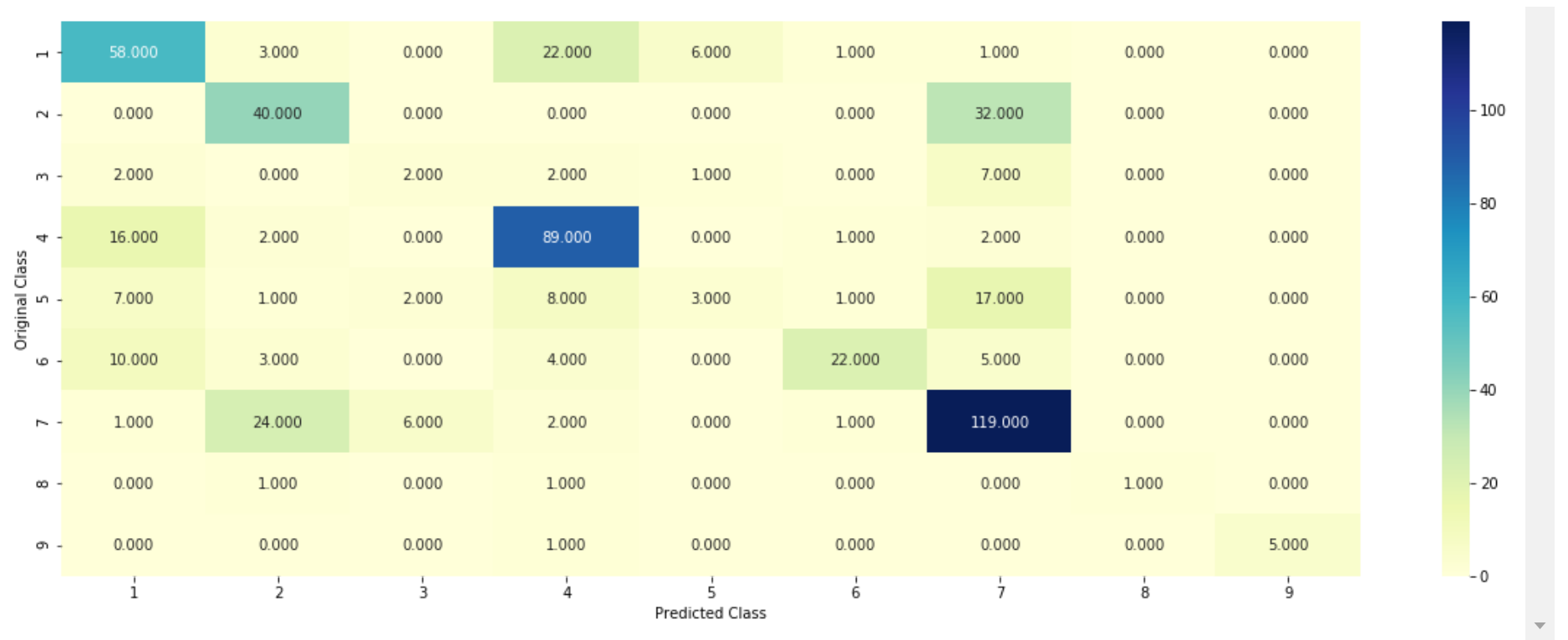

```
In [68]: # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

# methods of
# fit(X, y) : Fit the model using X as training data and y as target values
# predict(X):Predict the class labels for the provided data
# predict_proba(X):Return probability estimates for the test data X.
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-geometric-intu
#-----
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseCoding, cv_y, clf)
```

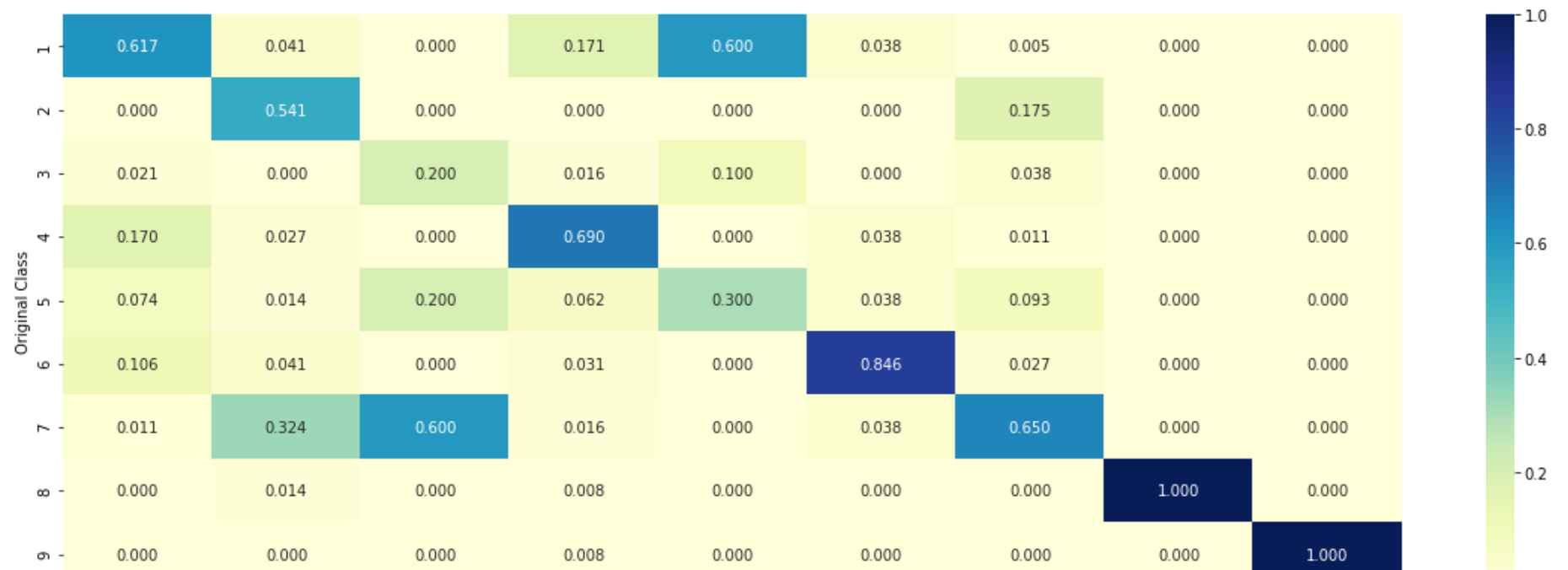
Log loss : 1.042853704031021

Number of mis-classified points : 0.36278195488721804

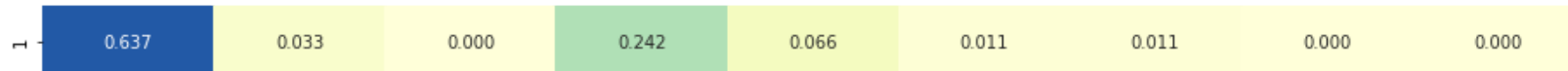
----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----



----- Recall matrix (Row sum=1) -----



4.2.3. Sample Query point -1

```
In [69]: clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

test_point_index = 1
predicted_cls = sig_clf.predict(test_x_responseCoding[0].reshape(1,-1))
print("Predicted Class :", predicted_cls[0])
print("Actual Class :", test_y[test_point_index])
neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
print("The ",alpha[best_alpha]," nearest neighbours of the test points belongs to classes",train_y[neighbors[1][0]])
print("Fequency of nearest points :",Counter(train_y[neighbors[1][0]]))
```

Predicted Class : 7

Actual Class : 7

The 11 nearest neighbours of the test points belongs to classes [1 7 2 2 5 7 2 2 2 7 2]

Fequency of nearest points : Counter({2: 6, 7: 3, 1: 1, 5: 1})

4.2.4. Sample Query Point-2

```
In [72]: clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

test_point_index = 201

predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
print("Predicted Class :", predicted_cls[0])
print("Actual Class :", test_y[test_point_index])
neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
print("the k value for knn is",alpha[best_alpha],"and the nearest neighbours of the test points belongs to classes",train_y[neighbors[1][0]])
print("Fequency of nearest points :",Counter(train_y[neighbors[1][0]]))
```

Predicted Class : 7

Actual Class : 3

the k value for knn is 11 and the nearest neighbours of the test points belongs to classes [7 7 7 5 5 7 7 7 7 7 7]

Fequency of nearest points : Counter({7: 9, 5: 2})

4.3. Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper paramter tuning

In [73]:

```
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier
# -----
# 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: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.Ca
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
#-----
# video link:
#-----

alpha = [10 ** x for x in range(-6, 3)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
```

```

# to avoid rounding error while multiplying probabilities we use log-probability estimates
print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.labels_))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.labels_))
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.labels_))

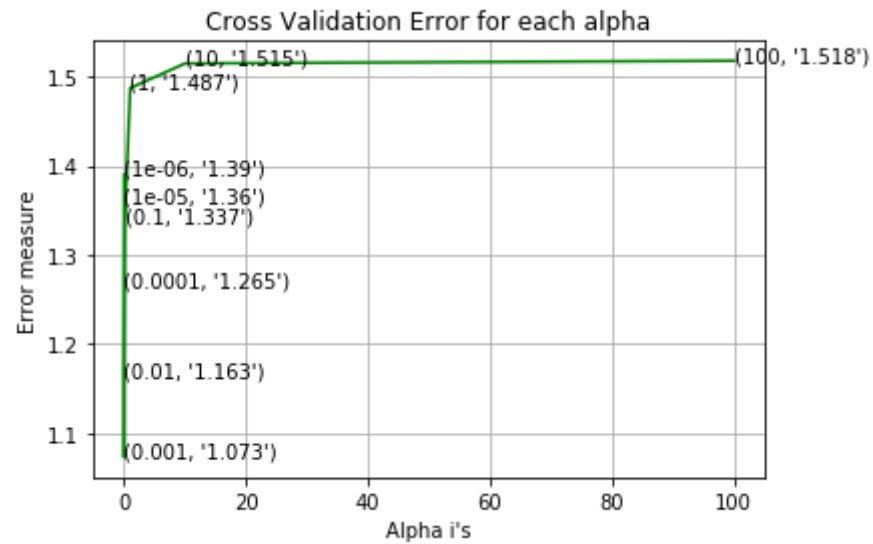
```

```

for alpha = 1e-06
Log Loss : 1.3901769755526199
for alpha = 1e-05
Log Loss : 1.3600404981745577
for alpha = 0.0001
Log Loss : 1.2653247361854167
for alpha = 0.001
Log Loss : 1.0725976007358062
for alpha = 0.01
Log Loss : 1.1631688428590035
for alpha = 0.1
Log Loss : 1.3365994821717635
for alpha = 1
Log Loss : 1.487359548323489

```

for alpha = 10
Log Loss : 1.5150588693299003
for alpha = 100
Log Loss : 1.5181067130716062



For values of best alpha = 0.001 The train log loss is: 0.6053141220312026
For values of best alpha = 0.001 The cross validation log loss is: 1.0725976007358062
For values of best alpha = 0.001 The test log loss is: 1.083944230571662

4.3.1.2. Testing the model with best hyper paramters


```
In [74]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier
# -----
# 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)

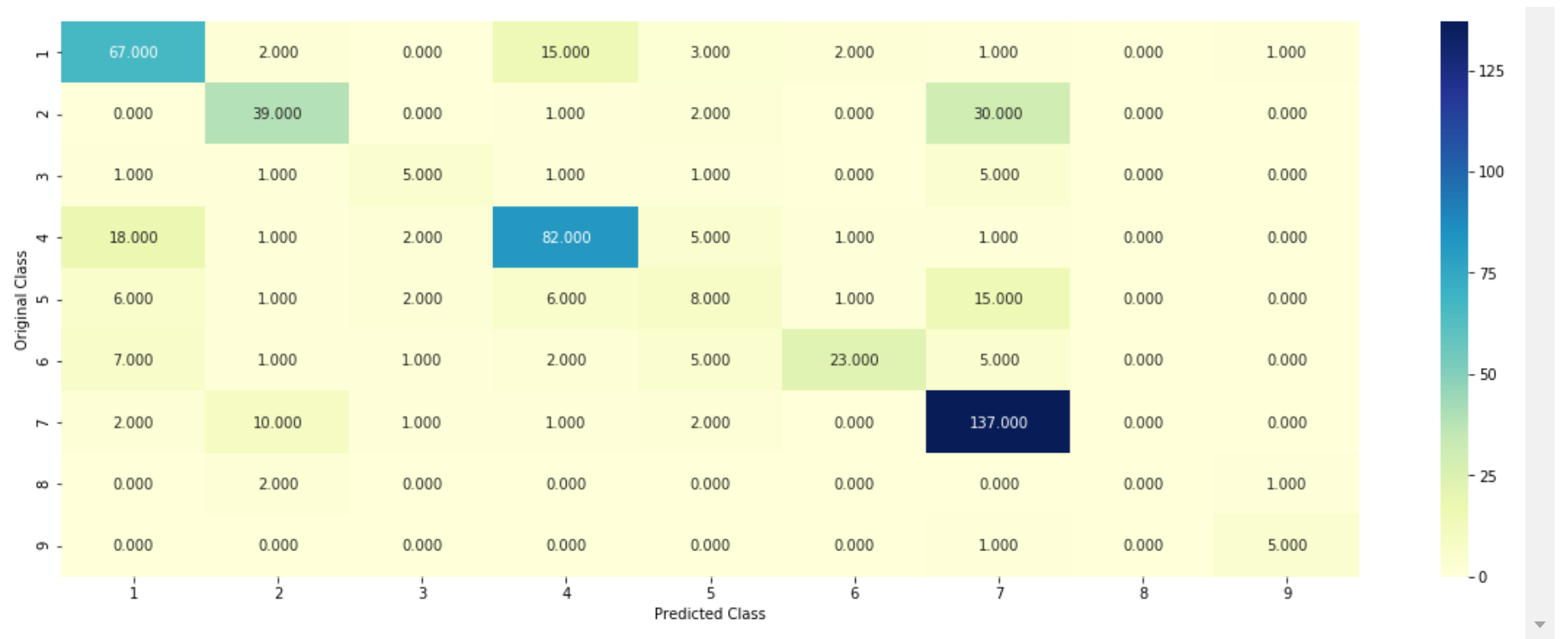
# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X)    Predict class labels for samples in X.

#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```

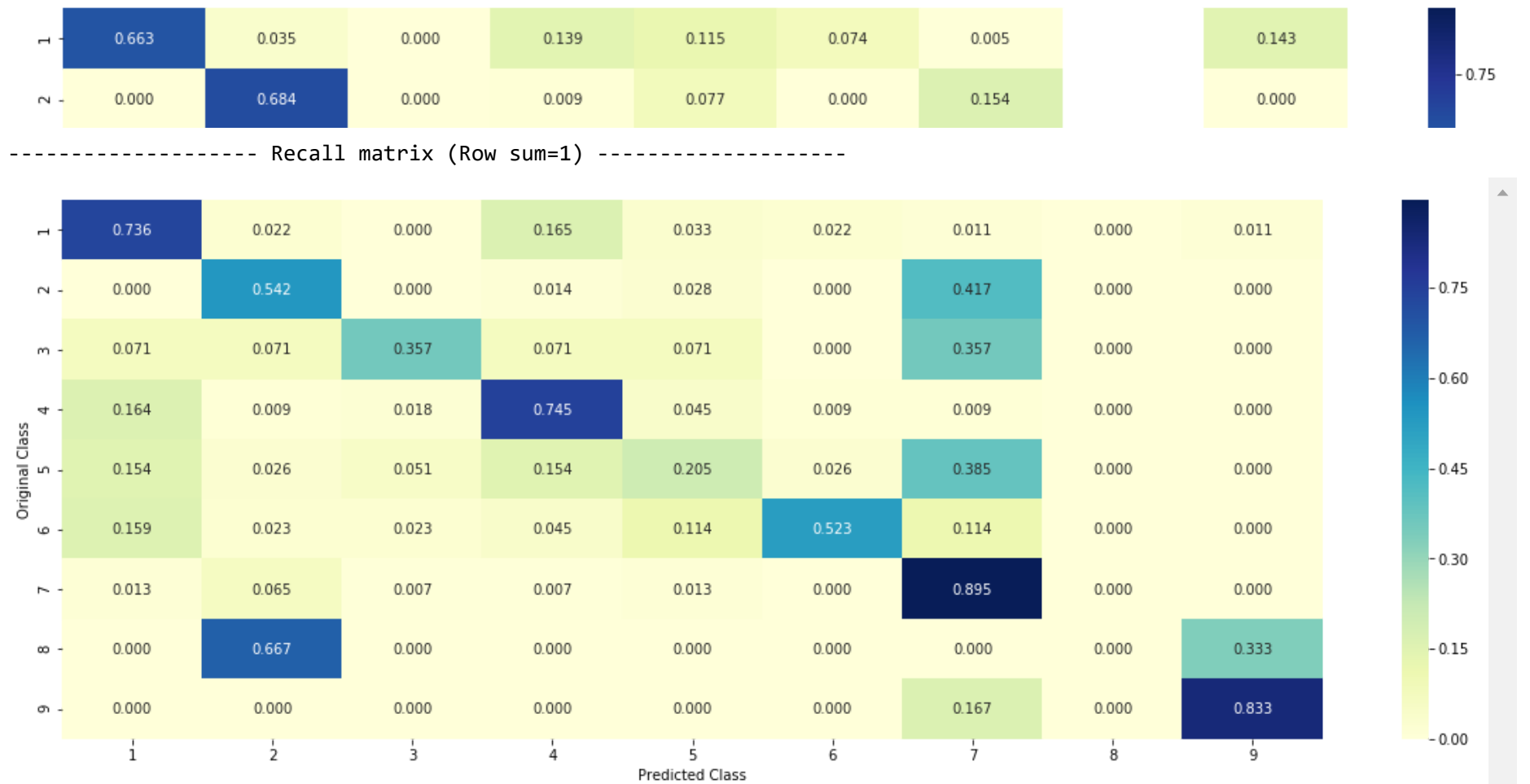
Log loss : 1.0725976007358062

Number of mis-classified points : 0.31203007518796994

----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----



4.3.1.3. Feature Importance

```

In [75]: def get_imp_feature_names(text, indices, removed_ind = []):
    word_present = 0
    tabulte_list = []
    incresingorder_ind = 0
    for i in indices:
        if i < train_gene_feature_onehotCoding.shape[1]:
            tabulte_list.append([incresingorder_ind, "Gene", "Yes"])
        elif i < 18:
            tabulte_list.append([incresingorder_ind, "Variation", "Yes"])
        if ((i > 17) & (i not in removed_ind)) :
            word = train_text_features[i]
            yes_no = True if word in text.split() else False
            if yes_no:
                word_present += 1
            tabulte_list.append([incresingorder_ind, train_text_features[i], yes_no])
            incresingorder_ind += 1
    print(word_present, "most important features are present in our query point")
    print("-"*50)
    print("The features that are most important of the ", predicted_cls[0], " class:")
    print (tabulate(tabulte_list, headers=["Index", 'Feature name', 'Present or Not']))

```

4.3.1.3.1. Correctly Classified point

```
In [76]: # from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding,train_y)
test_point_index = 1
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0167 0.1857 0.0102 0.0161 0.0138 0.0043 0.7394 0.0109 0.0029]]

Actual Class : 7

```
-----
152 Text feature [acquired] present in test data point [True]
227 Text feature [purple] present in test data point [True]
228 Text feature [colorectal] present in test data point [True]
256 Text feature [selective] present in test data point [True]
301 Text feature [administration] present in test data point [True]
304 Text feature [respond] present in test data point [True]
329 Text feature [pka] present in test data point [True]
341 Text feature [primary] present in test data point [True]
350 Text feature [evaluated] present in test data point [True]
402 Text feature [along] present in test data point [True]
419 Text feature [space] present in test data point [True]
424 Text feature [common] present in test data point [True]
454 Text feature [sry] present in test data point [True]
455 Text feature [level] present in test data point [True]
460 Text feature [frs2] present in test data point [True]
461 Text feature [belongs] present in test data point [True]
473 Text feature [calcium] present in test data point [True]
478 Text feature [tkis] present in test data point [True]
Out of the top 500 features 18 are present in query point
```

4.3.1.3.2. Incorrectly Classified point

```
In [77]: test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['
```

```
Predicted Class : 1
Predicted Class Probabilities: [[0.4964 0.072  0.0123 0.2856 0.0308 0.0076 0.0731 0.0175 0.0047]]
Actual Class : 4
-----
51 Text feature [participating] present in test data point [True]
131 Text feature [broad] present in test data point [True]
Out of the top 500 features 2 are present in query point
```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

```

In [78]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier
# -----
# 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: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.Ca
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
#-----
# video link:
#-----

alpha = [10 ** x for x in range(-6, 1)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))

```

```

    print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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

best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

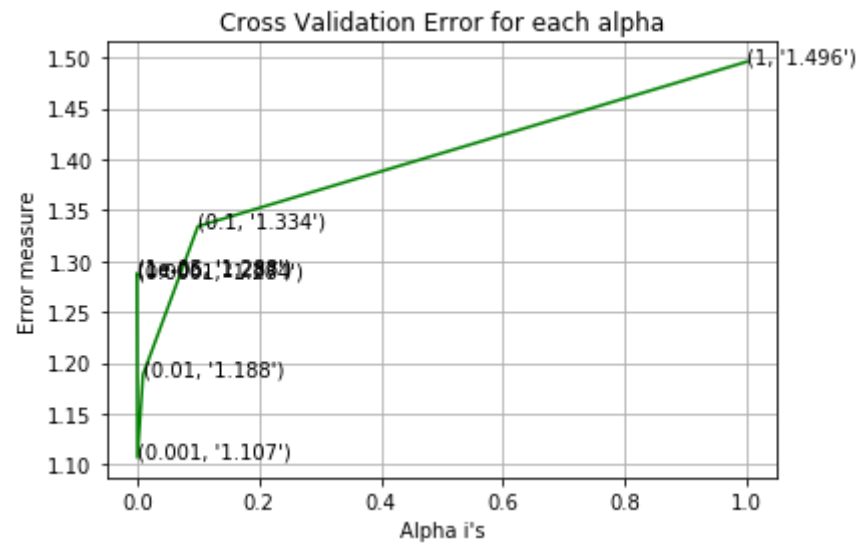
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, la
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.c

```

```

for alpha = 1e-06
Log Loss : 1.286770176304231
for alpha = 1e-05
Log Loss : 1.2880372377001756
for alpha = 0.0001
Log Loss : 1.284373920943486
for alpha = 0.001
Log Loss : 1.106969327419523
for alpha = 0.01
Log Loss : 1.188495284904659
for alpha = 0.1
Log Loss : 1.3342748144156729
for alpha = 1
Log Loss : 1.4960862543572748

```

For values of best alpha = 0.001 The train log loss is: 0.6050412503497927

For values of best alpha = 0.001 The cross validation log loss is: 1.106969327419523

For values of best alpha = 0.001 The test log loss is: 1.1121556788800262

4.3.2.2. Testing model with best hyper parameters

```
In [79]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier
# -----
# 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.

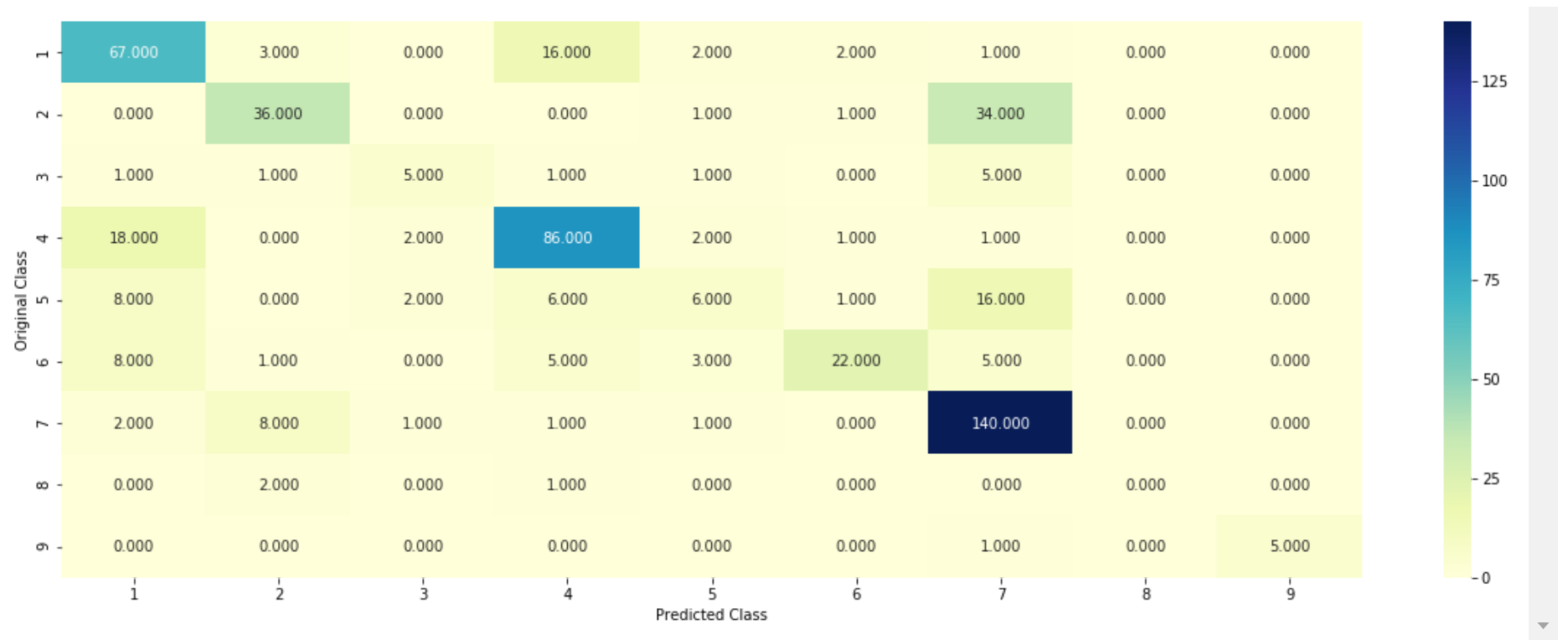
#-----
# video link:
#-----

clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```

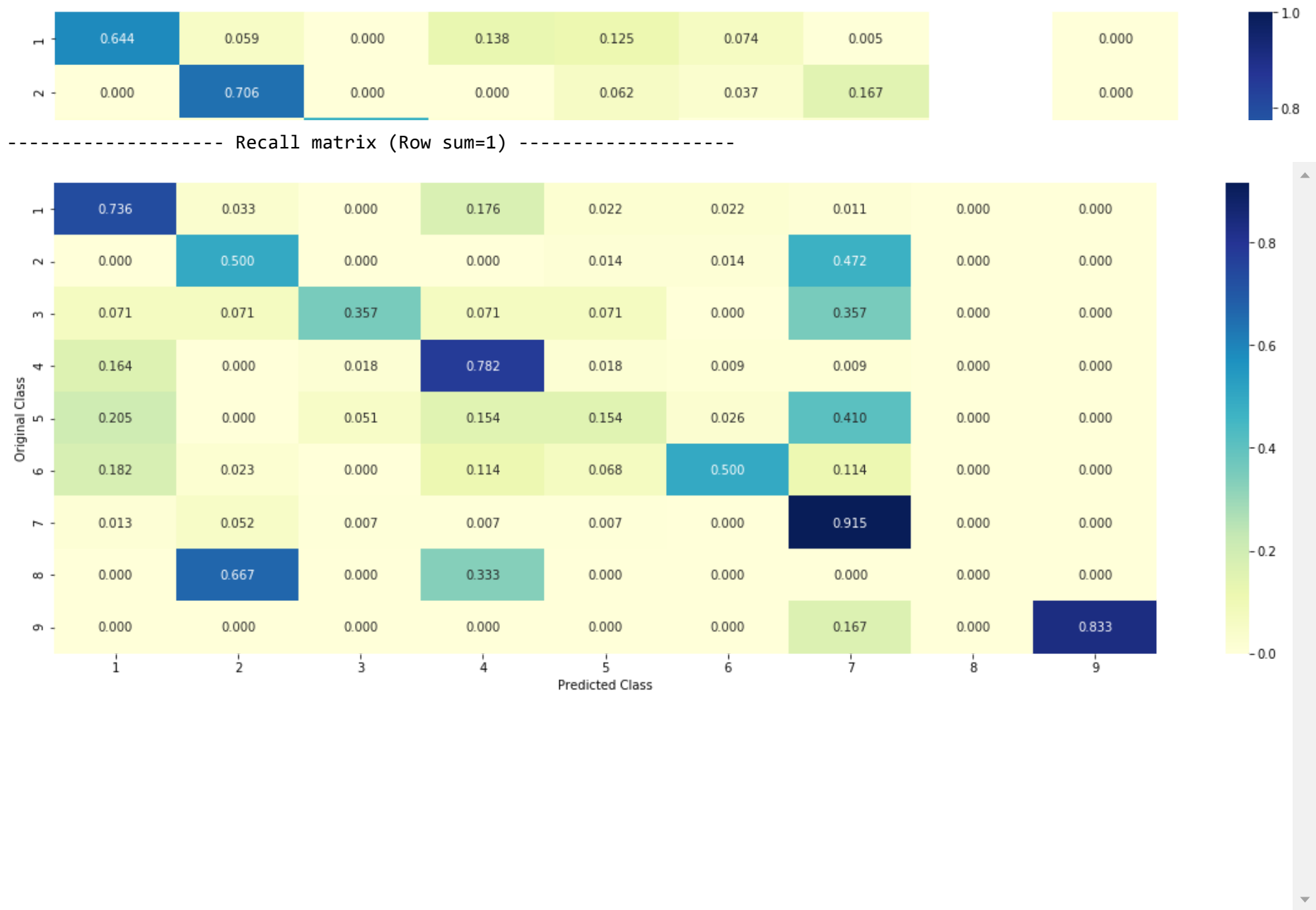
Log loss : 1.106969327419523

Number of mis-classified points : 0.3101503759398496

----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----



4.3.2.3. Feature Importance, Correctly Classified point

```
In [80]: clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding,train_y)
test_point_index = 1
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['
```

```
Predicted Class : 7
Predicted Class Probabilities: [[1.870e-02 1.642e-01 1.300e-03 1.870e-02 9.000e-03 3.100e-03 7.713e-01
 1.350e-02 1.000e-04]]
Actual Class : 7
```

```
-----
1 Text feature [glossary] present in test data point [True]
28 Text feature [19] present in test data point [True]
100 Text feature [became] present in test data point [True]
201 Text feature [males] present in test data point [True]
207 Text feature [modulated] present in test data point [True]
319 Text feature [smokers] present in test data point [True]
326 Text feature [examination] present in test data point [True]
331 Text feature [288] present in test data point [True]
332 Text feature [represent] present in test data point [True]
334 Text feature [factors] present in test data point [True]
361 Text feature [adaptor] present in test data point [True]
370 Text feature [myeloid] present in test data point [True]
382 Text feature [signals] present in test data point [True]
400 Text feature [website] present in test data point [True]
449 Text feature [types] present in test data point [True]
458 Text feature [artemin] present in test data point [True]
459 Text feature [persephin] present in test data point [True]
465 Text feature [repeated] present in test data point [True]
481 Text feature [becomes] present in test data point [True]
Out of the top 500 features 19 are present in query point
```

4.3.2.4. Feature Importance, Inorrectly Classified point

```
In [81]: test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_index], test_df['
```

```
Predicted Class : 1
Predicted Class Probabilities: [[4.557e-01 6.900e-02 2.100e-03 3.316e-01 1.920e-02 6.000e-03 9.640e-02
 1.970e-02 2.000e-04]]
Actual Class : 4
-----
156 Text feature [regulatory] present in test data point [True]
173 Text feature [terminator] present in test data point [True]
313 Text feature [responsive] present in test data point [True]
319 Text feature [pairs] present in test data point [True]
Out of the top 500 features 4 are present in query point
```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

```

In [82]: # read more about support vector machines with linear kernals here http://scikit-learn.org/stable/modules/generated/skle

# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_state=None)

# Some of methods of SVM()
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
# predict(X) Perform classification on samples in X.
# -----
# video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.Ca
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
#-----
# video link:
#-----

alpha = [10 ** x for x in range(-5, 3)]
cv_log_error_array = []
for i in alpha:
    print("for C =", i)
    # clf = SVC(C=i, kernel='linear', probability=True, class_weight='balanced')
    clf = SGDClassifier( class_weight='balanced', alpha=i, penalty='l2', loss='hinge', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))

```

```

    print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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

best_alpha = np.argmin(cv_log_error_array)
# clf = SVC(C=i,kernel='linear',probability=True, class_weight='balanced')
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, la
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.c

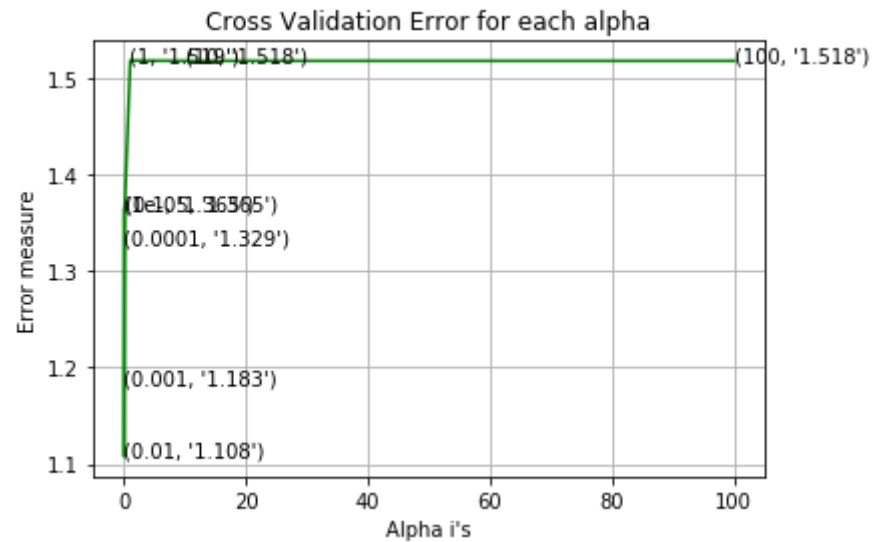
```

```

for C = 1e-05
Log Loss : 1.3647050095705298
for C = 0.0001
Log Loss : 1.329195084687183
for C = 0.001
Log Loss : 1.1832860555461862
for C = 0.01
Log Loss : 1.107547081775696
for C = 0.1
Log Loss : 1.3649491279140222
for C = 1
Log Loss : 1.5187307586360943
for C = 10

```


Log Loss : 1.5184486323400181
for C = 100
Log Loss : 1.518448898821054



For values of best alpha = 0.01 The train log loss is: 0.727383526677587
For values of best alpha = 0.01 The cross validation log loss is: 1.107547081775696
For values of best alpha = 0.01 The test log loss is: 1.1312736551858935

4.4.2. Testing model with best hyper parameters

In [83]: *# read more about support vector machines with linear kernals here [```

default parameters
SVC\(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_state=None\)

Some of methods of SVM\(\)
fit\(X, y, \[sample_weight\]\) Fit the SVM model according to the given training data.
predict\(X\) Perform classification on samples in X.

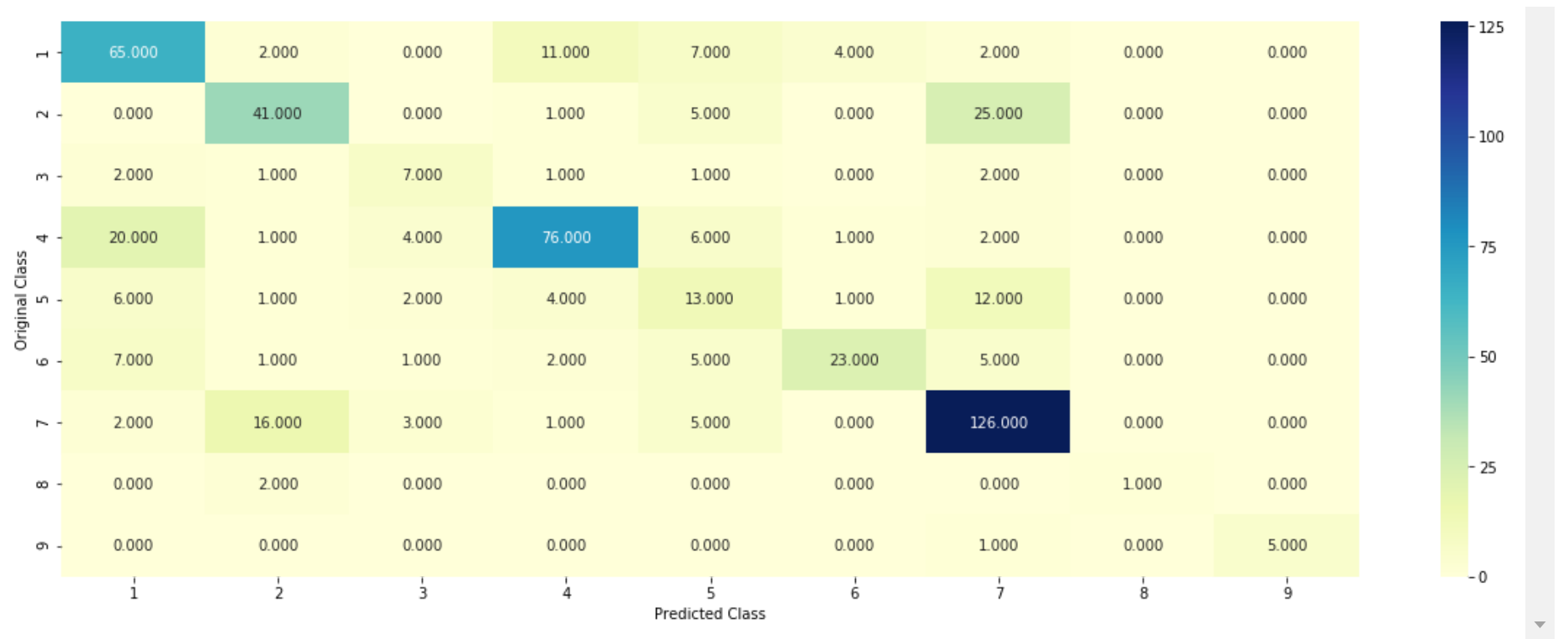
video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/

clf = SVC\(C=alpha\[best_alpha\],kernel='Linear',probability=True, class_weight='balanced'\)
clf = SGDClassifier\(alpha=alpha\[best_alpha\], penalty='l2', loss='hinge', random_state=42,class_weight='balanced'\)
predict_and_plot_confusion_matrix\(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf\)
```](http://scikit-learn.org/stable/modules/generated/skle</a></i></p></div><div data-bbox=)*

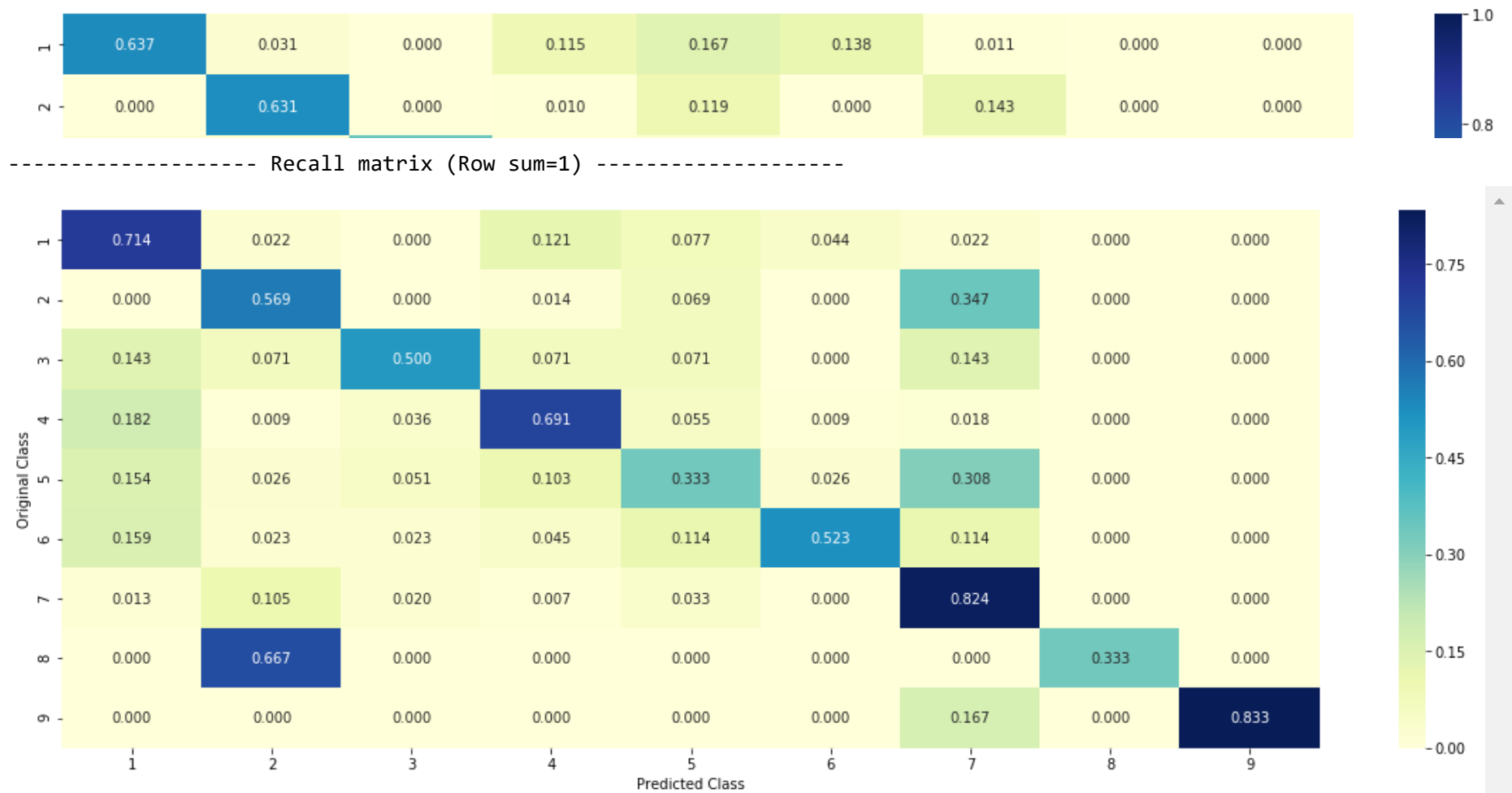
Log loss : 1.107547081775696

Number of mis-classified points : 0.32894736842105265

----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----



#### 4.3.3. Feature Importance

#### 4.3.3.1. For Correctly classified point

```
In [84]: clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
clf.fit(train_x_onehotCoding,train_y)
test_point_index = 1
test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:,no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0569 0.3293 0.0116 0.0683 0.0607 0.0409 0.4019 0.0208 0.0097]]

Actual Class : 7

-----

Out of the top 500 features 0 are present in query point

#### 4.3.3.2. For Incorrectly classified point

```
In [85]: test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(abs(-clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_index], test_df['
```

Predicted Class : 1

Predicted Class Probabilities: [[0.5202 0.0922 0.0074 0.1809 0.0476 0.0146 0.1176 0.0141 0.0054]]

Actual Class : 4

-----

Out of the top 500 features 0 are present in query point

## 4.5 Random Forest Classifier

### 4.5.1. Hyper paramter tuning (With One hot Encoding)

```

In [86]: # -----
default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
class_weight=None)

Some of methods of RandomForestClassifier()
fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
predict(X) Perform classification on samples in X.
predict_proba (X) Perform classification on samples in X.

some of attributes of RandomForestClassifier()
feature_importances_ : array of shape = [n_features]
The feature importances (the higher, the more important the feature).

video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-constructi

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.Ca

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

alpha = [100,200,500,1000,2000]
max_depth = [5, 10]
cv_log_error_array = []
for i in alpha:
 for j in max_depth:
 print("for n_estimators =", i,"and max depth = ", j)

```

```

clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
print("Log Loss :", log_loss(cv_y, sig_clf_probs))

'''fig, ax = plt.subplots()
features = np.dot(np.array(alpha)[:,None], np.array(max_depth)[None]).ravel()
ax.plot(features, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
 ax.annotate((alpha[int(i/2)], max_depth[int(i%2)], str(txt)), (features[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()
'''

best_alpha = np.argmin(cv_log_error_array)
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha/2)])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_))

```

```

for n_estimators = 100 and max depth = 5
Log Loss : 1.26498348765677
for n_estimators = 100 and max depth = 10
Log Loss : 1.2107261579827089
for n_estimators = 200 and max depth = 5
Log Loss : 1.2360101622602102
for n_estimators = 200 and max depth = 10
Log Loss : 1.1858265649908228

```



```
for n_estimators = 500 and max depth = 5
Log Loss : 1.226147462509591
for n_estimators = 500 and max depth = 10
Log Loss : 1.1778152763096252
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2290535390303687
for n_estimators = 1000 and max depth = 10
Log Loss : 1.1771476362705846
for n_estimators = 2000 and max depth = 5
Log Loss : 1.2235155702670426
for n_estimators = 2000 and max depth = 10
Log Loss : 1.1754807765324484
For values of best estimator = 2000 The train log loss is: 0.6547374900424958
For values of best estimator = 2000 The cross validation log loss is: 1.1754807765324482
For values of best estimator = 2000 The test log loss is: 1.1517109615567254
```

#### **4.5.2. Testing model with best hyper parameters (One Hot Encoding)**

```
In [87]: # -----
default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
class_weight=None)

Some of methods of RandomForestClassifier()
fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
predict(X) Perform classification on samples in X.
predict_proba (X) Perform classification on samples in X.

some of attributes of RandomForestClassifier()
feature_importances_ : array of shape = [n_features]
The feature importances (the higher, the more important the feature).

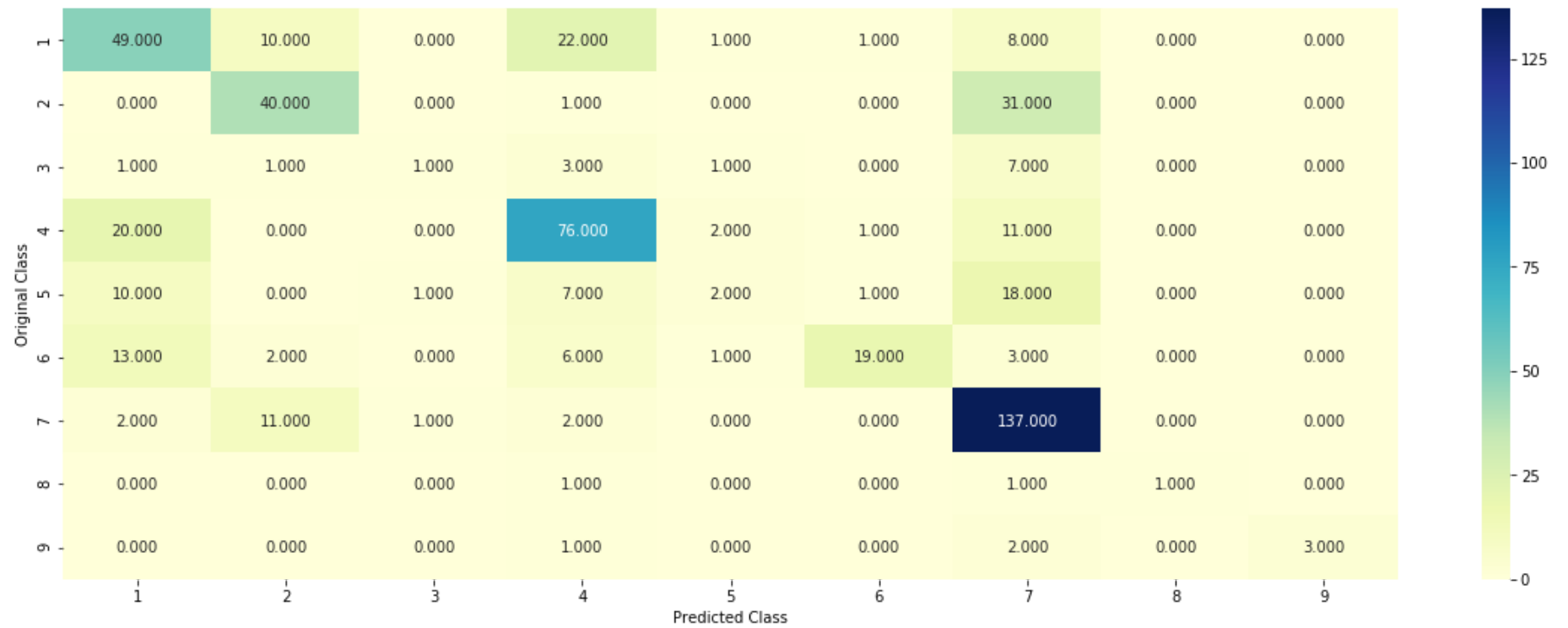
video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-constructi

clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf)
```

Log loss : 1.1754807765324484

Number of mis-classified points : 0.38345864661654133

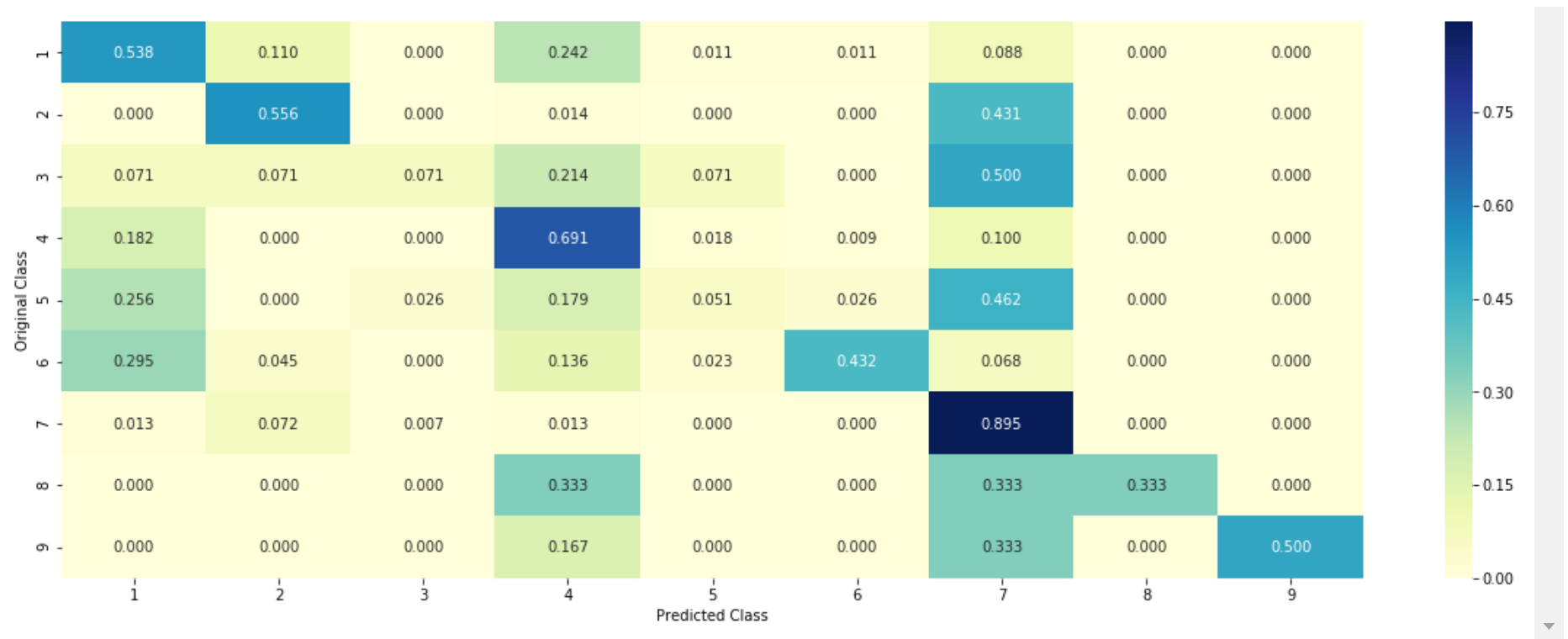
----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----



----- Recall matrix (Row sum=1) -----



### 4.5.3. Feature Importance

#### 4.5.3.1. Correctly Classified point

```
In [88]: # test_point_index = 10
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

test_point_index = 1
no_feature = 100
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_impfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index])
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0731 0.269 0.0153 0.0703 0.0451 0.0415 0.4748 0.0057 0.0051]]

Actual Class : 7

```

0 Text feature [kinase] present in test data point [True]
1 Text feature [activating] present in test data point [True]
2 Text feature [tyrosine] present in test data point [True]
3 Text feature [activation] present in test data point [True]
4 Text feature [suppressor] present in test data point [True]
5 Text feature [phosphorylation] present in test data point [True]
6 Text feature [constitutive] present in test data point [True]
7 Text feature [inhibitor] present in test data point [True]
8 Text feature [activated] present in test data point [True]
9 Text feature [treatment] present in test data point [True]
10 Text feature [inhibitors] present in test data point [True]
11 Text feature [akt] present in test data point [True]
12 Text feature [function] present in test data point [True]
14 Text feature [oncogenic] present in test data point [True]
16 Text feature [loss] present in test data point [True]
17 Text feature [pathogenic] present in test data point [True]
19 Text feature [missense] present in test data point [True]
20 Text feature [downstream] present in test data point [True]
21 Text feature [therapeutic] present in test data point [True]
22 Text feature [trials] present in test data point [True]
```

23 Text feature [constitutively] present in test data point [True]  
24 Text feature [inhibited] present in test data point [True]  
27 Text feature [patients] present in test data point [True]  
28 Text feature [receptor] present in test data point [True]  
29 Text feature [growth] present in test data point [True]  
30 Text feature [serum] present in test data point [True]  
33 Text feature [advanced] present in test data point [True]  
34 Text feature [activate] present in test data point [True]  
35 Text feature [treated] present in test data point [True]  
36 Text feature [cells] present in test data point [True]  
38 Text feature [inhibition] present in test data point [True]  
41 Text feature [functional] present in test data point [True]  
42 Text feature [proliferation] present in test data point [True]  
43 Text feature [months] present in test data point [True]  
44 Text feature [variants] present in test data point [True]  
46 Text feature [drug] present in test data point [True]  
48 Text feature [kinases] present in test data point [True]  
50 Text feature [proteins] present in test data point [True]  
51 Text feature [cell] present in test data point [True]  
52 Text feature [expressing] present in test data point [True]  
53 Text feature [resistance] present in test data point [True]  
54 Text feature [transforming] present in test data point [True]  
55 Text feature [efficacy] present in test data point [True]  
56 Text feature [protein] present in test data point [True]  
59 Text feature [ligand] present in test data point [True]  
63 Text feature [clinical] present in test data point [True]  
65 Text feature [autophosphorylation] present in test data point [True]  
66 Text feature [extracellular] present in test data point [True]  
68 Text feature [factor] present in test data point [True]  
71 Text feature [functions] present in test data point [True]  
75 Text feature [phosphorylated] present in test data point [True]  
80 Text feature [survival] present in test data point [True]  
81 Text feature [variant] present in test data point [True]  
86 Text feature [classified] present in test data point [True]  
87 Text feature [egfr] present in test data point [True]  
89 Text feature [sensitivity] present in test data point [True]  
92 Text feature [respond] present in test data point [True]  
94 Text feature [effective] present in test data point [True]  
95 Text feature [affected] present in test data point [True]  
96 Text feature [patient] present in test data point [True]  
97 Text feature [pathway] present in test data point [True]

98 Text feature [tki] present in test data point [True]  
Out of the top 100 features 62 are present in query point

#### **4.5.3.2. Incorrectly Classified point**



```
In [90]: test_point_index = 201
no_feature = 100
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_impfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index])
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0355 0.1143 0.0142 0.0382 0.0432 0.0286 0.7172 0.0049 0.0039]]

Actual Class : 3

```

0 Text feature [kinase] present in test data point [True]
1 Text feature [activating] present in test data point [True]
2 Text feature [tyrosine] present in test data point [True]
3 Text feature [activation] present in test data point [True]
5 Text feature [phosphorylation] present in test data point [True]
6 Text feature [constitutive] present in test data point [True]
7 Text feature [inhibitor] present in test data point [True]
8 Text feature [activated] present in test data point [True]
9 Text feature [treatment] present in test data point [True]
10 Text feature [inhibitors] present in test data point [True]
11 Text feature [akt] present in test data point [True]
14 Text feature [oncogenic] present in test data point [True]
16 Text feature [loss] present in test data point [True]
18 Text feature [signaling] present in test data point [True]
20 Text feature [downstream] present in test data point [True]
23 Text feature [constitutively] present in test data point [True]
24 Text feature [inhibited] present in test data point [True]
27 Text feature [patients] present in test data point [True]
28 Text feature [receptor] present in test data point [True]
29 Text feature [growth] present in test data point [True]
30 Text feature [serum] present in test data point [True]
35 Text feature [treated] present in test data point [True]
36 Text feature [cells] present in test data point [True]
38 Text feature [inhibition] present in test data point [True]
42 Text feature [proliferation] present in test data point [True]
43 Text feature [months] present in test data point [True]
```

44 Text feature [variants] present in test data point [True]  
45 Text feature [amplification] present in test data point [True]  
50 Text feature [proteins] present in test data point [True]  
51 Text feature [cell] present in test data point [True]  
52 Text feature [expressing] present in test data point [True]  
53 Text feature [resistance] present in test data point [True]  
54 Text feature [transforming] present in test data point [True]  
56 Text feature [protein] present in test data point [True]  
57 Text feature [neutral] present in test data point [True]  
59 Text feature [ligand] present in test data point [True]  
60 Text feature [ic50] present in test data point [True]  
61 Text feature [starved] present in test data point [True]  
65 Text feature [autophosphorylation] present in test data point [True]  
66 Text feature [extracellular] present in test data point [True]  
68 Text feature [factor] present in test data point [True]  
72 Text feature [harboring] present in test data point [True]  
73 Text feature [dose] present in test data point [True]  
76 Text feature [lines] present in test data point [True]  
77 Text feature [daily] present in test data point [True]  
78 Text feature [assays] present in test data point [True]  
80 Text feature [survival] present in test data point [True]  
81 Text feature [variant] present in test data point [True]  
86 Text feature [classified] present in test data point [True]  
89 Text feature [sensitivity] present in test data point [True]  
91 Text feature [imatinib] present in test data point [True]  
94 Text feature [effective] present in test data point [True]  
95 Text feature [affected] present in test data point [True]  
96 Text feature [patient] present in test data point [True]  
98 Text feature [tki] present in test data point [True]  
99 Text feature [interleukin] present in test data point [True]  
Out of the top 100 features 56 are present in query point

#### 4.5.3. Hyper paramter tuning (With Response Coding)

```

In [91]: # -----
default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
class_weight=None)

Some of methods of RandomForestClassifier()
fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
predict(X) Perform classification on samples in X.
predict_proba (X) Perform classification on samples in X.

some of attributes of RandomForestClassifier()
feature_importances_ : array of shape = [n_features]
The feature importances (the higher, the more important the feature).

video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-constructi

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.Ca

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

alpha = [10,50,100,200,500,1000]
max_depth = [2,3,5,10]
cv_log_error_array = []
for i in alpha:
 for j in max_depth:
 print("for n_estimators =", i,"and max depth = ", j)

```

```

clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
print("Log Loss :", log_loss(cv_y, sig_clf_probs))
...
fig, ax = plt.subplots()
features = np.dot(np.array(alpha)[:,None], np.array(max_depth)[None]).ravel()
ax.plot(features, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
 ax.annotate((alpha[int(i/4)], max_depth[int(i%4)], str(txt)), (features[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()
...

best_alpha = np.argmin(cv_log_error_array)
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_depth=max_depth[int(best_alpha/4)])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_))

```

```

for n_estimators = 10 and max depth = 2
Log Loss : 1.9356531486743864
for n_estimators = 10 and max depth = 3
Log Loss : 1.6712232791197597
for n_estimators = 10 and max depth = 5
Log Loss : 1.412801502258631
for n_estimators = 10 and max depth = 10
Log Loss : 1.9825434355517357

```

```
for n_estimators = 50 and max depth = 2
Log Loss : 1.6382448414807238
for n_estimators = 50 and max depth = 3
Log Loss : 1.4745661360418654
for n_estimators = 50 and max depth = 5
Log Loss : 1.388978913629135
for n_estimators = 50 and max depth = 10
Log Loss : 1.6719297597115235
for n_estimators = 100 and max depth = 2
Log Loss : 1.5613860777993436
for n_estimators = 100 and max depth = 3
Log Loss : 1.477703893334632
for n_estimators = 100 and max depth = 5
Log Loss : 1.2833449142538584
for n_estimators = 100 and max depth = 10
Log Loss : 1.7175206236814797
for n_estimators = 200 and max depth = 2
Log Loss : 1.6145592118414251
for n_estimators = 200 and max depth = 3
Log Loss : 1.536061130018752
for n_estimators = 200 and max depth = 5
Log Loss : 1.3356913728935886
for n_estimators = 200 and max depth = 10
Log Loss : 1.6545899795078416
for n_estimators = 500 and max depth = 2
Log Loss : 1.7076474475954626
for n_estimators = 500 and max depth = 3
Log Loss : 1.5716302705731642
for n_estimators = 500 and max depth = 5
Log Loss : 1.396267615364467
for n_estimators = 500 and max depth = 10
Log Loss : 1.7238723632613444
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6699165011543078
for n_estimators = 1000 and max depth = 3
Log Loss : 1.5781890578364297
for n_estimators = 1000 and max depth = 5
Log Loss : 1.364714514645983
for n_estimators = 1000 and max depth = 10
Log Loss : 1.7185625870927776
For values of best alpha = 100 The train log loss is: 0.05262663545806976
```

For values of best alpha = 100 The cross validation log loss is: 1.2833449142538584  
For values of best alpha = 100 The test log loss is: 1.3297641376432576

#### **4.5.4. Testing model with best hyper parameters (Response Coding)**

```
In [92]: # -----
default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
class_weight=None)

Some of methods of RandomForestClassifier()
fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
predict(X) Perform classification on samples in X.
predict_proba (X) Perform classification on samples in X.

some of attributes of RandomForestClassifier()
feature_importances_ : array of shape = [n_features]
The feature importances (the higher, the more important the feature).

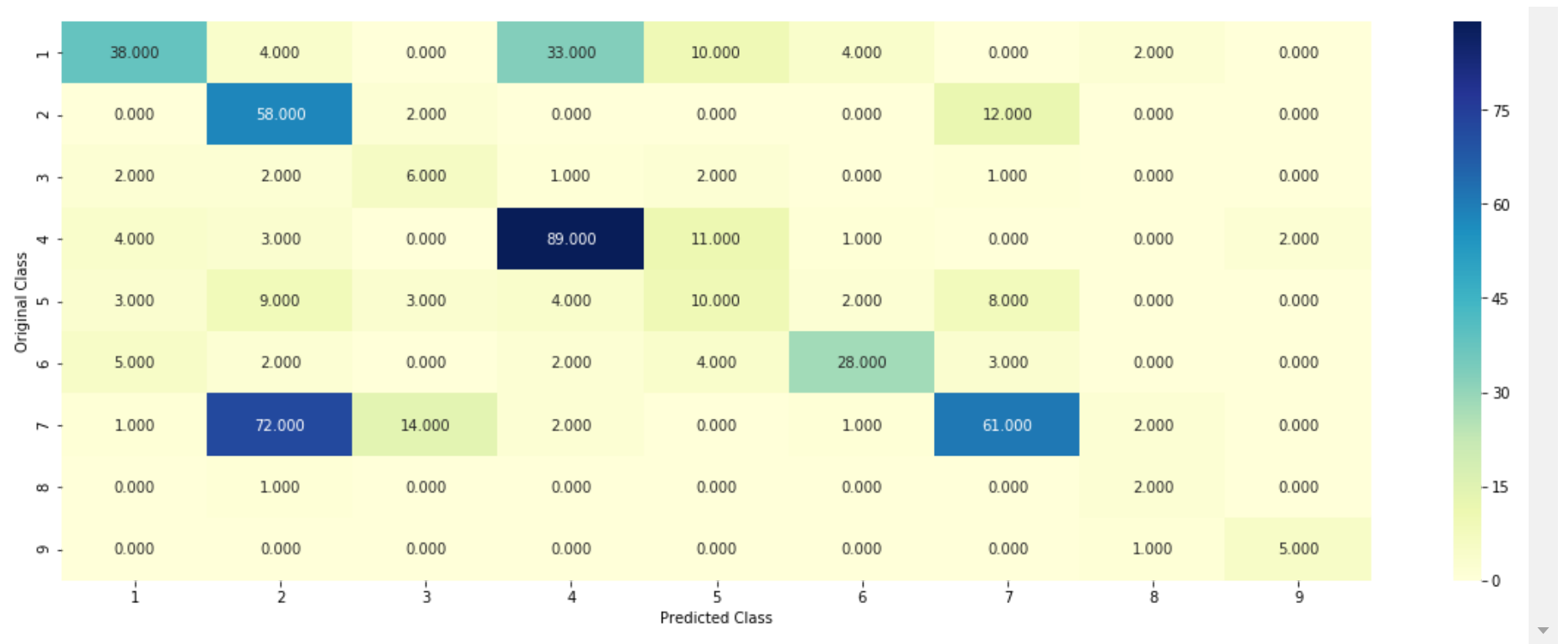
video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-constructi

clf = RandomForestClassifier(max_depth=max_depth[int(best_alpha%4)], n_estimators=alpha[int(best_alpha/4)], criterion='g
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y,cv_x_responseCoding,cv_y, clf)
```

Log loss : 1.2833449142538584

Number of mis-classified points : 0.4417293233082707

----- Confusion matrix -----



----- Precision matrix (Column Sum=1) -----





#### 4.5.5. Feature Importance

##### 4.5.5.1. Correctly Classified point



```

In [95]: clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_depth=max_depth[int(best_alpha
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

test_point_index = 100
no_feature = 27
predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_responseCoding[test_point_index].reshape(1
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
for i in indices:
 if i<9:
 print("Gene is important feature")
 elif i<18:
 print("Variation is important feature")
 else:
 print("Text is important feature")

```

```

Predicted Class : 4
Predicted Class Probabilities: [[0.2112 0.018 0.0953 0.566 0.026 0.0397 0.0064 0.0152 0.0222]]
Actual Class : 4

Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Text is important feature
Variation is important feature
Text is important feature
Gene is important feature
Text is important feature
Text is important feature
Text is important feature
Gene is important feature

```

Variation is important feature  
Gene is important feature  
Text is important feature  
Gene is important feature  
Variation is important feature  
Gene is important feature  
Text is important feature  
Gene is important feature  
Text is important feature  
Gene is important feature  
Variation is important feature  
Text is important feature  
Gene is important feature

#### **4.5.5.2. Incorrectly Classified point**

```
In [94]: test_point_index = 201
predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_responseCoding[test_point_index].reshape(1,-1)).ravel(), 5))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
for i in indices:
 if i<9:
 print("Gene is important feature")
 elif i<18:
 print("Variation is important feature")
 else:
 print("Text is important feature")
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0236 0.1286 0.2284 0.0252 0.0316 0.0582 0.3763 0.0829 0.0452]]

Actual Class : 3

```

Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Text is important feature
Variation is important feature
Text is important feature
Gene is important feature
Text is important feature
Text is important feature
Text is important feature
Gene is important feature
Variation is important feature
Gene is important feature
Text is important feature
Gene is important feature
Variation is important feature
Gene is important feature
Text is important feature
```

Gene is important feature  
Text is important feature  
Gene is important feature  
Variation is important feature  
Text is important feature  
Gene is important feature

## **4.7 Stack the models**

### **4.7.1 testing with hyper parameter tuning**

```

In [96]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier

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: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----

read more about support vector machines with linear kernals here http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC

default parameters
SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_state=None)

Some of methods of SVM()
fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
predict(X) Perform classification on samples in X.

video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/

read more about support vector machines with linear kernals here http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier

default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
class_weight=None)

Some of methods of RandomForestClassifier()
fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
predict(X) Perform classification on samples in X.

```

```

predict_proba (X) Perform classification on samples in X.

some of attributes of RandomForestClassifier()
feature_importances_ : array of shape = [n_features]
The feature importances (the higher, the more important the feature).

video link: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-constructi

clf1 = SGDClassifier(alpha=0.001, penalty='l2', loss='log', class_weight='balanced', random_state=0)
clf1.fit(train_x_onehotCoding, train_y)
sig_clf1 = CalibratedClassifierCV(clf1, method="sigmoid")

clf2 = SGDClassifier(alpha=1, penalty='l2', loss='hinge', class_weight='balanced', random_state=0)
clf2.fit(train_x_onehotCoding, train_y)
sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")

clf3 = MultinomialNB(alpha=0.001)
clf3.fit(train_x_onehotCoding, train_y)
sig_clf3 = CalibratedClassifierCV(clf3, method="sigmoid")

sig_clf1.fit(train_x_onehotCoding, train_y)
print("Logistic Regression : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf1.predict_proba(cv_x_onehotCoding))))
sig_clf2.fit(train_x_onehotCoding, train_y)
print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf2.predict_proba(cv_x_onehotCoding))))
sig_clf3.fit(train_x_onehotCoding, train_y)
print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predict_proba(cv_x_onehotCoding))))
print("-"*50)
alpha = [0.0001,0.001,0.01,0.1,1,10]
best_alpha = 999
for i in alpha:
 lr = LogisticRegression(C=i)
 sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probabilities=True)
 sclf.fit(train_x_onehotCoding, train_y)
 print("Stacking Classifier : for the value of alpha: %f Log Loss: %0.3f" % (i, log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))))
 log_error = log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))
 if best_alpha > log_error:
 best_alpha = log_error

```



```
Logistic Regression : Log Loss: 1.08
Support vector machines : Log Loss: 1.52
Naive Bayes : Log Loss: 1.28
```

```

Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.178
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.035
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.494
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.078
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.133
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.340
```

#### **4.7.2 testing the model with the best hyper parameters**

```
In [97]: lr = LogisticRegression(C=0.1)
sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probabilities=True)
sclf.fit(train_x_onehotCoding, train_y)

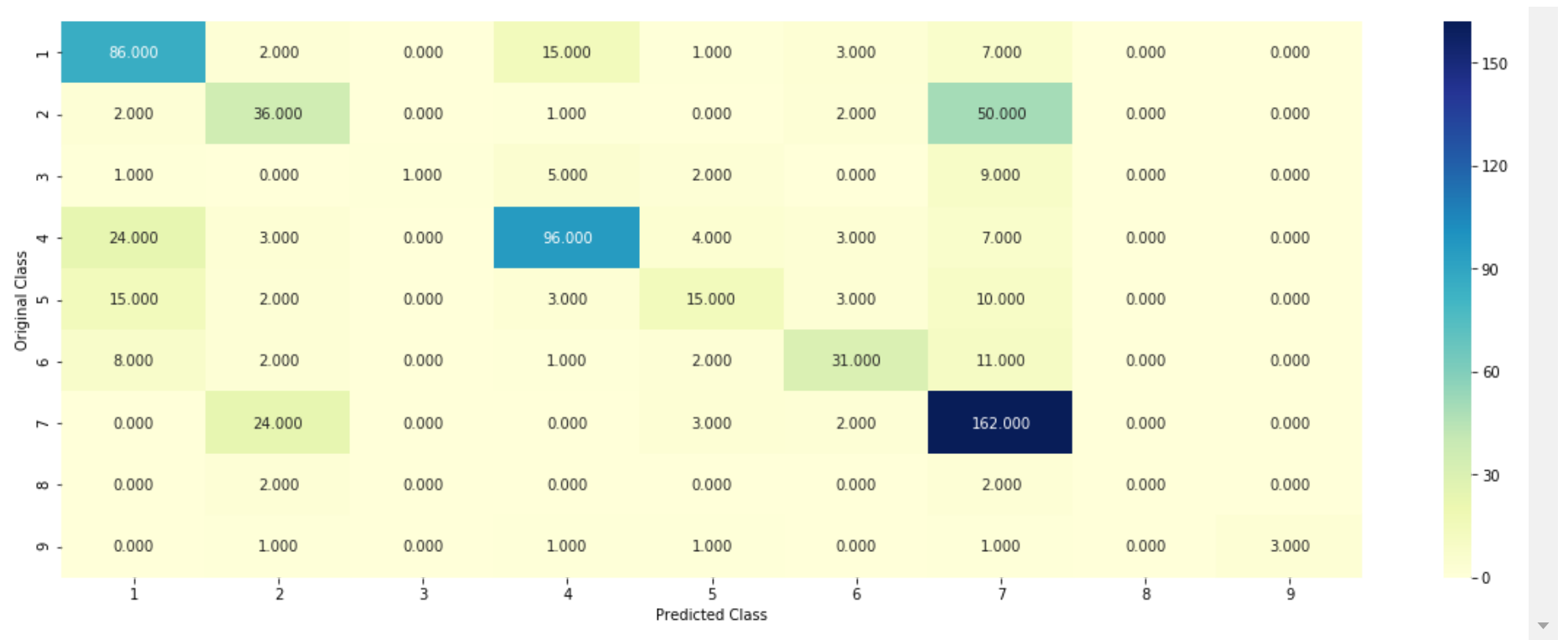
log_error = log_loss(train_y, sclf.predict_proba(train_x_onehotCoding))
print("Log loss (train) on the stacking classifier :",log_error)

log_error = log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))
print("Log loss (CV) on the stacking classifier :",log_error)

log_error = log_loss(test_y, sclf.predict_proba(test_x_onehotCoding))
print("Log loss (test) on the stacking classifier :",log_error)

print("Number of missclassified point :", np.count_nonzero((sclf.predict(test_x_onehotCoding)- test_y))/test_y.shape[0])
plot_confusion_matrix(test_y=test_y, predict_y=sclf.predict(test_x_onehotCoding))
```

```
Log loss (train) on the stacking classifier : 0.6474394636512784
Log loss (CV) on the stacking classifier : 1.077948983414518
Log loss (test) on the stacking classifier : 1.1358631666055556
Number of missclassified point : 0.3533834586466165
----- Confusion matrix -----
```



----- Precision matrix (Column Sum=1) -----



#### 4.7.3 Maximum Voting classifier

```
In [98]: #Refer:http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingClassifier.html
from sklearn.ensemble import VotingClassifier
vclf = VotingClassifier(estimators=[('lr', sig_clf1), ('svc', sig_clf2), ('rf', sig_clf3)], voting='soft')
vclf.fit(train_x_onehotCoding, train_y)
print("Log loss (train) on the VotingClassifier :", log_loss(train_y, vclf.predict_proba(train_x_onehotCoding)))
print("Log loss (CV) on the VotingClassifier :", log_loss(cv_y, vclf.predict_proba(cv_x_onehotCoding)))
print("Log loss (test) on the VotingClassifier :", log_loss(test_y, vclf.predict_proba(test_x_onehotCoding)))
print("Number of missclassified point :", np.count_nonzero((vclf.predict(test_x_onehotCoding)- test_y))/test_y.shape[0])
plot_confusion_matrix(test_y=test_y, predict_y=vclf.predict(test_x_onehotCoding))
```

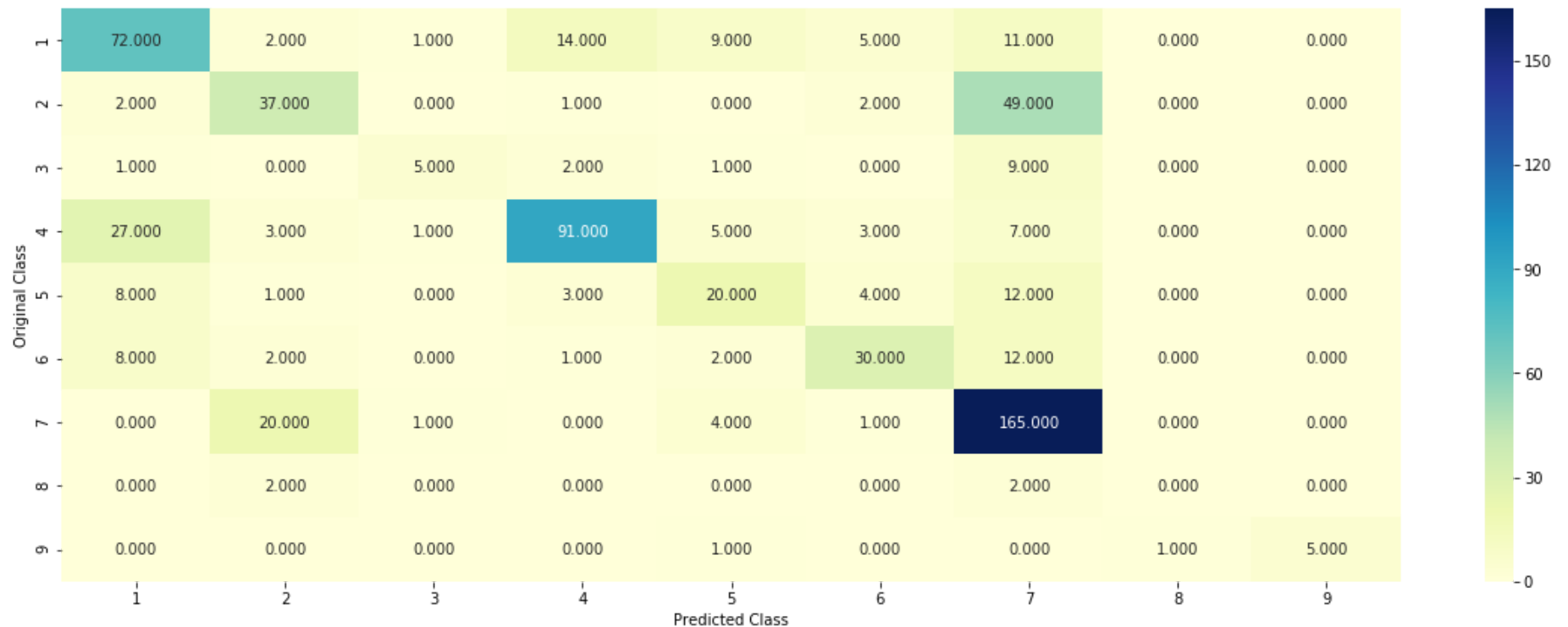
Log loss (train) on the VotingClassifier : 0.8742064739568054

Log loss (CV) on the VotingClassifier : 1.1501179936812944

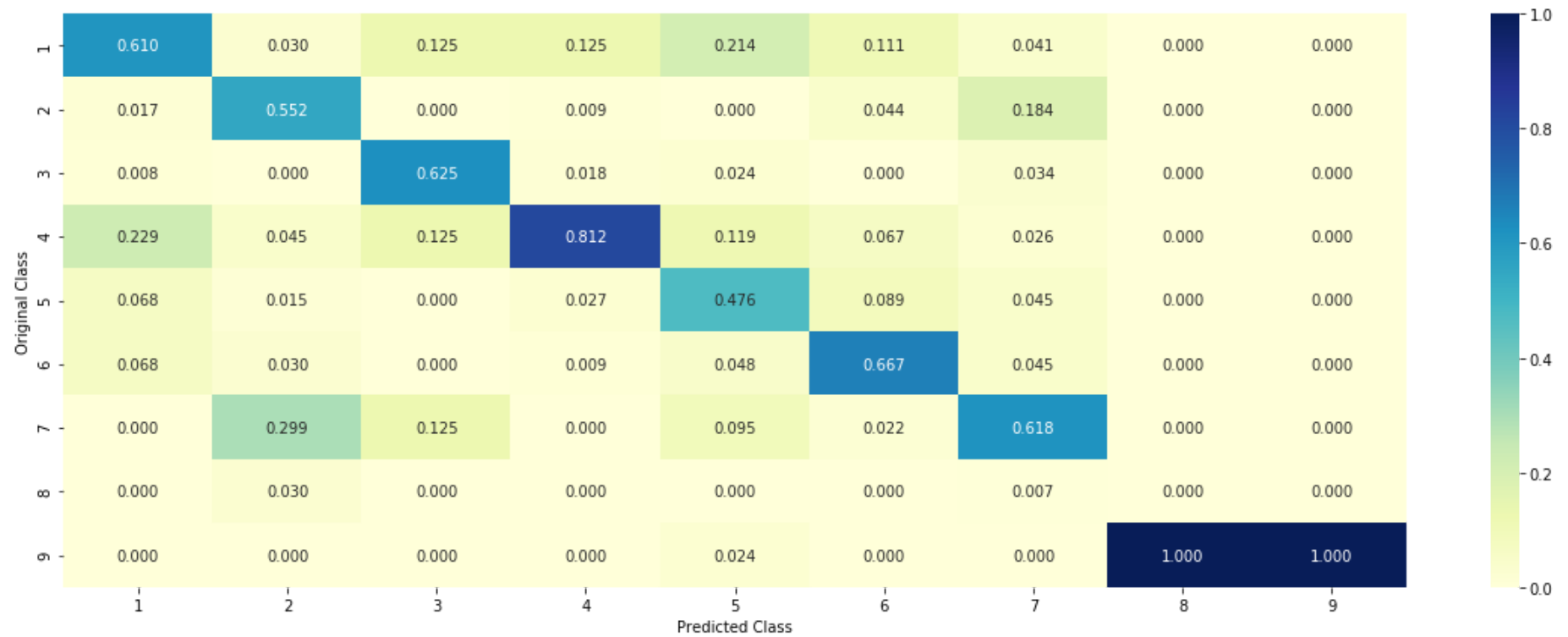
Log loss (test) on the VotingClassifier : 1.1831239143567645

Number of missclassified point : 0.3609022556390977

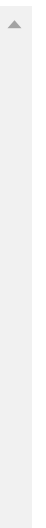
----- Confusion matrix -----

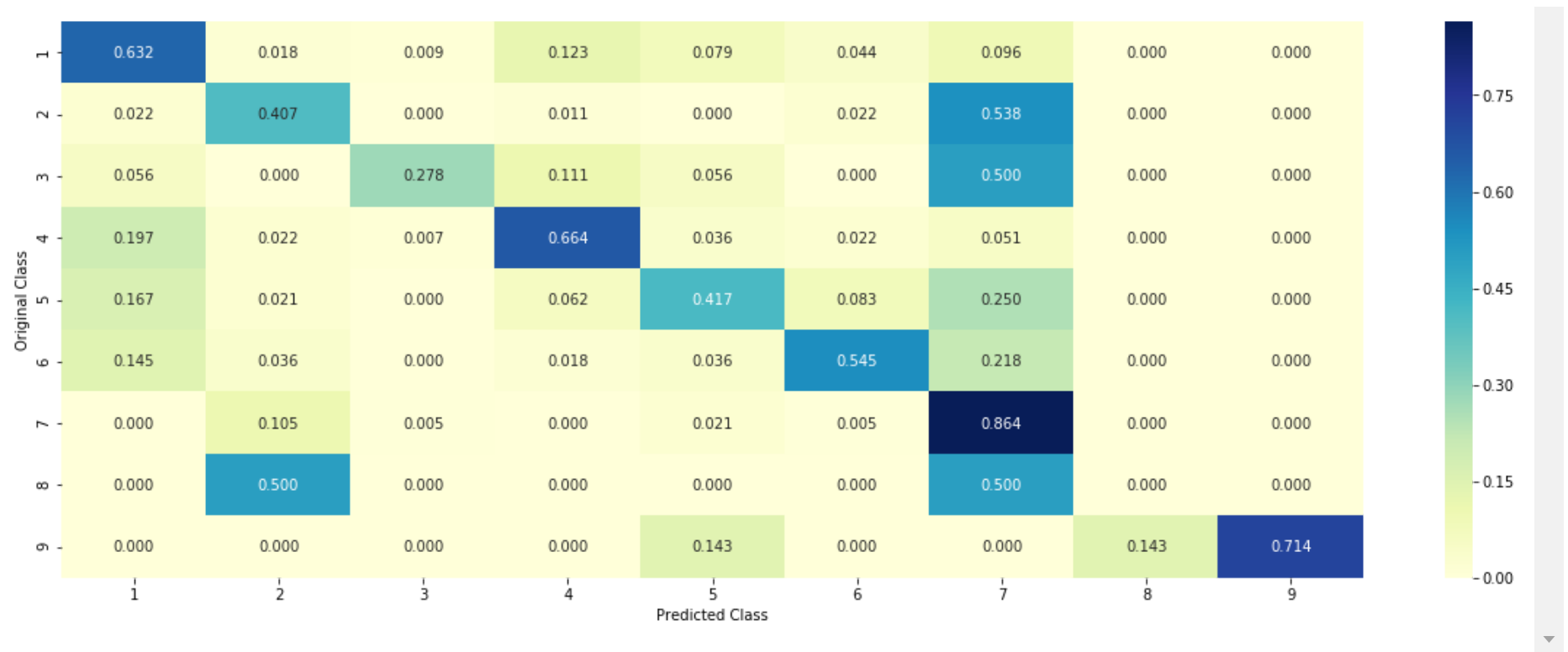


----- Precision matrix (Column Sum=1) -----



----- Recall matrix (Row sum=1) -----





## Comparing the scores from all the above models

Performance of various model on replacing Countvectorizer with TFIDF Vectorizer

```
In [101]: # http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
x.field_names = ["Model-Name", "Train loss", "CV loss", "Test Loss", "% Misclassified"]
x.add_row(["Naive Bayes", "0.944", "1.223", "1.237", "41.16"])
x.add_row(["KNN", "0.647", "1.042", "1.078", "36.27"])
x.add_row(["LR with class Balancing", "0.605", "1.072", "1.083", "31.20"])
x.add_row(["LR without class Balancing", "0.605", "1.106", "1.112", "31.01"])
x.add_row(["Linear SVM", "0.727", "1.107", "1.131", "32.89"])
x.add_row(["RF(One Hot Encoding)", "0.654", "1.175", "1.151", "38.34"])
x.add_row(["RF(Response coding)", "0.05", "1.283", "1.329", "44.17"])
x.add_row(["Stacking Classifier", "0.647", "1.017", "1.135", "35.33"])
x.add_row(["Max voting classifier", "0.874", "1.150", "1.183", "36.09"])
print(x)
```

| Model-Name                 | Train loss | CV loss | Test Loss | % Misclassified |
|----------------------------|------------|---------|-----------|-----------------|
| Naive Bayes                | 0.944      | 1.223   | 1.237     | 41.16           |
| KNN                        | 0.647      | 1.042   | 1.078     | 36.27           |
| LR with class Balancing    | 0.605      | 1.072   | 1.083     | 31.20           |
| LR without class Balancing | 0.605      | 1.106   | 1.112     | 31.01           |
| Linear SVM                 | 0.727      | 1.107   | 1.131     | 32.89           |
| RF(One Hot Encoding)       | 0.654      | 1.175   | 1.151     | 38.34           |
| RF(Response coding)        | 0.05       | 1.283   | 1.329     | 44.17           |
| Stacking Classifier        | 0.647      | 1.017   | 1.135     | 35.33           |
| Max voting classifier      | 0.874      | 1.150   | 1.183     | 36.09           |

## Conclusion / Observations :-

1. We got the minimum test loss with KNN which is 1.078 but the percentage of misclassified points are on the higher side being 36.27%.
2. We got the minimum misclassified points with Logistic Regression without class balancing with 31.01% with test log loss of 1.112.



3. Random Forest with Response coding gave the maximum misclassified points with 44.17% of misclassified points. Also there is severe overfitting with this model where the train loss is 0.05 while the test loss is 1.329 which is highest amongst all.
4. So for Task-1 where we replaced Countvectorizer with TFIDF vectorizer Logistic Regression without class balancing gave us the best result.