

Personalized cancer diagnosis

Task 3

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>)

2. <https://www.youtube.com/watch?v=UwbuW7oK8rk> (<https://www.youtube.com/watch?v=UwbuW7oK8rk>)
3. <https://www.youtube.com/watch?v=qxXRKVompl8> (<https://www.youtube.com/watch?v=qxXRKVompl8>)

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 [1]: 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
```

C:\Users\Himanshu Pc\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was de

precated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

C:\Users\Himanshu Pc\Anaconda3\lib\site-packages\sklearn\ensemble\weight_boosting.py:29: DeprecationWarning: numpy.core.umath_tests is an internal NumPy module and should not be imported. It will be removed in a future NumPy release.

```
from numpy.core.umath_tests import inner1d
```

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

```
In [2]: 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[2]:

	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 [3]: # 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[3]:

	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 [4]: 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[4]: True

3.1.3. Preprocessing of text

```

In [5]: # 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 [6]: #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 : 121.01383429999987 seconds

```



```
In [7]: #merging both gene_variations and text data based on ID
result = pd.merge(data, data_text,on='ID', how='left')
result.head()
```

Out[7]:

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

Out[8]:

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

```
In [10]: result[result['ID']==1109]
```

Out[10]:

	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 [11]: 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 [12]: 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 [13]: # 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()

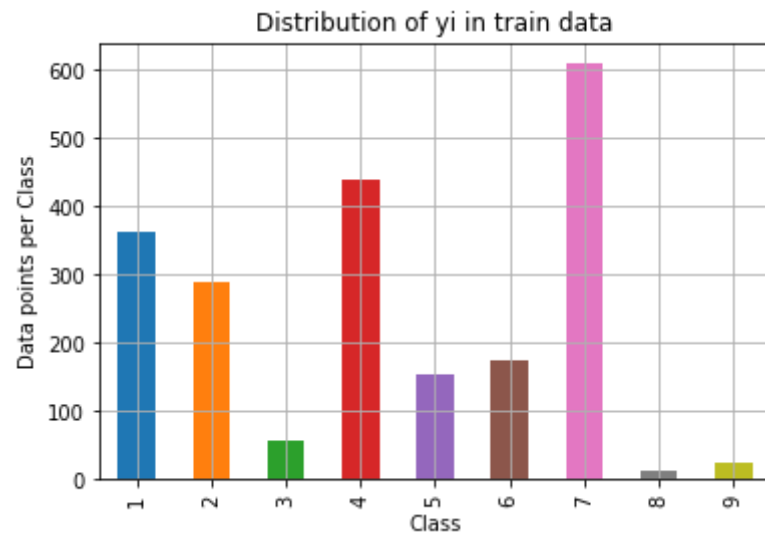
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

```

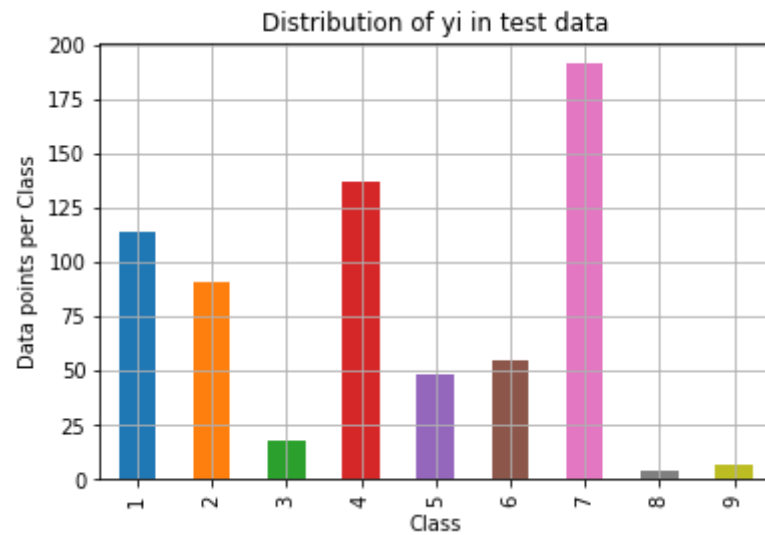
```
plt.grid()
plt.show()

# 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
```

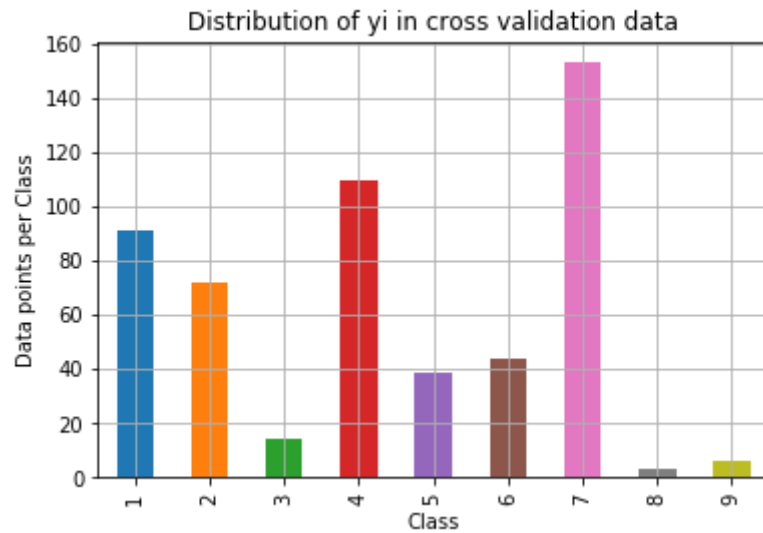


Number of data points in class 7 : 609 (28.672 %)
 Number of data points in class 4 : 439 (20.669 %)

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 [14]: # 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 [15]: # 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))

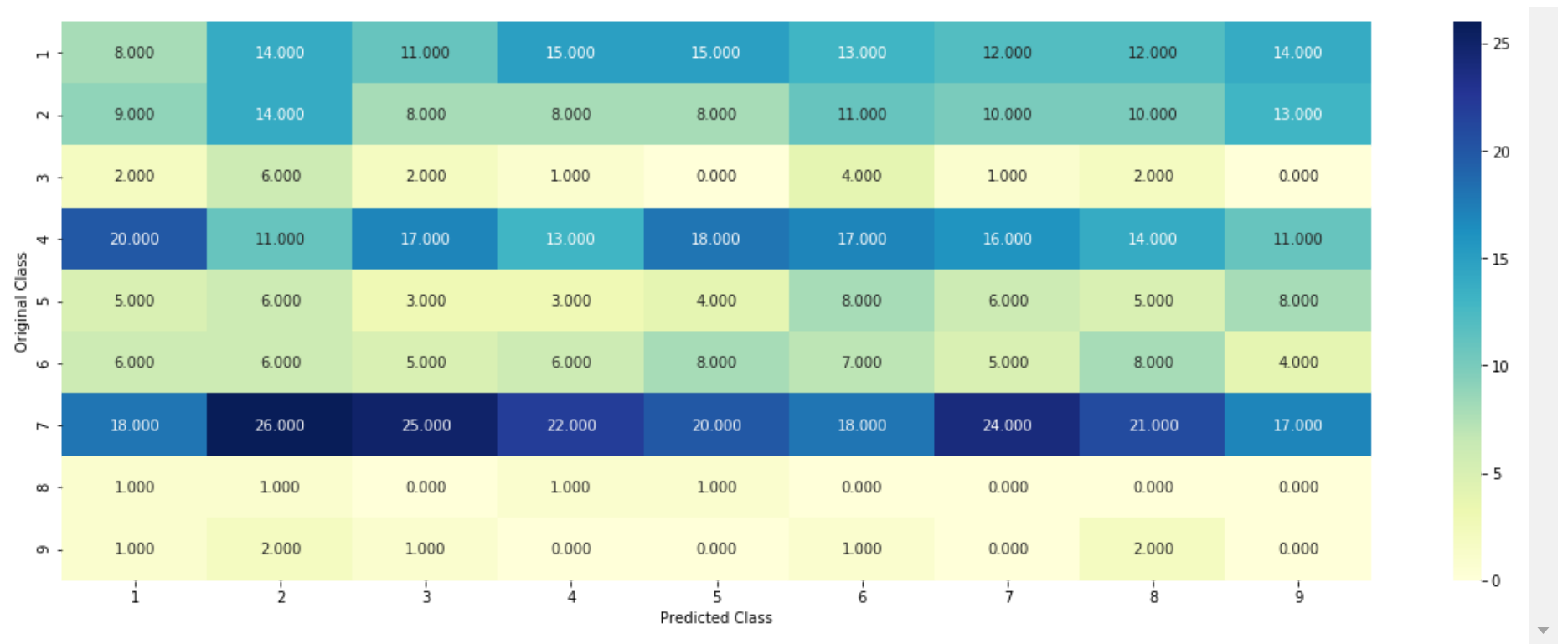
predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

```

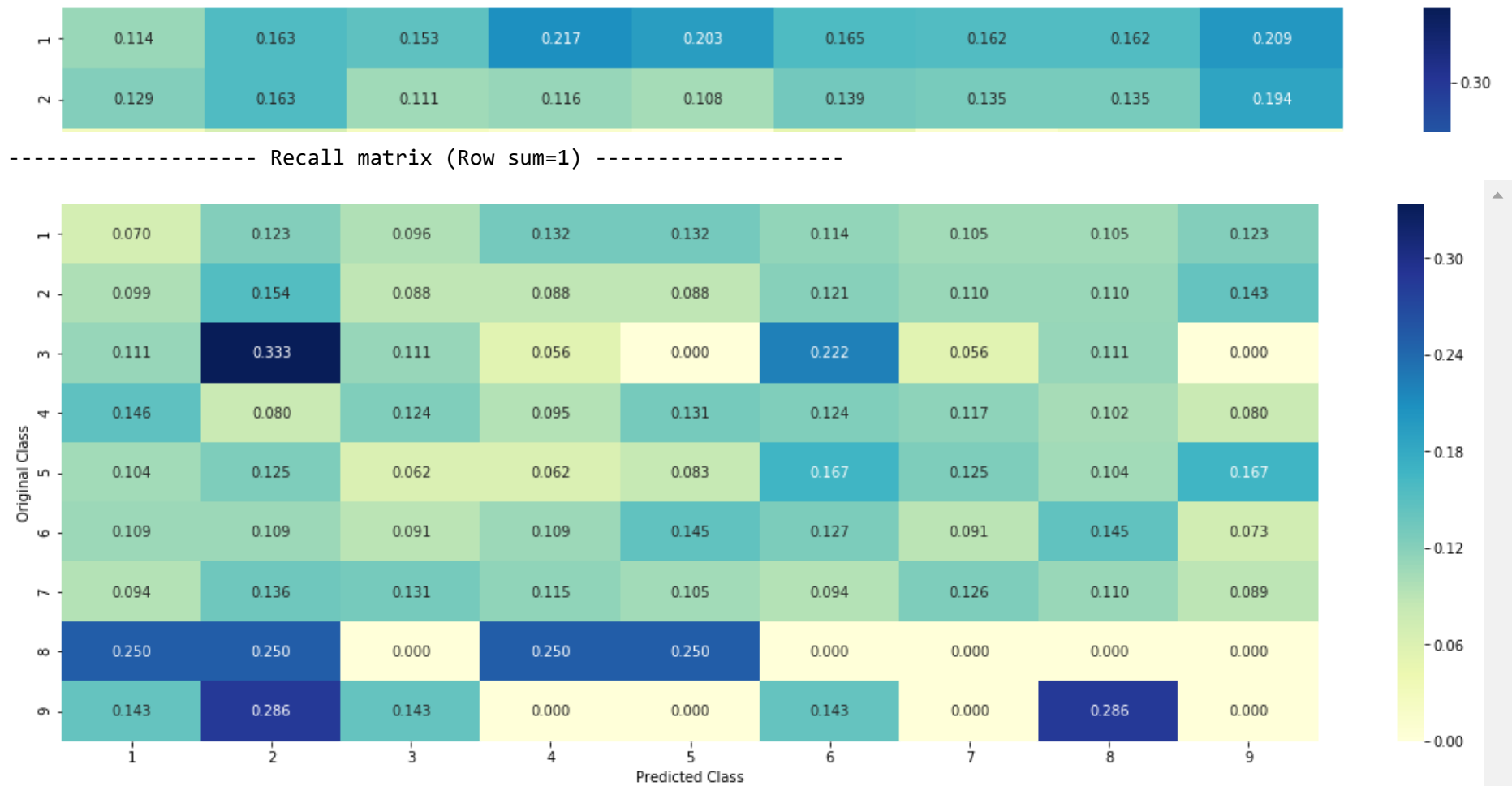
Log loss on Cross Validation Data using Random Model 2.4700741925010603

Log loss on Test Data using Random Model 2.4707308127024215

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



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



3.3 Univariate Analysis

```

In [16]: # 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, 0.0933333333
#      ...
#      }
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?

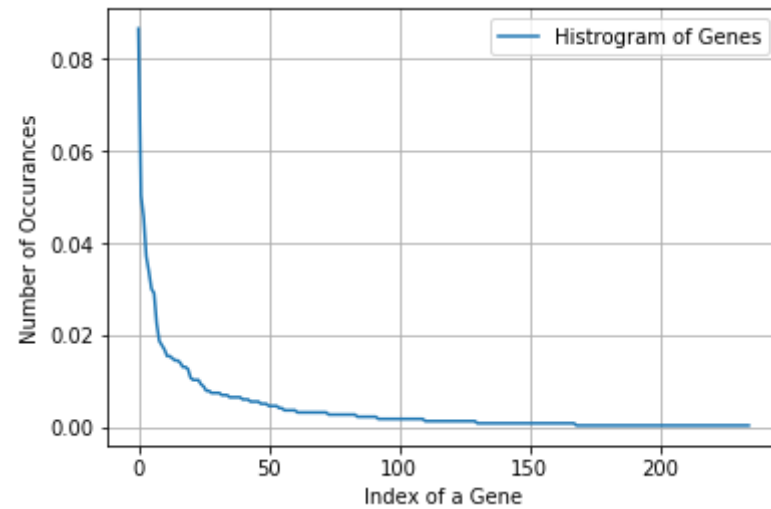
```
In [17]: unique_genes = train_df['Gene'].value_counts()
print('Number of Unique Genes :', unique_genes.shape[0])
# the top 10 genes that occurred most
print(unique_genes.head(10))
```

```
Number of Unique Genes : 235
BRCA1      184
TP53       107
EGFR        97
PTEN        79
BRCA2       72
KIT         64
BRAF        62
ERBB2       48
PDGFRA      40
PIK3CA      38
Name: Gene, dtype: int64
```

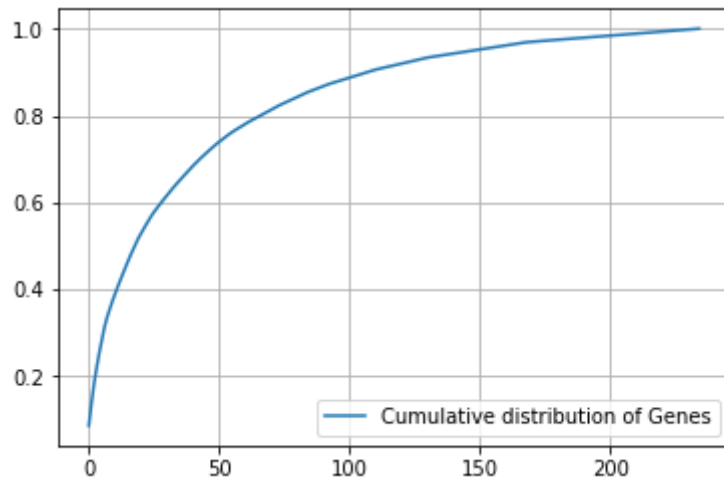
```
In [18]: print("Ans: There are", unique_genes.shape[0] , "different categories of genes in the train data, and they are distributed
```

Ans: There are 235 different categories of genes in the train data, and they are distributed as follows


```
In [19]: 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 [20]: 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 [21]: #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 [22]: 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 [23]: # one-hot encoding of Gene feature.  
from sklearn.feature_extraction.text import TfidfVectorizer  
gene_vectorizer = CountVectorizer()  
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 [24]: train_df['Gene'].head()
```

```
Out[24]: 2070    TET2  
1053    TSC2  
1572    ALK  
1958    ATM  
1533    ALK  
Name: Gene, dtype: object
```

```
In [25]: gene_vectorizer.get_feature_names()
```

```
'bcor',  
'braf',  
'brca1',  
'brca2',  
'brd4',  
'brip1',  
'btk',  
'card11',  
'carm1',  
'casp8',  
'cb1',  
'ccnd1',  
'ccnd2',  
'ccnd3',  
'ccne1',  
'cdh1',  
  
'cdk12',  
'cdk4',  
'cdk6',  
...
```

```
In [26]: 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, 234)
```

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 [27]: 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)

```

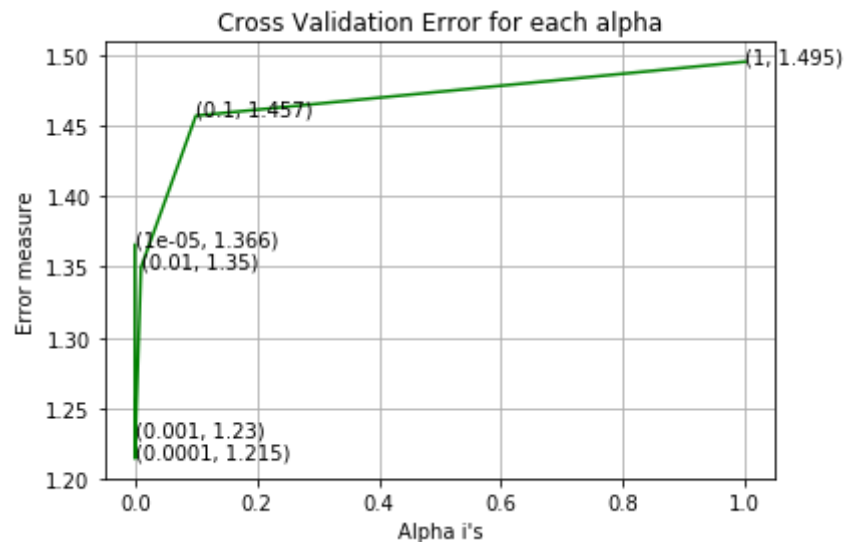
```

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.c
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, la
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.c

```

For values of alpha = 1e-05 The log loss is: 1.3655905574728293
 For values of alpha = 0.0001 The log loss is: 1.2145105032283185
 For values of alpha = 0.001 The log loss is: 1.2299353809814233
 For values of alpha = 0.01 The log loss is: 1.349613574883199
 For values of alpha = 0.1 The log loss is: 1.4569790379878775
 For values of alpha = 1 The log loss is: 1.4951698164284835



For values of best alpha = 0.0001 The train log loss is: 1.0395717753275246
For values of best alpha = 0.0001 The cross validation log loss is: 1.2145105032283185
For values of best alpha = 0.0001 The test log loss is: 1.2329495396074415

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 [28]: 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 235 genes in train dataset?

Ans

1. In test data 652 out of 665 : 98.04511278195488
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?

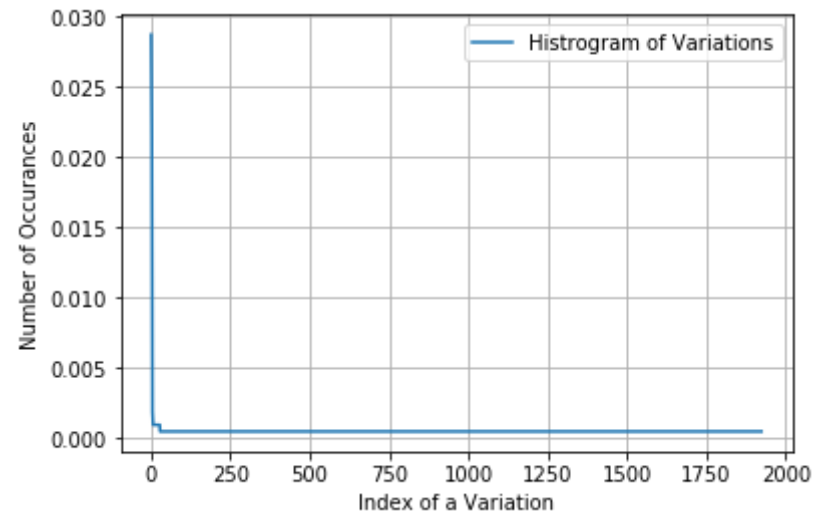
```
In [29]: unique_variations = train_df['Variation'].value_counts()
print('Number of Unique Variations :', unique_variations.shape[0])
# the top 10 variations that occurred most
print(unique_variations.head(10))
```

```
Number of Unique Variations : 1924
Truncating_Mutations      61
Deletion                  52
Amplification             46
Fusions                   18
Overexpression            4
G12V                      3
G13D                      2
R170W                     2
C618R                     2
T58I                      2
Name: Variation, dtype: int64
```

```
In [30]: print("Ans: There are", unique_variations.shape[0] , "different categories of variations in the train data, and they are
```

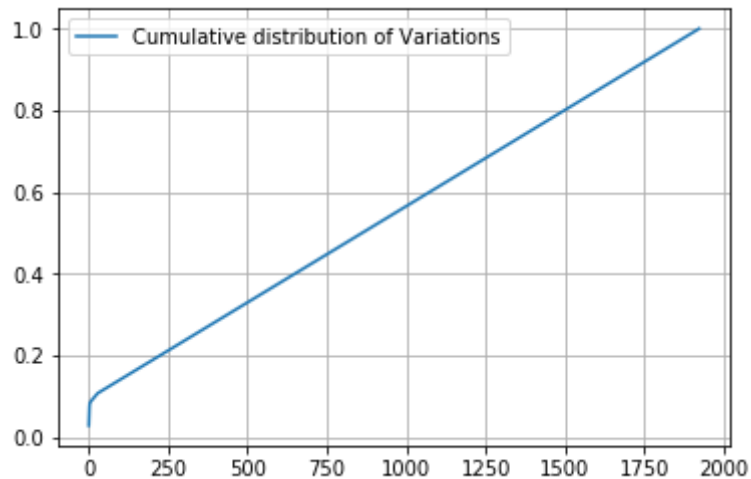
Ans: There are 1924 different categories of variations in the train data, and they are distributed as follows


```
In [31]: 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 [32]: c = np.cumsum(h)
print(c)
plt.plot(c,label='Cumulative distribution of Variations')
plt.grid()
plt.legend()
plt.show()

[0.0287194  0.05320151 0.07485876 ... 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 [33]: # 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 [34]: 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 [35]: # one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer()
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 [36]: 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, 1958)

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

Let's build a model just like the earlier!

```
In [37]: 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()
```

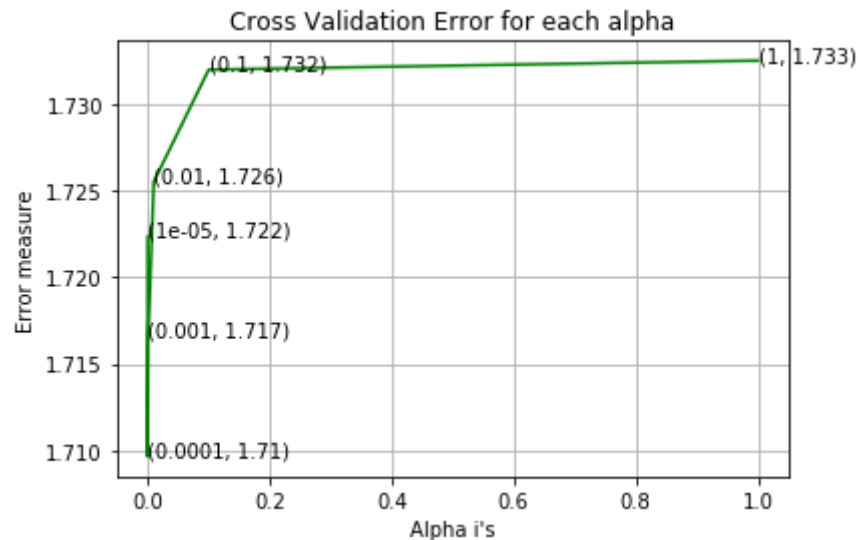
```

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: 1.7223483779532718
 For values of alpha = 0.0001 The log loss is: 1.7096403650875167
 For values of alpha = 0.001 The log loss is: 1.7165998187299627
 For values of alpha = 0.01 The log loss is: 1.7255156082696825
 For values of alpha = 0.1 The log loss is: 1.7319859875934729
 For values of alpha = 1 The log loss is: 1.7325091506556003



For values of best alpha = 0.0001 The train log loss is: 0.7686267185056543
For values of best alpha = 0.0001 The cross validation log loss is: 1.7096403650875167
For values of best alpha = 0.0001 The test log loss is: 1.704671483416769

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 [38]: print("Q12. How many data points are covered by total ", unique_variations.shape[0], " genes in test and cross validation")
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 1924 genes in test and cross validation data sets?

Ans

1. In test data 64 out of 665 : 9.624060150375941
2. In cross validation data 59 out of 532 : 11.090225563909774

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 [39]: # 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 [40]: 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 [41]: # building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = CountVectorizer(ngram_range = (1, 2), min_df=3)
train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])
test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
cv_text_feature_onehotCoding = text_vectorizer.transform(cv_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 : 763311


```

In [42]: 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 [43]: # 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 [44]: #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()))
```

```
In [45]: # Number of words for a given frequency.
print(Counter(sorted_text_occur))
```

```
Counter({3: 145234, 4: 97729, 6: 63840, 5: 63374, 7: 45612, 8: 36196, 9: 35306, 10: 25736, 11: 24468, 12: 18799, 14: 15035, 15: 14878, 13: 14686, 16: 10682, 18: 7372, 20: 7136, 19: 6998, 17: 6938, 22: 5177, 21: 5125, 24: 4968, 28: 4861, 27: 4503, 23: 4440, 38: 3567, 25: 3392, 26: 3107, 32: 3087, 30: 2856, 29: 2693, 31: 2570, 54: 2312, 33: 2254, 44: 2194, 34: 1843, 36: 1811, 39: 1762, 35: 1760, 40: 1647, 37: 1489, 51: 1367, 42: 1359, 41: 1347, 69: 1332, 45: 1266, 43: 1155, 48: 1103, 46: 1084, 47: 931, 49: 892, 56: 868, 50: 861, 55: 845, 52: 843, 53: 780, 60: 748, 57: 707, 58: 672, 59: 647, 64: 639, 62: 625, 70: 615, 63: 609, 66: 596, 61: 575, 65: 556, 72: 545, 88: 522, 71: 489, 68: 479, 76: 434, 73: 428, 67: 425, 75: 421, 78: 391, 77: 390, 80: 389, 74: 386, 79: 344, 81: 342, 84: 326, 83: 310, 89: 308, 82: 308, 86: 304, 90: 302, 108: 294, 85: 292, 87: 290, 91: 287, 93: 278, 92: 272, 96: 268, 94: 263, 102: 256, 97: 238, 98: 234, 100: 223, 104: 222, 95: 213, 99: 199, 110: 195, 103: 192, 114: 189, 101: 187, 107: 182, 111: 171, 109: 171, 132: 167, 113: 166, 112: 164, 105: 161, 106: 160, 120: 158, 117: 158, 115: 155, 124: 153, 118: 150, 134: 147, 123: 147, 116: 142, 138: 141, 136: 140, 126: 139, 119: 135, 122: 134, 121: 131, 128: 128, 129: 126, 131: 124, 135: 123, 142: 122, 145: 118, 125: 118, 127: 116, 155: 115, 152: 114, 139: 114, 133: 114, 130: 114, 148: 110, 146: 107, 162: 106, 140: 106, 141: 104, 144: 102, 143: 100, 137: 99, 156: 96, 151: 96, 150: 96, 153: 95, 176: 91, 168: 89, 160: 89, 154: 86, 174: 83, 157: 82, 149: 82, 166: 81, 147: 81, 161: 80, 159: 78, 192: 77, 164: 77, 169: 75, 165: 74, 158: 73, 171: 72, 163: 70, 172: 69, 170: 68, 173: 67, 199: 66, 175: 66, 190: 65, 184: 64, 187: 63, 167: 62, 212: 59, 195: 59, 200: 58, 183: 58, 182: 57, 179: 56, 178: 56, 209: 55, 202: 55, 191: 55, 188: 55, 186: 55, 177: 55, 205: 54, 216: 53, 203: 52, 194: 52, 185: 52, 222: 51, 215: 51, 181: 51, 224: 50, 194: 50, 193: 50, 189: 50, 210: 49, 208: 49, 198: 49, 270: 48, 226: 48, 213: 48, 207: 48, 228: 47, 217: 47, 206: 47, 180: 47, 196: 46, 231: 45, 211: 45, 226: 44, 201: 44, 239: 43, 230: 43, 204: 43, 246: 41, 240: 41, 235: 41, 225: 41, 220: 41, 219: 41, 214: 41, 242: 40, 257: 39, 248: 39, 221: 38, 250: 38, 243: 37, 242: 37, 241: 37, 234: 37, 233: 37, 232: 37, 231: 37, 230: 37, 229: 37, 228: 37, 227: 37, 226: 37, 225: 37, 224: 37, 223: 37, 222: 37, 221: 37, 220: 37, 219: 37, 218: 37, 217: 37, 216: 37, 215: 37, 214: 37, 213: 37, 212: 37, 211: 37, 210: 37, 209: 37, 208: 37, 207: 37, 206: 37, 205: 37, 204: 37, 203: 37, 202: 37, 201: 37, 200: 37, 199: 37, 198: 37, 197: 37, 196: 37, 195: 37, 194: 37, 193: 37, 192: 37, 191: 37, 190: 37, 189: 37, 188: 37, 187: 37, 186: 37, 185: 37, 184: 37, 183: 37, 182: 37, 181: 37, 180: 37, 179: 37, 178: 37, 177: 37, 176: 37, 175: 37, 174: 37, 173: 37, 172: 37, 171: 37, 170: 37, 169: 37, 168: 37, 167: 37, 166: 37, 165: 37, 164: 37, 163: 37, 162: 37, 161: 37, 160: 37, 159: 37, 158: 37, 157: 37, 156: 37, 155: 37, 154: 37, 153: 37, 152: 37, 151: 37, 150: 37, 149: 37, 148: 37, 147: 37, 146: 37, 145: 37, 144: 37, 143: 37, 142: 37, 141: 37, 140: 37, 139: 37, 138: 37, 137: 37, 136: 37, 135: 37, 134: 37, 133: 37, 132: 37, 131: 37, 130: 37, 129: 37, 128: 37, 127: 37, 126: 37, 125: 37, 124: 37, 123: 37, 122: 37, 121: 37, 120: 37, 119: 37, 118: 37, 117: 37, 116: 37, 115: 37, 114: 37, 113: 37, 112: 37, 111: 37, 110: 37, 109: 37, 108: 37, 107: 37, 106: 37, 105: 37, 104: 37, 103: 37, 102: 37, 101: 37, 100: 37, 99: 37, 98: 37, 97: 37, 96: 37, 95: 37, 94: 37, 93: 37, 92: 37, 91: 37, 90: 37, 89: 37, 88: 37, 87: 37, 86: 37, 85: 37, 84: 37, 83: 37, 82: 37, 81: 37, 80: 37, 79: 37, 78: 37, 77: 37, 76: 37, 75: 37, 74: 37, 73: 37, 72: 37, 71: 37, 70: 37, 69: 37, 68: 37, 67: 37, 66: 37, 65: 37, 64: 37, 63: 37, 62: 37, 61: 37, 60: 37, 59: 37, 58: 37, 57: 37, 56: 37, 55: 37, 54: 37, 53: 37, 52: 37, 51: 37, 50: 37, 49: 37, 48: 37, 47: 37, 46: 37, 45: 37, 44: 37, 43: 37, 42: 37, 41: 37, 40: 37, 39: 37, 38: 37, 37: 37, 36: 37, 35: 37, 34: 37, 33: 37, 32: 37, 31: 37, 30: 37, 29: 37, 28: 37, 27: 37, 26: 37, 25: 37, 24: 37, 23: 37, 22: 37, 21: 37, 20: 37, 19: 37, 18: 37, 17: 37, 16: 37, 15: 37, 14: 37, 13: 37, 12: 37, 11: 37, 10: 37, 9: 37, 8: 37, 7: 37, 6: 37, 5: 37, 4: 37, 3: 37, 2: 37, 1: 37, 0: 37, -1: 37, -2: 37, -3: 37, -4: 37, -5: 37, -6: 37, -7: 37, -8: 37, -9: 37, -10: 37, -11: 37, -12: 37, -13: 37, -14: 37, -15: 37, -16: 37, -17: 37, -18: 37, -19: 37, -20: 37, -21: 37, -22: 37, -23: 37, -24: 37, -25: 37, -26: 37, -27: 37, -28: 37, -29: 37, -30: 37, -31: 37, -32: 37, -33: 37, -34: 37, -35: 37, -36: 37, -37: 37, -38: 37, -39: 37, -40: 37, -41: 37, -42: 37, -43: 37, -44: 37, -45: 37, -46: 37, -47: 37, -48: 37, -49: 37, -50: 37, -51: 37, -52: 37, -53: 37, -54: 37, -55: 37, -56: 37, -57: 37, -58: 37, -59: 37, -60: 37, -61: 37, -62: 37, -63: 37, -64: 37, -65: 37, -66: 37, -67: 37, -68: 37, -69: 37, -70: 37, -71: 37, -72: 37, -73: 37, -74: 37, -75: 37, -76: 37, -77: 37, -78: 37, -79: 37, -80: 37, -81: 37, -82: 37, -83: 37, -84: 37, -85: 37, -86: 37, -87: 37, -88: 37, -89: 37, -90: 37, -91: 37, -92: 37, -93: 37, -94: 37, -95: 37, -96: 37, -97: 37, -98: 37, -99: 37, -100: 37, -101: 37, -102: 37, -103: 37, -104: 37, -105: 37, -106: 37, -107: 37, -108: 37, -109: 37, -110: 37, -111: 37, -112: 37, -113: 37, -114: 37, -115: 37, -116: 37, -117: 37, -118: 37, -119: 37, -120: 37, -121: 37, -122: 37, -123: 37, -124: 37, -125: 37, -126: 37, -127: 37, -128: 37, -129: 37, -130: 37, -131: 37, -132: 37, -133: 37, -134: 37, -135: 37, -136: 37, -137: 37, -138: 37, -139: 37, -140: 37, -141: 37, -142: 37, -143: 37, -144: 37, -145: 37, -146: 37, -147: 37, -148: 37, -149: 37, -150: 37, -151: 37, -152: 37, -153: 37, -154: 37, -155: 37, -156: 37, -157: 37, -158: 37, -159: 37, -160: 37, -161: 37, -162: 37, -163: 37, -164: 37, -165: 37, -166: 37, -167: 37, -168: 37, -169: 37, -170: 37, -171: 37, -172: 37, -173: 37, -174: 37, -175: 37, -176: 37, -177: 37, -178: 37, -179: 37, -180: 37, -181: 37, -182: 37, -183: 37, -184: 37, -185: 37, -186: 37, -187: 37, -188: 37, -189: 37, -190: 37, -191: 37, -192: 37, -193: 37, -194: 37, -195: 37, -196: 37, -197: 37, -198: 37, -199: 37, -200: 37, -201: 37, -202: 37, -203: 37, -204: 37, -205: 37, -206: 37, -207: 37, -208: 37, -209: 37, -210: 37, -211: 37, -212: 37, -213: 37, -214: 37, -215: 37, -216: 37, -217: 37, -218: 37, -219: 37, -220: 37, -221: 37, -222: 37, -223: 37, -224: 37, -225: 37, -226: 37, -227: 37, -228: 37, -229: 37, -230: 37, -231: 37, -232: 37, -233: 37, -234: 37, -235: 37, -236: 37, -237: 37, -238: 37, -239: 37, -240: 37, -241: 37, -242: 37, -243: 37, -244: 37, -245: 37, -246: 37, -247: 37, -248: 37, -249: 37, -250: 37, -251: 37, -252: 37, -253: 37, -254: 37, -255: 37, -256: 37, -257: 37, -258: 37, -259: 37, -260: 37, -261: 37, -262: 37, -263: 37, -264: 37, -265: 37, -266: 37, -267: 37, -268: 37, -269: 37, -270: 37, -271: 37, -272: 37, -273: 37, -274: 37, -275: 37, -276: 37, -277: 37, -278: 37, -279: 37, -280: 37, -281: 37, -282: 37, -283: 37, -284: 37, -285: 37, -286: 37, -287: 37, -288: 37, -289: 37, -290: 37, -291: 37, -292: 37, -293: 37, -294: 37, -295: 37, -296: 37, -297: 37, -298: 37, -299: 37, -300: 37, -301: 37, -302: 37, -303: 37, -304: 37, -305: 37, -306: 37, -307: 37, -308: 37, -309: 37, -310: 37, -311: 37, -312: 37, -313: 37, -314: 37, -315: 37, -316: 37, -317: 37, -318: 37, -319: 37, -320: 37, -321: 37, -322: 37, -323: 37, -324: 37, -325: 37, -326: 37, -327: 37, -328: 37, -329: 37, -330: 37, -331: 37, -332: 37, -333: 37, -334: 37, -335: 37, -336: 37, -337: 37, -338: 37, -339: 37, -340: 37, -341: 37, -342: 37, -343: 37, -344: 37, -345: 37, -346: 37, -347: 37, -348: 37, -349: 37, -350: 37, -351: 37, -352: 37, -353: 37, -354: 37, -355: 37, -356: 37, -357: 37, -358: 37, -359: 37, -360: 37, -361: 37, -362: 37, -363: 37, -364: 37, -365: 37, -366: 37, -367: 37, -368: 37, -369: 37, -370: 37, -371: 37, -372: 37, -373: 37, -374: 37, -375: 37, -376: 37, -377: 37, -378: 37, -379: 37, -380: 37, -381: 37, -382: 37, -383: 37, -384: 37, -385: 37, -386: 37, -387: 37, -388: 37, -389: 37, -390: 37, -391: 37, -392: 37, -393: 37, -394: 37, -395: 37, -396: 37, -397: 37, -398: 37, -399: 37, -400: 37, -401: 37, -402: 37, -403: 37, -404: 37, -405: 37, -406: 37, -407: 37, -408: 37, -409: 37, -410: 37, -411: 37, -412: 37, -413: 37, -414: 37, -415: 37, -416: 37, -417: 37, -418: 37, -419: 37, -420: 37, -421: 37, -422: 37, -423: 37, -424: 37, -425: 37, -426: 37, -427: 37, -428: 37, -429: 37, -430: 37, -431: 37, -432: 37, -433: 37, -434: 37, -435: 37, -436: 37, -437: 37, -438: 37, -439: 37, -440: 37, -441: 37, -442: 37, -443: 37, -444: 37, -445: 37, -446: 37, -447: 37, -448: 37, -449: 37, -450: 37, -451: 37, -452: 37, -453: 37, -454: 37, -455: 37, -456: 37, -457: 37, -458: 37, -459: 37, -460: 37, -461: 37, -462: 37, -463: 37, -464: 37, -465: 37, -466: 37, -467: 37, -468: 37, -469: 37, -470: 37, -471: 37, -472: 37, -473: 37, -474: 37, -475: 37, -476: 37, -477: 37, -478: 37, -479: 37, -480: 37, -481: 37, -482: 37, -483: 37, -484: 37, -485: 37, -486: 37, -487: 37, -488: 37, -489: 37, -490: 37, -491: 37, -492: 37, -493: 37, -494: 37, -495: 37, -496: 37, -497: 37, -498: 37, -499: 37, -500: 37, -501: 37, -502: 37, -503: 37, -504: 37, -505: 37, -506: 37, -507: 37, -508: 37, -509: 37, -510: 37, -511: 37, -512: 37, -513: 37, -514: 37, -515: 37, -516: 37, -517: 37, -518: 37, -519: 37, -520: 37, -521: 37, -522: 37, -523: 37, -524: 37, -525: 37, -526: 37, -527: 37, -528: 37, -529: 37, -530: 37, -531: 37, -532: 37, -533: 37, -534: 37, -535: 37, -536: 37, -537: 37, -538: 37, -539: 37, -540: 37, -541: 37, -542: 37, -543: 37, -544: 37, -545: 37, -546: 37, -547: 37, -548: 37, -549: 37, -550: 37, -551: 37, -552: 37, -553: 37, -554: 37, -555: 37, -556: 37, -557: 37, -558: 37, -559: 37, -560: 37, -561: 37, -562: 37, -563: 37, -564: 37, -565: 37, -566: 37, -567: 37, -568: 37, -569: 37, -570: 37, -571: 37, -572: 37, -573: 37, -574: 37, -575: 37, -576: 37, -577: 37, -578: 37, -579: 37, -580: 37, -581: 37, -582: 37, -583: 37, -584: 37, -585: 37, -586: 37, -587: 37, -588: 37, -589: 37, -590: 37, -591: 37, -592: 37, -593: 37, -594: 37, -595: 37, -596: 37, -597: 37, -598: 37, -599: 37, -600: 37, -601: 37, -602: 37, -603: 37, -604: 37, -605: 37, -606: 37, -607: 37, -608: 37, -609: 37, -610: 37, -611: 37, -612: 37, -613: 37, -614: 37, -615: 37, -616: 37, -617: 37, -618: 37, -619: 37, -620: 37, -621: 37, -622: 37, -623: 37, -624: 37, -625: 37, -626: 37, -627: 37, -628: 37, -629: 37, -630: 37, -631: 37, -632: 37, -633: 37, -634: 37, -635: 37, -636: 37, -637: 37, -638: 37, -639: 37, -640: 37, -641: 37, -642: 37, -643: 37, -644: 37, -645: 37, -646: 37, -647: 37, -648: 37, -649: 37, -650: 37, -651: 37, -652: 37, -653: 37, -654: 37, -655: 37, -656: 37, -657: 37, -658: 37, -659: 37, -660: 37, -661: 37, -662: 37, -663: 37, -664: 37, -665: 37, -666: 37, -667: 37, -668: 37, -669: 37, -670: 37, -671: 37, -672: 37, -673: 37, -674: 37, -675: 37, -676: 37, -677: 37, -678: 37, -679: 37, -680: 37, -681: 37, -682: 37, -683: 37, -684: 37, -685: 37, -686: 37, -687: 37, -688: 37, -689: 37, -690: 37, -691: 37, -692: 37, -693: 37, -694: 37, -695: 37, -696: 37, -697: 37, -698: 37, -699: 37, -700: 37, -701: 37, -702: 37, -703: 37, -704: 37, -705: 37, -706: 37, -707: 37, -708: 37, -709: 37, -710: 37, -711: 37, -712: 37, -713: 37, -714: 37, -715: 37, -716: 37, -717: 37, -718: 37, -719: 37, -720: 37, -721: 37, -722: 37, -723: 37, -724: 37, -725: 37, -726: 37, -727: 37, -728: 37, -729: 37, -730: 37, -731: 37, -732: 37, -733: 37, -734: 37, -735: 37, -736: 37, -737: 37, -738: 37, -739: 37, -740: 37, -741: 37, -742: 37, -743: 37, -744: 37, -745: 37, -746: 37, -747: 37, -748: 37, -749: 37, -750: 37, -751: 37, -752: 37, -753: 37, -754: 37, -755: 37, -756: 37, -757: 37, -758: 37, -759: 37, -760: 37, -761: 37, -762: 37, -763: 37, -764: 37, -765: 37, -766: 37, -767: 37, -768: 37, -769: 37, -770: 37, -771: 37, -772: 37, -773: 37, -774: 37, -775: 37, -776: 37, -777: 37, -778: 37, -779: 37, -780: 37, -781: 37, -782: 37, -783: 37, -784: 37, -785: 37, -786: 37, -787: 37, -788: 37, -789: 37, -790: 37, -791: 37, -792: 37, -793: 37, -794: 37, -795: 37, -796: 37, -797: 37, -798: 37, -799: 37, -800: 37, -801: 37, -802: 37, -803: 37, -804: 37, -805: 37, -806: 37, -807: 37, -808: 37, -809: 37, -810: 37, -811: 37, -812: 37, -813: 37, -814: 37, -815: 37, -816: 37, -817: 37, -818: 37, -819: 37, -820: 37, -821: 37, -822: 37, -823: 37, -824: 37, -825: 37, -826: 37, -827: 37, -828: 37, -829: 37, -830: 37, -831: 37, -832: 37, -833: 37, -834: 37, -835: 37, -836: 37, -837: 37, -838: 37, -839: 37, -840: 37, -841: 37, -842: 37, -843: 37, -844: 37, -845: 37, -846: 37, -847: 37, -848: 37, -849: 37, -850: 37, -851: 37, -852: 37, -853: 37, -854: 37, -855: 37, -856: 37, -857: 37, -858: 37, -859: 37, -860: 37, -861: 37, -862: 37, -863: 37, -864: 37, -865: 37, -866: 37, -867: 37, -868: 37, -869: 37, -870: 37, -871: 37, -872: 37, -873: 37, -874: 37, -875: 37, -876: 37, -877: 37, -878: 37, -879: 37, -880: 37, -881: 37, -882: 37, -883: 37, -884: 37, -885: 37, -886: 37, -887: 37, -888: 37, -889: 37, -890: 37, -891: 37, -892: 37, -893: 37, -894: 37, -895: 37, -
```

```

In [46]: # 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)

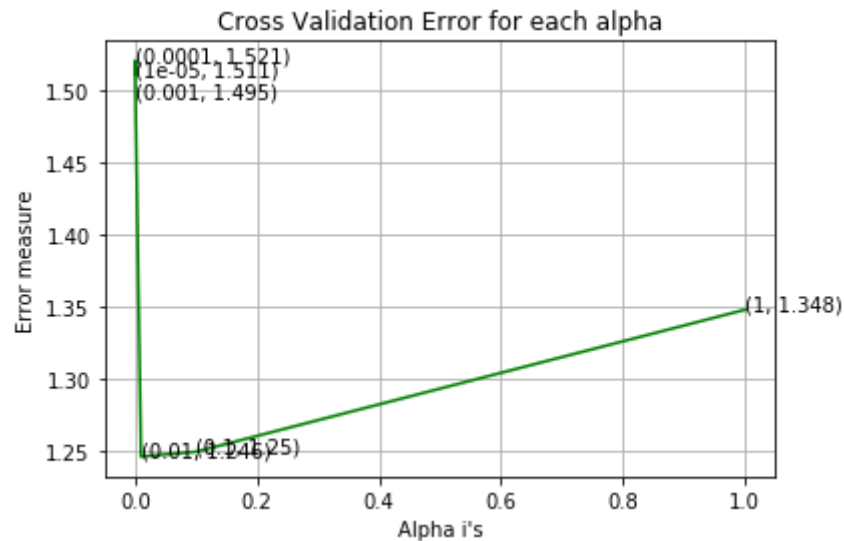
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "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: 1.5109472994988364
For values of alpha = 0.0001 The log loss is: 1.5206385256032233
For values of alpha = 0.001 The log loss is: 1.4953117435158458
For values of alpha = 0.01 The log loss is: 1.2463635372719784
For values of alpha = 0.1 The log loss is: 1.2497437969512837
For values of alpha = 1 The log loss is: 1.3481095036291508

```



For values of best alpha = 0.01 The train log loss is: 0.8809704073864684
For values of best alpha = 0.01 The cross validation log loss is: 1.2463635372719784
For values of best alpha = 0.01 The test log loss is: 1.2144859178402967

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

Ans. Yes, it seems like!

```
In [47]: def get_intersec_text(df):
          df_text_vec = CountVectorizer(ngram_range = (1, 2), 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 [48]: 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")
```

```
92.34 % of word of test data appeared in train data
96.561 % of word of Cross Validation appeared in train data
```

4. Machine Learning Models

In [49]: *#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 [50]: **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 [51]: # 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_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 [52]: *# 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']))
```

In [53]:

```
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, 765503)
(number of data points * number of features) in test data = (665, 765503)
(number of data points * number of features) in cross validation data = (532, 765503)
```

Base Line Model

Logistic Regression

With Class balancing

```
In [54]: #Data preparation for ML models.
#Misc. functionns 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 willl 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)
```

4.3.1.1. Hyper paramter tuning

In [55]:

```
# 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
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

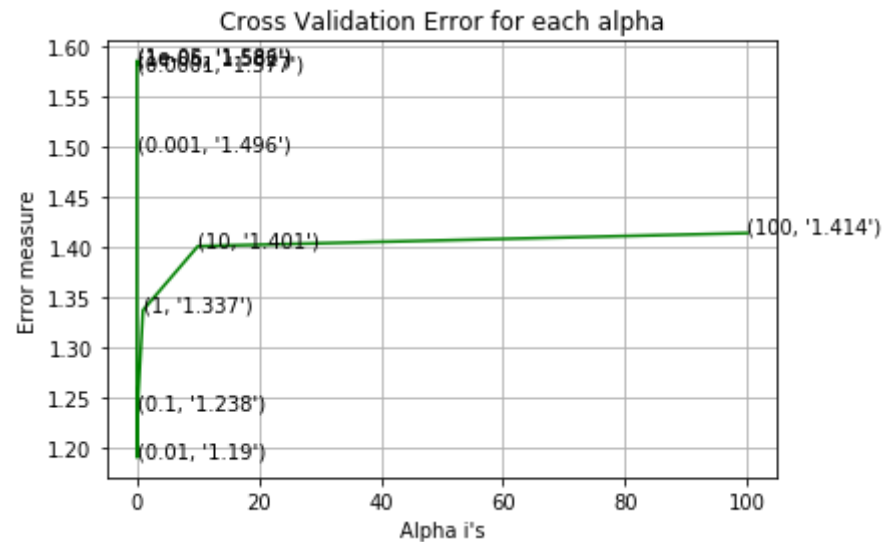
```

```

for alpha = 1e-06
Log Loss : 1.5822198989453882
for alpha = 1e-05
Log Loss : 1.5857302679673981
for alpha = 0.0001
Log Loss : 1.5766750638591878
for alpha = 0.001
Log Loss : 1.4963469808765173
for alpha = 0.01
Log Loss : 1.1899526019434379
for alpha = 0.1
Log Loss : 1.2382407312339097
for alpha = 1

```

Log Loss : 1.3368538802238692
for alpha = 10
Log Loss : 1.4007317580679004
for alpha = 100
Log Loss : 1.4137606889530232



For values of best alpha = 0.01 The train log loss is: 0.8611028579256755
For values of best alpha = 0.01 The cross validation log loss is: 1.1899526019434379
For values of best alpha = 0.01 The test log loss is: 1.1747573830372053

4.3.1.2. Testing the model with best hyper paramters

```
In [56]: # 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)

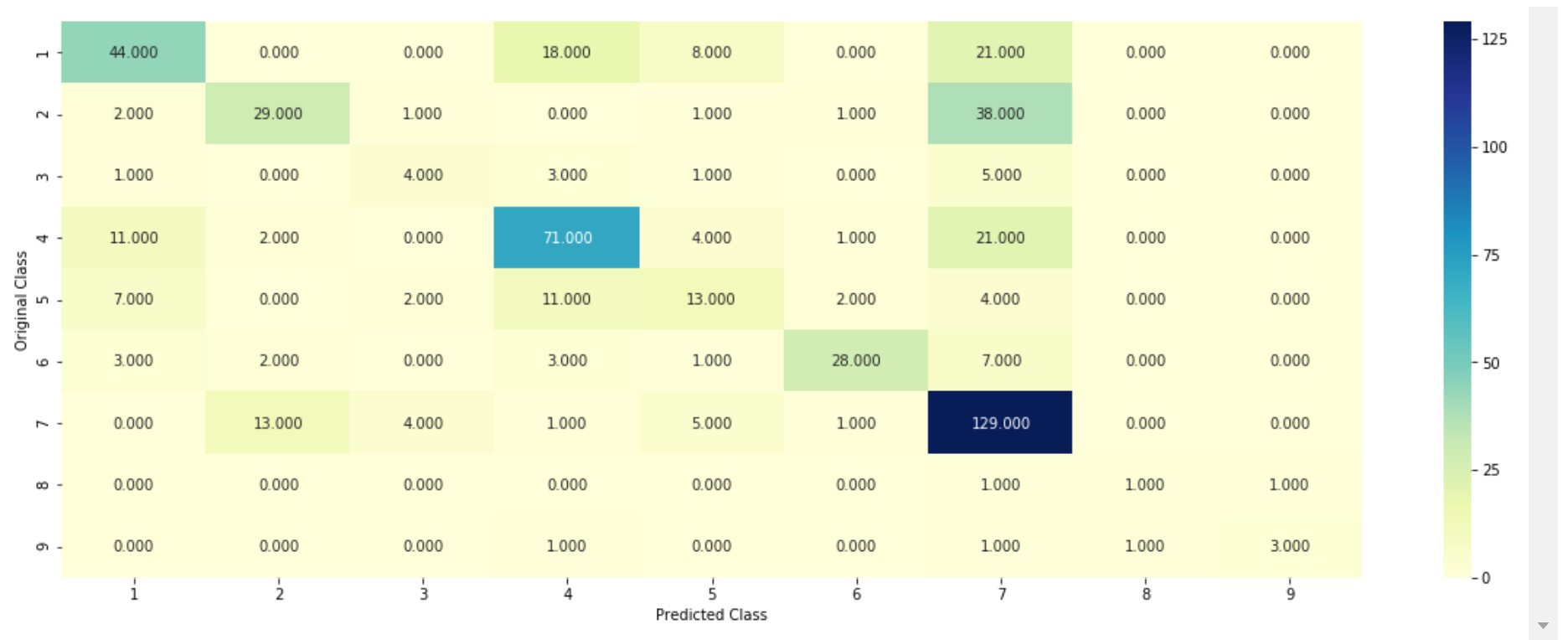
# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X)    Predict class labels for samples in X.

#-----
# 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)
```

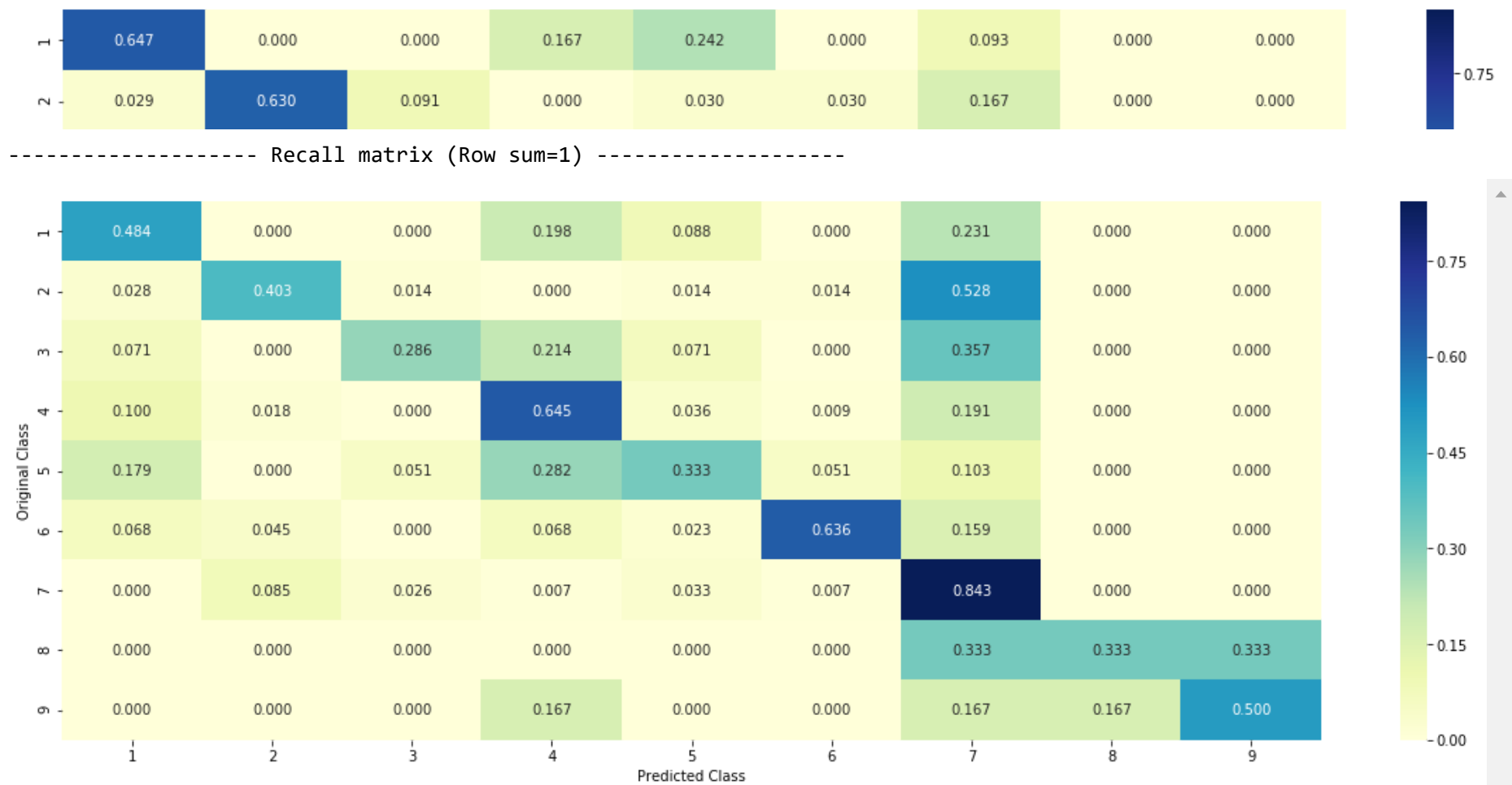
Log loss : 1.1899526019434379

Number of mis-classified points : 0.39473684210526316

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



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



4.3.1.3. Feature Importance

```

In [57]: # 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 = CountVectorizer()
    var_count_vec = CountVectorizer()
    text_count_vec = CountVectorizer(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 [{}]" .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 [{}]" .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 [{}]" .format(word, yes_no))

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

```

4.3.1.3.1. Correctly Classified point

```
In [58]: # 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]
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.107  0.0478 0.0118 0.7586 0.0163 0.0055 0.0307 0.0146 0.0077]]
Actual Class : 1
```

4.3.1.3.2. Incorrectly Classified point

```
In [59]: test_point_index = 5
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]
```

```
Predicted Class : 7
Predicted Class Probabilities: [[7.700e-03 3.950e-02 3.200e-03 1.200e-03 2.100e-03 5.000e-04 9.365e-01
 6.300e-03 3.000e-03]]
Actual Class : 7
```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

```

In [60]: # 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)

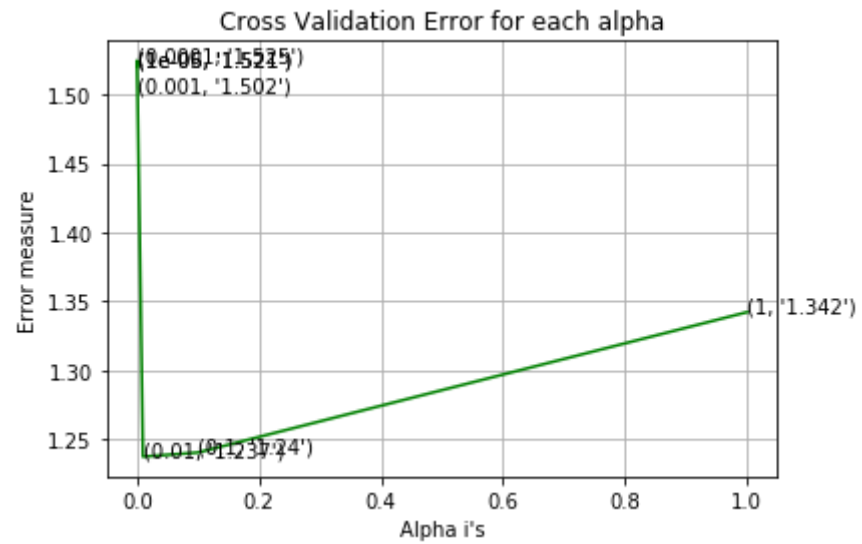
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation 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 : 1.5212345008713353
for alpha = 1e-05
Log Loss : 1.52132422031717
for alpha = 0.0001
Log Loss : 1.52473330258948
for alpha = 0.001
Log Loss : 1.5022523517239366
for alpha = 0.01
Log Loss : 1.2371110900467106
for alpha = 0.1
Log Loss : 1.2400878063650833
for alpha = 1
Log Loss : 1.3420480885251493

```



For values of best alpha = 0.01 The train log loss is: 0.8565374057562842

For values of best alpha = 0.01 The cross validation log loss is: 1.2371110900467106

For values of best alpha = 0.01 The test log loss is: 1.2026094299197474

4.3.2.2. Testing model with best hyper parameters

```
In [61]: # 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.

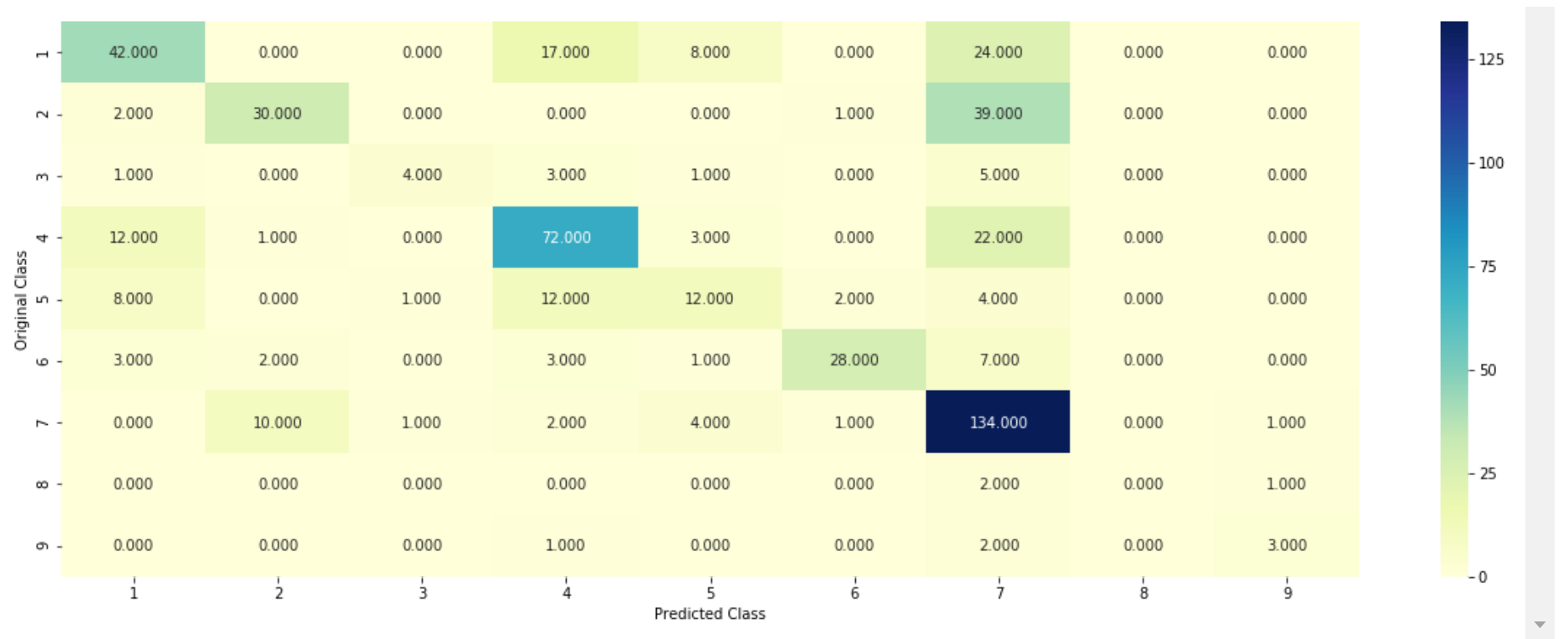
#-----
# video link:
#-----

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)
```

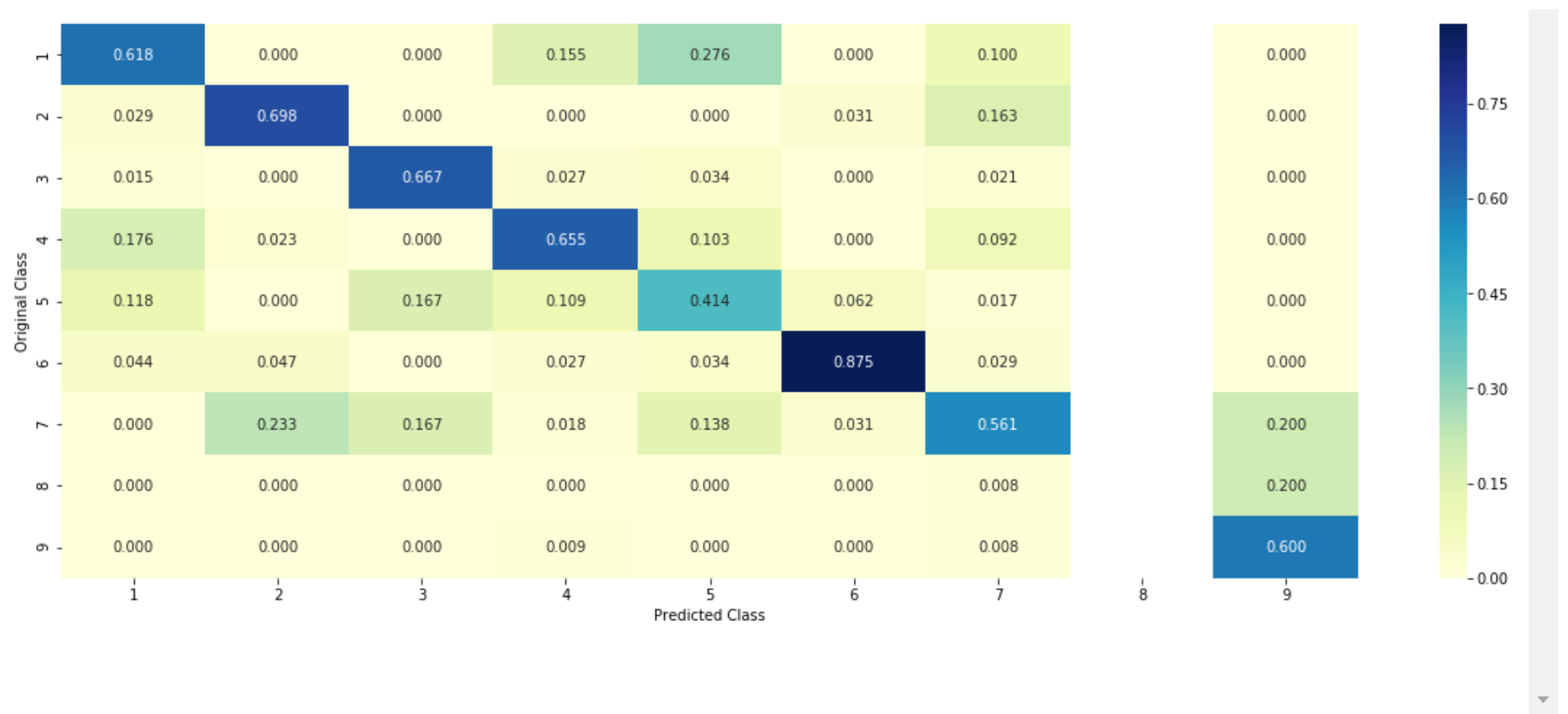
Log loss : 1.2371110900467106

Number of mis-classified points : 0.3890977443609023

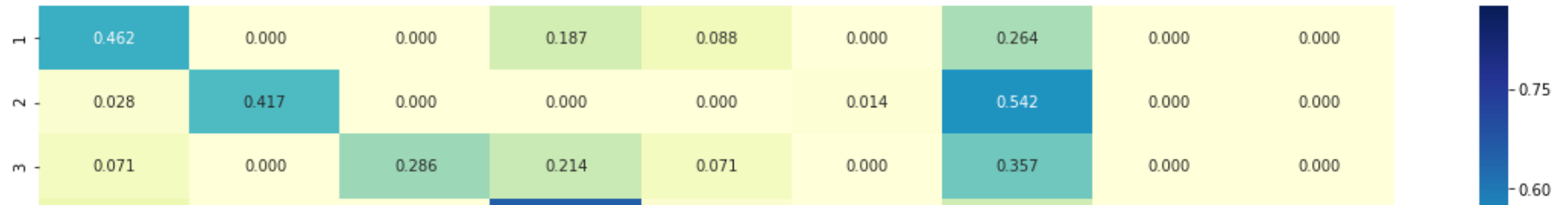
----- Confusion matrix -----



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



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



4.3.2.3. Feature Importance, Correctly Classified point

```
In [62]: 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]
```

Predicted Class : 4

Predicted Class Probabilities: [[0.086 0.0426 0.0013 0.8103 0.0059 0.0019 0.0433 0.0087 0.]]

Actual Class : 1

4.3.2.4. Feature Importance, Incorrectly Classified point

```
In [63]: test_point_index = 5
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]
```

```
Predicted Class : 7
Predicted Class Probabilities: [[1.320e-02 4.970e-02 3.000e-04 2.700e-03 1.000e-03 3.000e-04 9.262e-01
 6.700e-03 0.000e+00]]
Actual Class : 7
```

Comparing the scores from all the above models

Logistic regression with CountVectorizer Features, including both unigrams and bigrams

```
In [65]: # 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(["LR with class Balancing", "0.861", "1.189", "1.174", "39.47"])
x.add_row(["LR without class Balancing", "0.856", "1.237", "1.202", "38.90"])
print(x)
```

Model-Name	Train loss	CV loss	Test Loss	% Misclassified
LR with class Balancing	0.861	1.189	1.174	39.47
LR without class Balancing	0.856	1.237	1.202	38.90

Conclusion / Observations :-

-
1. For Task 3 we applied countvectorizer with unigrams & bi-grams on Logistic Regression.
 2. On Applying Logistic-Regression with Class-Balancing we got a log-loss of 1.174 for test data with 39.47% of mis-classified points.
 3. On Applying Logistic-Regression without Class-Balancing we got a log-loss of 1.202 for test data with 38.90% of mis- classified points.

In []: