

# Personalized Cancer diagnosis using Genetic Data (BOW Unigrams + Bigrams)

## 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=qxXRVompl8> (<https://www.youtube.com/watch?v=qxXRVompl8>)

### 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 [3]:

```
from IPython.core.display import display, HTML
display(HTML("<style>.container { width:100% !important; }</style>"))

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 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 [4]:

```
data = pd.read_csv('training/training_variants')
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[4]:

	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 [5]:

```
# note '|' the separator in this file
data_text =pd.read_csv("training/training_text",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[5]:

	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...

### 3.1.3. Preprocessing of text

In [6]:

```
# 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 [7]:

```
#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 : 211.58687500000002 seconds

In [8]:

```
#Merging both gene_variations and text data based on ID
result = pd.merge(data, data_text,on='ID', how='left')
result.head()
```

Out[8]:

	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 [9]:

```
#Check for null entries
result[result.isnull().any(axis=1)]
```

Out[9]:

	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 [10]:

```
#Replace the null entries in Text by Gene + Variation Data
result.loc[result['TEXT'].isnull(),'TEXT'] = result['Gene'] + ' '+result['Variation']
```

In [11]:

```
#Display row after null replacement
result[result['ID']==1109]
```

Out[11]:

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

In [12]:

```
result.columns
```

Out[12]:

Index(['ID', 'Gene', 'Variation', 'Class', 'TEXT'], dtype='object')

In [13]:

```
#Save the final dataframe to "result.csv"
result.to_csv("result.csv",columns=result.columns,index=None)
```

In [14]:

```
#Load the pre-processed data frame
result=pd.read_csv('result.csv')
result.head()
```

Out[14]:

	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...

### 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 [15]:

```
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' [stratify=y_train]
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 [16]:

```
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 [17]:

```
# 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 = 'rbkymc'
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()

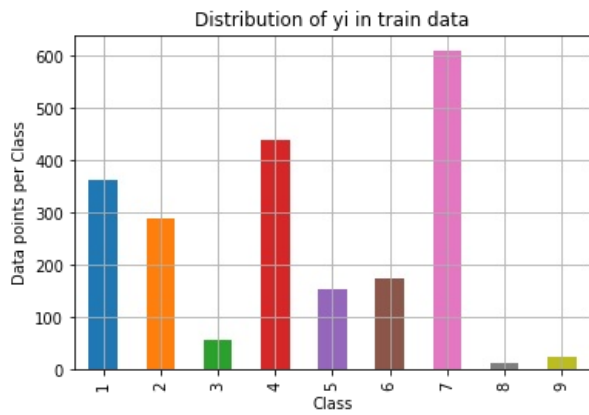
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':', train_class_distribution.values[i], '(', np.round((train_class_distribution.values[i]/train_df.shape[0]*100), 3), '%)')

print('-'*80)
my_colors = 'rbkymc'
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_distribution.values[i]/test_df.shape[0]*100), 3), '%)')

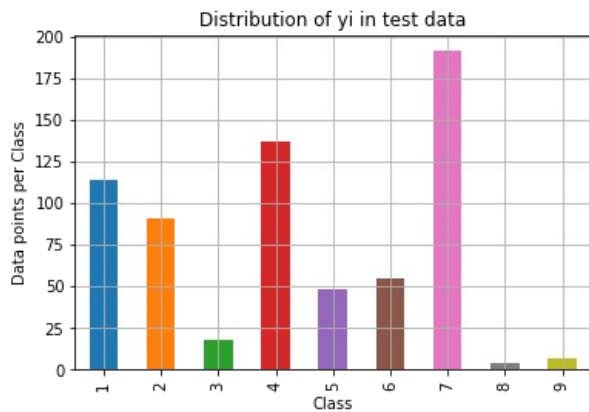
print('-'*80)
my_colors = 'rbkymc'
cv_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in cross validation data')
plt.grid()
plt.show()

# 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_distribution.values[i]/cv_df.shape[0]*100), 3), '%)')
```



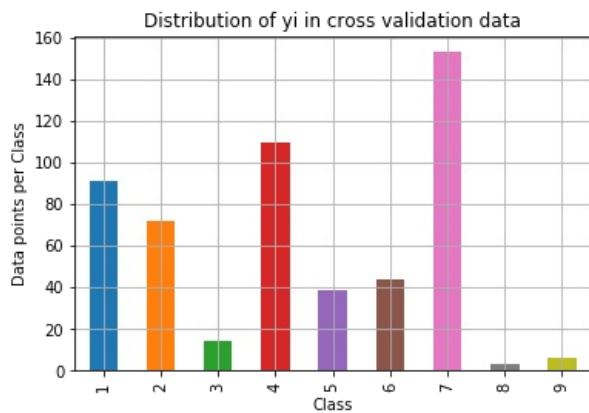
Number of data points in class 7 : 609 ( 28.672 %)  
 Number of data points in class 4 : 439 ( 20.669 %)  
 Number of data points in class 1 : 363 ( 17.09 %)  
 Number of data points in class 2 : 289 ( 13.606 %)  
 Number of data points in class 6 : 176 ( 8.286 %)  
 Number of data points in class 5 : 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 %)



## Key Take-away:

We can see the class labels are distributed fairly equally in all the three dataset - train, cross validation and test data. If it hadn't been the case our models won't perform good. The distributions has to be similar in all the three datasets for all the 9 classes.

## 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 [18]:

```
# 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(axis = 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(axis = 0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                       [3/4, 4/6]]

    #C = Actual confusion matrix
    #B = Recall matrix
    #C = Precision Matrix
    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 (Column 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 [19]:

```
# 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))
print()

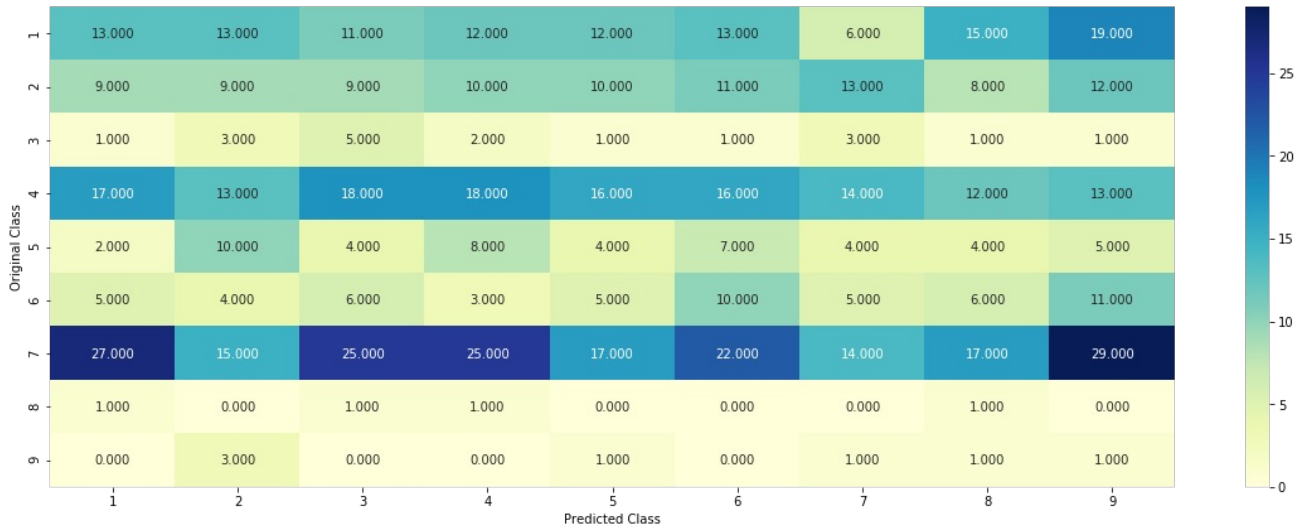
# Test-Set error.
#we create a output array that has exactly same as the test data
test_predicted_y = np.zeros((test_data_len,9))
for i in range(test_data_len):
    rand_probs = np.random.rand(1,9)
    test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs))))[0])
print("Log loss on Test Data using Random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

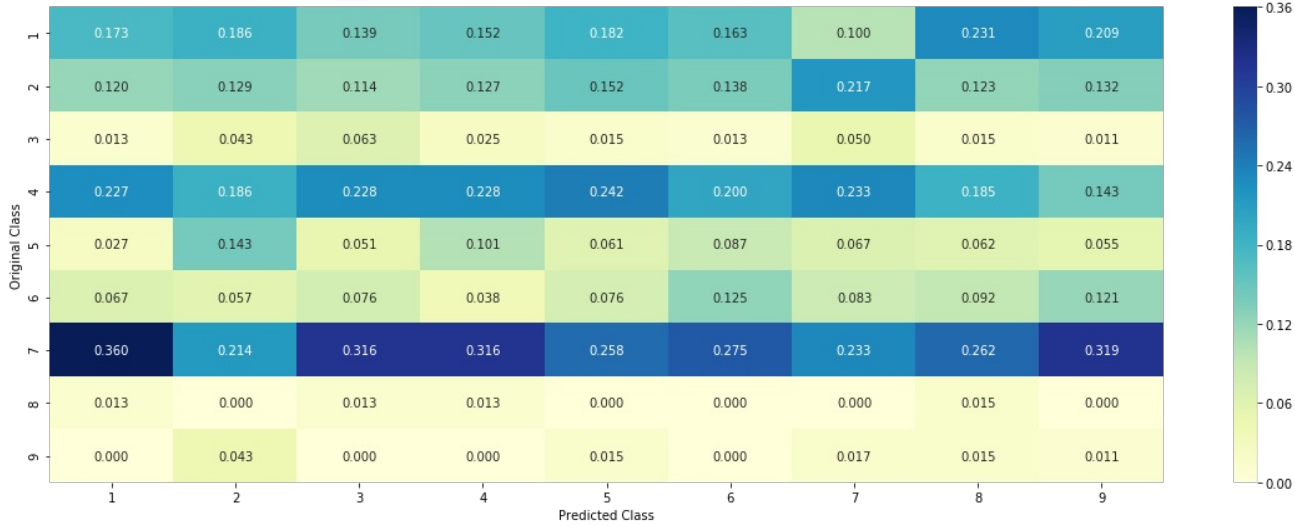
Log loss on Cross Validation Data using Random Model 2.5732273744205525

Log loss on Test Data using Random Model 2.4326444436611143

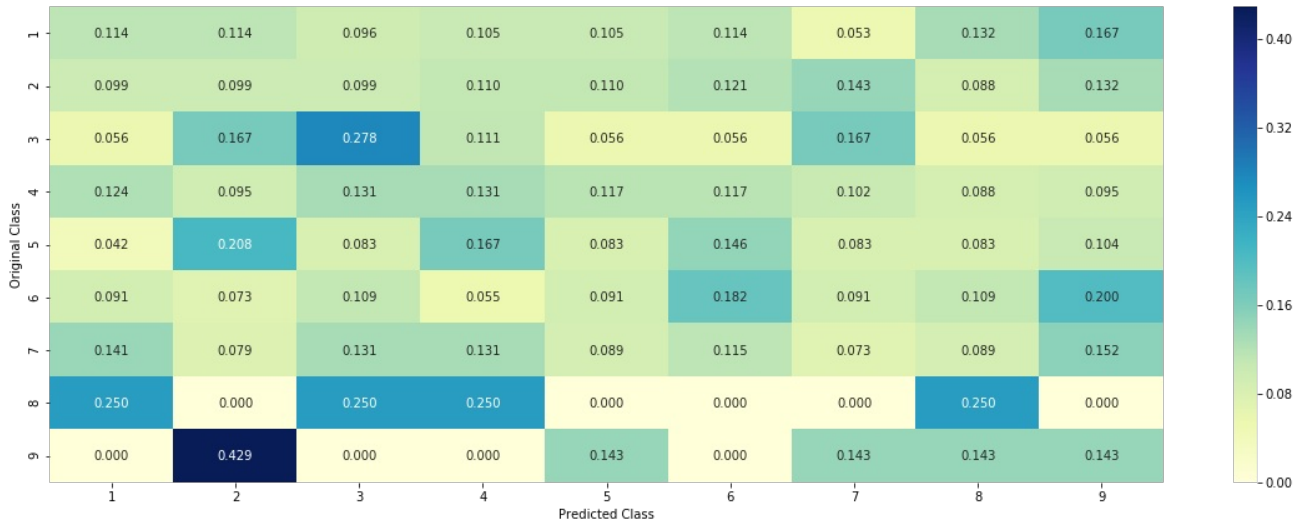
----- Confusion matrix -----



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



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



### 3.3 Univariate Analysis

In [20]:

```
# 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 a 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 time it
occurred in total data+90*alpha)
# gv_dict is like a look up table, for every gene it store 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, 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
    # i = feature name, denominator = feature count
    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 the
            dataset of class 1(Let)
            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, 0.06818181818181817, 0.13636363636363635, 0.25, 0.
    .19318181818181818, 0.03787878787878788, 0.03787878787878788, 0.03787878787878788],
    #   'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.27040816326530615, 0.06122
    4489795918366, 0.066326530612244902, 0.051020408163265307, 0.051020408163265307, 0.056122448979591837],
    #   'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.06818181818181817, 0.0681818181818177
    , 0.0625, 0.34659090909090912, 0.0625, 0.056818181818181816],

```

```
# 'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.0787878787878782, 0.139
3939393939394, 0.34545454545454546, 0.060606060606060608, 0.060606060606060608, 0.060606060606060608],
# 'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.46540880503144655, 0.0754
71698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081761006289, 0.062893081761006289],
# 'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295, 0.06622
5165562913912, 0.066225165562913912, 0.27152317880794702, 0.066225165562913912, 0.066225165562913912],
# 'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333333333333334, 0.073333333333333334, 0.0933
3333333333338, 0.080000000000000002, 0.29999999999999999, 0.066666666666666666, 0.066666666666666666],
# ...
# }
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 wil
l add the feature to gv_fea
# if not we will add [1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9] 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([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
# gv_fea.append([-1,-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 [21]:

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

```
Number of Unique Genes in train data: 233
BRCA1      169
TP53       103
EGFR        96
BRCA2       84
PTEN        81
KIT         67
BRAF        53
ERBB2       50
ALK         44
PDGFRA      37
Name: Gene, dtype: int64
```

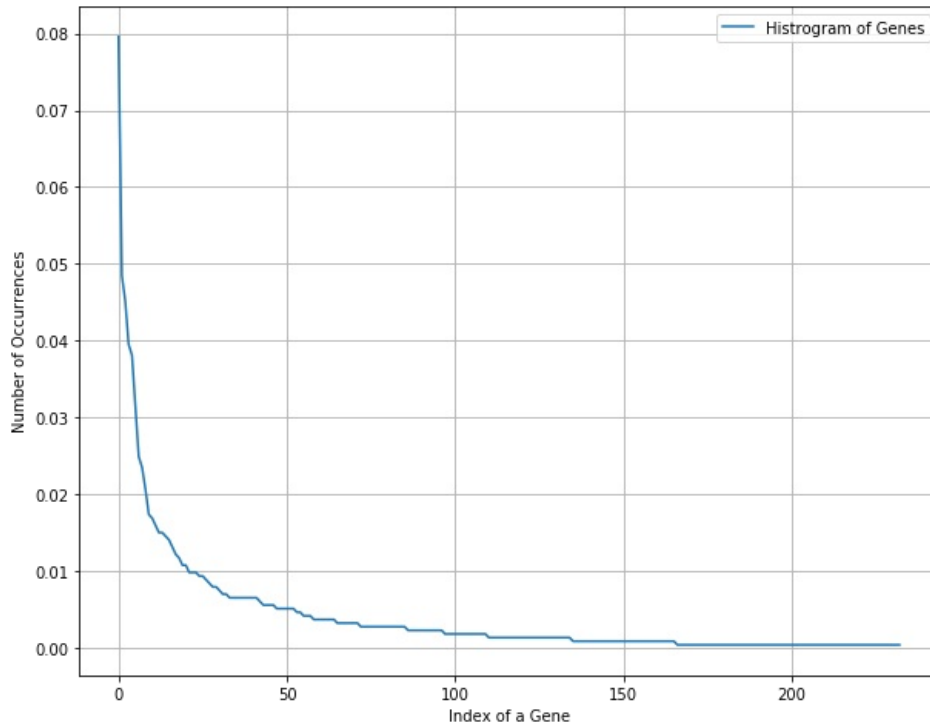
In [22]:

```
print("Ans: There are", unique_genes.shape[0], "different categories of genes in the training data, and they are
distributed as follows",)
```

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

In [23]:

```
#Histogram of distribution of genes
s = sum(unique_genes.values);
h = unique_genes.values/s;
plt.figure(figsize=(10,8))
plt.plot(h, label="Histogram of Genes")
plt.xlabel('Index of a Gene')
plt.ylabel('Number of Occurrences')
plt.legend()
plt.grid()
plt.show()
```

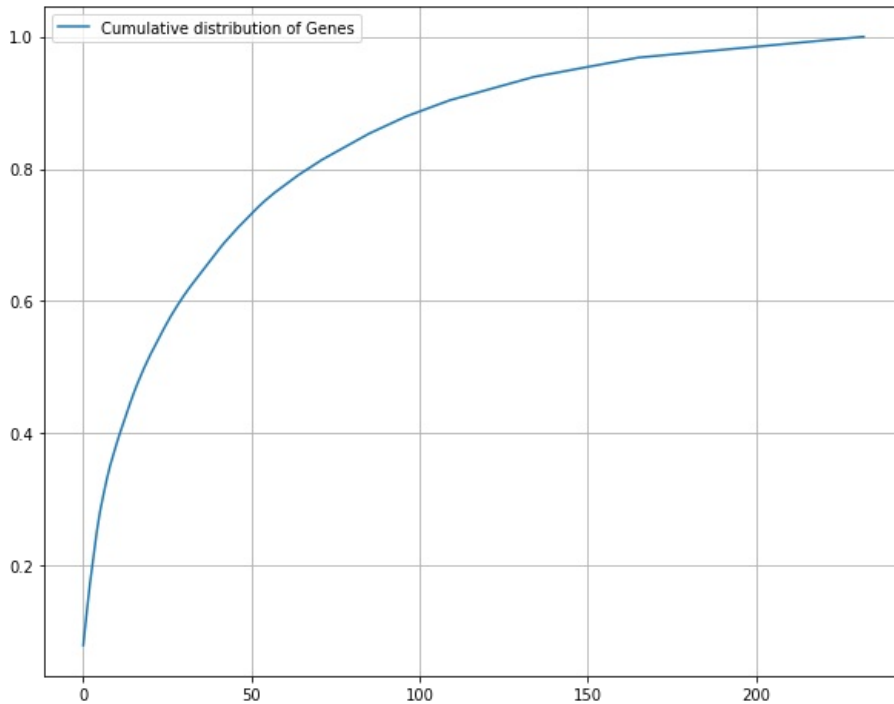


**Key Take-away from the above Histogram:**

The dsitributions of the gene is very skewed towards the left. There are almost 50 genes which occurs most frequently, and there are almost 200 odd genes which occurs less frequently than the first 50.

In [24]:

```
#Get the CDF of gene
c = np.cumsum(h)
plt.figure(figsize=(10,8))
plt.plot(c,label='Cumulative distribution of Genes')
plt.grid()
plt.legend()
plt.show()
```



#### Key Take-away from the above CDF Plot:

The top 50 genes contribute to almost 75% of our data. The rest of the 185 genes contribute to only 25% of the data. This says that the top 50% gene mutation type occurs most frequently. There are also many gene/mutation pair types which occurs very less frequently.

### 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 [25]:

```
#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 [26]:

```
print("train_gene_feature_responseCoding is converted to a feature vector using response coding method. The shape of gene feature:", train_gene_feature_responseCoding.shape)
print("test_gene_feature_responseCoding is converted to a feature vector using response coding method. The shape of gene feature:", test_gene_feature_responseCoding.shape)
print("cv_gene_feature_responseCoding is converted to a feature vector using response coding method. The shape of gene feature:", cv_gene_feature_responseCoding.shape)
```

train\_gene\_feature\_responseCoding is converted to a feature vector using response coding method. The shape of gene feature: (2124, 9)  
test\_gene\_feature\_responseCoding is converted to a feature vector using response coding method. The shape of gene feature: (665, 9)  
cv\_gene\_feature\_responseCoding is converted to a feature vector using response coding method. The shape of gene feature: (532, 9)

In [27]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer(ngram_range=(1,2))
train_gene_feature_ngrams = gene_vectorizer.fit_transform(train_df['Gene'])
test_gene_feature_ngrams = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_ngrams = gene_vectorizer.transform(cv_df['Gene'])

print("train_gene_feature_ngrams is converted to a feature vector using one-hot coding method. The shape of gene feature:", train_gene_feature_ngrams.shape)
print("test_gene_feature_ngrams is converted to a feature vector using one-hot coding method. The shape of gene feature:", test_gene_feature_ngrams.shape)
print("cv_gene_feature_ngrams is converted to a feature vector using one-hot coding method. The shape of gene feature:", cv_gene_feature_ngrams.shape)
```

train\_gene\_feature\_ngrams is converted to a feature vector using one-hot coding method. The shape of gene feature: (2124, 233)  
test\_gene\_feature\_ngrams is converted to a feature vector using one-hot coding method. The shape of gene feature: (665, 233)  
cv\_gene\_feature\_ngrams is converted to a feature vector using one-hot coding method. The shape of gene feature: (532, 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 [28]:

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

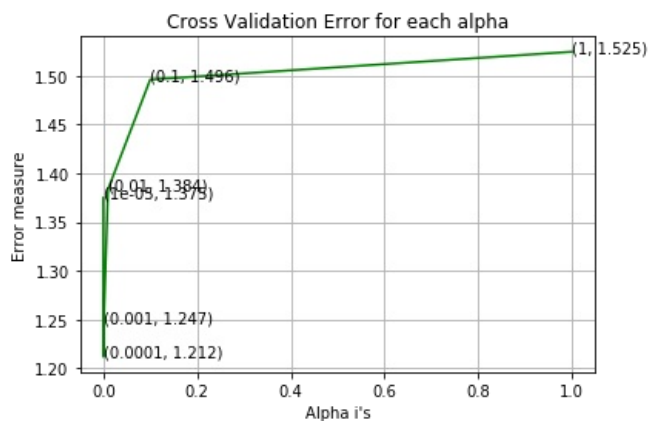
cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_gene_feature_ngrams, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_gene_feature_ngrams, y_train)
    predict_y = sig_clf.predict_proba(cv_gene_feature_ngrams)
    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_ngrams, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_gene_feature_ngrams, y_train)

predict_y = sig_clf.predict_proba(train_gene_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_gene_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_gene_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

For values of alpha = 1e-05 The log loss is: 1.3749198199513637  
For values of alpha = 0.0001 The log loss is: 1.211911334844622  
For values of alpha = 0.001 The log loss is: 1.246651170042778  
For values of alpha = 0.01 The log loss is: 1.3835669245352178  
For values of alpha = 0.1 The log loss is: 1.4961402200334912  
For values of alpha = 1 The log loss is: 1.5247159760397542



For values of best alpha = 0.0001 The train log loss is: 1.0251757570548385  
For values of best alpha = 0.0001 The cross validation log loss is: 1.211911334844622  
For values of best alpha = 0.0001 The test log loss is: 1.2006717886375393

**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 [29]:

```
print("Q6. How many data points in Test and CV datasets are covered by the ", unique_genes.shape[0], " genes in t
rain dataset?\n")

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 640 out of 665 : 96.2406015037594
2. In cross validation data 518 out of 532 : 97.36842105263158

### 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 [30]:

```
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 : 1929
Truncating_Mutations      61
Deletion                  47
Amplification              45
Fusions                   21
Overexpression             3
G12V                      3
E17K                      3
R841K                     2
EWSR1-ETV1_Fusion         2
S308A                     2
Name: Variation, dtype: int64
```

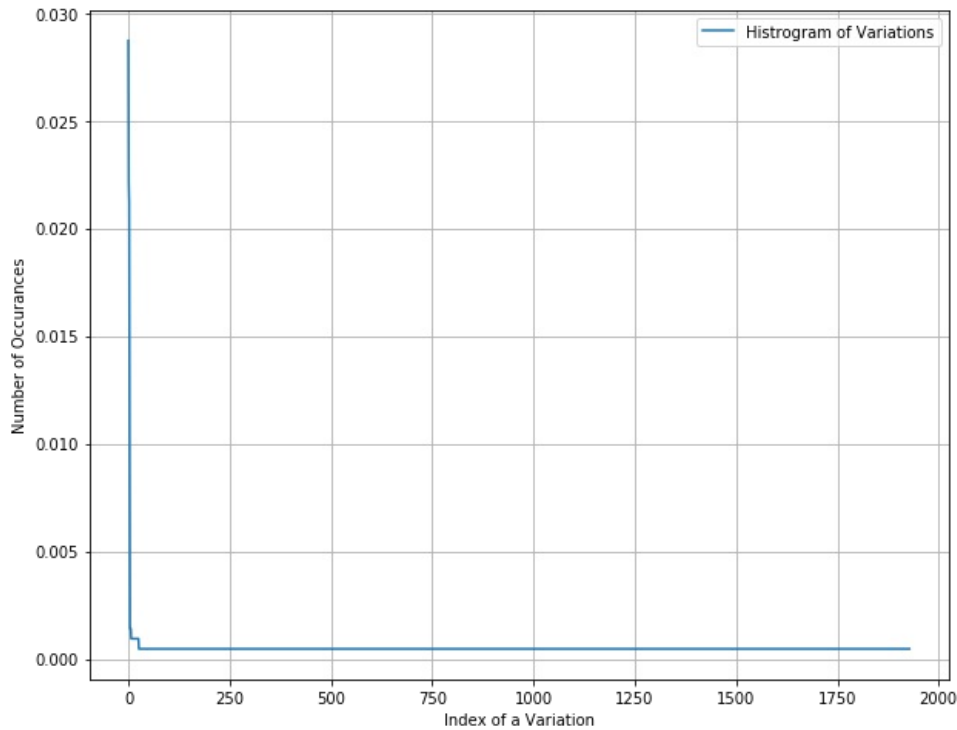
In [31]:

```
print("Ans: There are", unique_variations.shape[0] , "different categories of variations in the train data, and th
ey are distributed as follows",)
```

Ans: There are 1929 different categories of variations in the train data, and they are distributed as follows

In [32]:

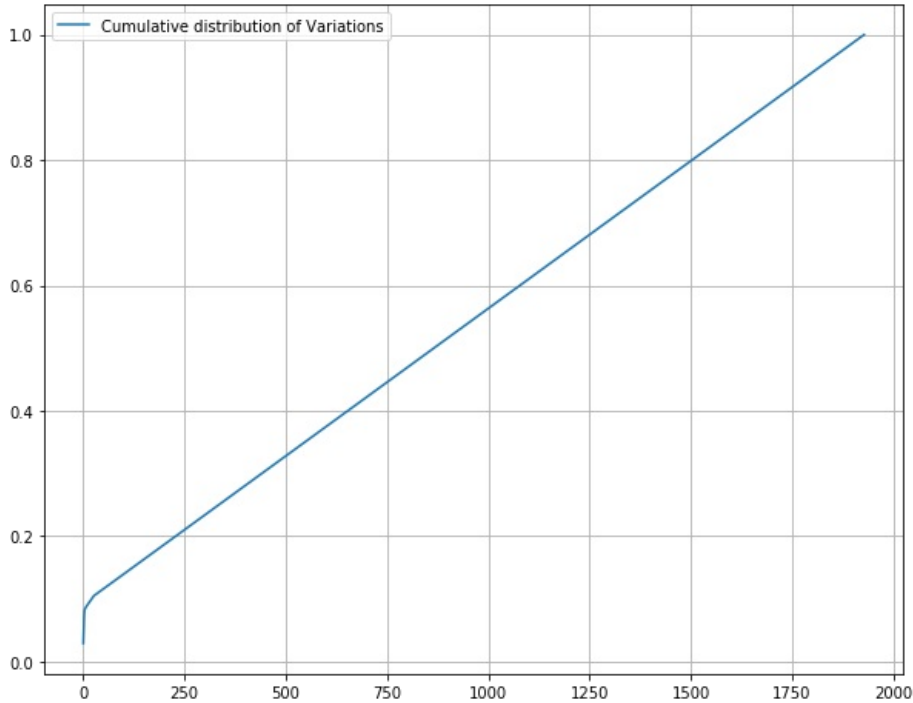
```
#Histogram for the distribution of the fetaure Variation
plt.figure(figsize=(10,8))
s = sum(unique_variations.values);
h = unique_variations.values/s;
plt.plot(h, label="Histogram of Variations")
plt.xlabel('Index of a Variation')
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```



In [33]:

```
#CDF Of variation feature
plt.figure(figsize=(10,8))
c = np.cumsum(h)
print(c)
plt.plot(c,label='Cumulative distribution of Variations')
plt.grid()
plt.legend()
plt.show()
```

```
[0.0287194  0.05084746 0.0720339 ... 0.99905838 0.99952919 1. ... ]
```



### Key Take-aways:

This CDF suggest that almost 80% of the dataset is explained by 1500 variations (out of a total number of 1927). If you see a CDF which is almost a straight line, you can quickly tell that most variations occur once or maybe twice in the entire training data.

### 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

### Very Important note about response coding: Avoiding response leakage.

1. We have to be extremely careful not to use the test and cross validate data for response coding. This is because we don't the issue of data leakage.
2. Suppose we have a variant V2 present in Test / Cross Val dataset. But V2 is not present in Train. So in that case, while building the response bales for V2, we will just assign equal probability values to each of the 9 array values. Proba =  $1/9$  for each og them.
3. We will take the help of laplace smoothing in order to achieve this. Without laplace smoothing, we would get a 0/0 error.
4. We are seeing the data that is present in the cross validation and test data during the time of training. So we are literally leaking the information that is present in test/cv data at the time of training.

This should be strictly avoided, as we do not want a data leakage issue.

In [34]:

```
# alpha is used for laplace smoothing (Naive Bayes)
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))

print("train_variation_feature_responseCoding is a converted feature using the response coding method. The shape of Variation feature:", train_variation_feature_responseCoding.shape)
print("test_variation_feature_responseCoding is a converted feature using the response coding method. The shape of Variation feature:", test_variation_feature_responseCoding.shape)
print("cv_variation_feature_responseCoding is a converted feature using the response coding method. The shape of Variation feature:", cv_variation_feature_responseCoding.shape)
```

train\_variation\_feature\_responseCoding is a converted feature using the response coding method. The shape of Variation feature: (2124, 9)  
test\_variation\_feature\_responseCoding is a converted feature using the response coding method. The shape of Variation feature: (665, 9)  
cv\_variation\_feature\_responseCoding is a converted feature using the response coding method. The shape of Variation feature: (532, 9)

In [35]:

```
# one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer(ngram_range=(1,2))
train_variation_feature_ngrams = variation_vectorizer.fit_transform(train_df['Variation'])
test_variation_feature_ngrams = variation_vectorizer.transform(test_df['Variation'])
cv_variation_feature_ngrams = variation_vectorizer.transform(cv_df['Variation'])

print("train_variation_feature_ngrams is converted feature using the one-hot encoding method. The shape of Variation feature:", train_variation_feature_ngrams.shape)
print("test_variation_feature_ngrams is converted feature using the one-hot encoding method. The shape of Variation feature:", test_variation_feature_ngrams.shape)
print("cv_variation_feature_ngrams is converted feature using the one-hot encoding method. The shape of Variation feature:", cv_variation_feature_ngrams.shape)
```

train\_variation\_feature\_ngrams is converted feature using the one-hot encoding method. The shape of Variation feature: (2124, 2056)  
test\_variation\_feature\_ngrams is converted feature using the one-hot encoding method. The shape of Variation feature: (665, 2056)  
cv\_variation\_feature\_ngrams is converted feature using the one-hot encoding method. The shape of Variation feature: (532, 2056)

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

Let's build a model just like the earlier!

In [36]:

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

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

    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_variation_feature_ngrams, y_train)
    predict_y = sig_clf.predict_proba(cv_variation_feature_ngrams)

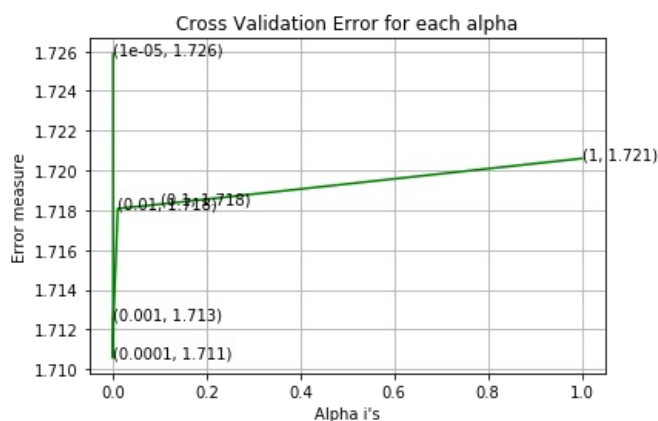
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

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

predict_y = sig_clf.predict_proba(train_variation_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_variation_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_variation_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

For values of alpha = 1e-05 The log loss is: 1.7258780480035605  
For values of alpha = 0.0001 The log loss is: 1.7105678916951321  
For values of alpha = 0.001 The log loss is: 1.7125529923634253  
For values of alpha = 0.01 The log loss is: 1.7180833852657083  
For values of alpha = 0.1 The log loss is: 1.7183001775825715  
For values of alpha = 1 The log loss is: 1.7206023151639096



For values of best alpha = 0.0001 The train log loss is: 0.7494343679674323  
For values of best alpha = 0.0001 The cross validation log loss is: 1.7105678916951321  
For values of best alpha = 0.0001 The test log loss is: 1.700329708727258

## Conclusion:

1. The variation feature is certainly useful in determining the class labels, as it has brought down the log loss of a random model by a significant level.
2. However, if we compare the test and cross validation log loss of variation features to that of the gene features, we see that test/cross validation log loss for variations is quite high as compared to genes. This suggests that the variation feature might be unstable.
3. Also, there is a significant difference in log loss for train and test/cross-validation. This means the model is overfitting using only the variation feature.
4. Having said all these, we will still keep the variation feature as it has managed to decrease the log-loss of a random model. We will use this feature along with other features to see how it behaves when we build our final model.

**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.

**Q12 .** How many data points are covered by total 1924 variations in test and cross validation data sets?

In [37]:

```
print("Q12. How many data points are covered by total ", unique_variations.shape[0], " variations in test and cross validation data sets?\n")
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 1929 variations in test and cross validation data sets?

Ans.

1. In test data 67 out of 665 : 10.075187969924812
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 [38]:

```
# cls_text is a data frame
# for every row in data frame consider the 'TEXT' column
# 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 [39]:

```
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 [40]:

```
# building a BOW text vectorizer with all the words
text_vectorizer = CountVectorizer(ngram_range=(1,2))
train_text_feature_ngrams = text_vectorizer.fit_transform(train_df['TEXT'])
# getting all the feature names (words)
train_text_features= text_vectorizer.get_feature_names()

#Copy the train_text_feature_ngrams into a separate variable
train_text_feature_ngrams_ones = train_text_feature_ngrams.copy()
train_text_feature_ngrams_ones.data = train_text_feature_ngrams_ones.data / train_text_feature_ngrams_ones.data

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

#Convert all the float values to integers
train_text_fea_counts = [int(col_sum[i]) for i in range(0,len(col_sum))]

# 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 unigrams + bigrams in train data :", len(train_text_features))
```

Total number of unique unigrams + bigrams in train data : 2362629

In [41]:

```
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 [42]:

```
#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 [43]:

```
# 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 [44]:

```
# don't forget to normalize every feature
train_text_feature_ngrams = normalize(train_text_feature_ngrams, axis=0)

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

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



In [45]:

```
#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 [46]:

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

```
Counter({1: 1255340, 2: 349150, 3: 180882, 4: 103541, 5: 67010, 6: 60578, 7: 48188, 8: 45542, 9: 284
24, 10: 25775, 11: 21673, 12: 18618, 13: 12116, 14: 11543, 17: 9361, 15: 8960, 16: 7869, 18: 6631, 1
9: 6026, 28: 5426, 20: 4208, 25: 4197, 22: 4149, 21: 4129, 23: 3764, 26: 3540, 24: 3432, 42: 3025, 2
9: 2401, 27: 2359, 52: 2130, 44: 2126, 30: 2029, 32: 1829, 31: 1785, 33: 1558, 53: 1357, 34: 1357, 6
6: 1307, 35: 1290, 36: 1189, 37: 1166, 45: 1083, 43: 1054, 38: 1045, 39: 984, 40: 911, 41: 879, 47:
826, 46: 819, 54: 748, 48: 703, 49: 625, 55: 603, 50: 596, 51: 582, 56: 573, 57: 548, 58: 535, 67: 5
07, 59: 494, 60: 466, 61: 455, 63: 434, 68: 408, 69: 382, 62: 374, 70: 370, 64: 366, 71: 334, 65: 32
9, 72: 328, 74: 323, 73: 316, 75: 287, 76: 265, 78: 258, 77: 256, 79: 242, 80: 230, 83: 227, 81: 222
, 82: 214, 84: 213, 85: 208, 86: 200, 91: 191, 87: 185, 90: 177, 94: 175, 88: 173, 95: 166, 92: 157,
89: 155, 93: 147, 97: 146, 96: 140, 101: 138, 99: 138, 106: 136, 104: 136, 102: 136, 103: 130, 114:
124, 107: 124, 98: 123, 108: 120, 105: 120, 117: 112, 110: 109, 100: 109, 111: 108, 115: 105, 113: 1
05, 109: 101, 119: 99, 123: 98, 118: 97, 116: 96, 120: 95, 122: 92, 131: 90, 112: 88, 127: 84, 130:
83, 139: 78, 133: 77, 126: 77, 128: 75, 125: 75, 121: 75, 137: 73, 138: 72, 140: 71, 147: 70, 129: 7
0, 135: 69, 136: 68, 132: 66, 144: 65, 134: 65, 124: 64, 142: 62, 154: 60, 151: 59, 145: 58, 156: 56
, 153: 56, 168: 55, 148: 55, 143: 55, 164: 54, 150: 54, 162: 53, 157: 52, 141: 51, 159: 50, 152: 50,
173: 49, 155: 47, 166: 45, 146: 45, 160: 44, 163: 43, 149: 43, 184: 42, 161: 42, 170: 41, 169: 41, 1
65: 41, 158: 41, 182: 40, 190: 39, 185: 38, 172: 38, 178: 37, 174: 37, 171: 36, 229: 35, 180: 35, 17
9: 35, 205: 34, 197: 34, 195: 34, 192: 34, 167: 34, 187: 33, 245: 32, 207: 31, 189: 31, 183: 31, 211
: 30, 191: 30, 186: 30, 181: 30, 177: 30, 193: 29, 236: 28, 231: 28, 215: 28, 204: 28, 217: 27, 209:
27, 188: 27, 176: 27, 213: 26, 203: 26, 234: 25, 225: 25, 206: 25, 199: 25, 198: 25, 194: 25, 175: 2
5, 257: 24, 251: 24, 232: 24, 230: 24, 212: 24, 289: 23, 247: 23, 227: 23, 226: 23, 208: 23, 201: 23
, 200: 23, 228: 22, 224: 22, 222: 22, 214: 22, 196: 22, 276: 21, 267: 21, 221: 21, 220: 21, 219: 21,
210: 21, 287: 20, 271: 20, 265: 20, 250: 20, 242: 20, 223: 20, 317: 19, 274: 19, 237: 19, 235: 19, 2
33: 19, 350: 18, 272: 18, 255: 18, 249: 18, 240: 18, 218: 18, 202: 18, 304: 17, 303: 17, 300: 17, 27
9: 17, 277: 17, 260: 17, 243: 17, 216: 17, 364: 16, 357: 16, 340: 16, 290: 16, 273: 16, 256: 16, 248
: 16, 246: 16, 238: 16, 379: 15, 346: 15, 344: 15, 321: 15, 284: 15, 282: 15, 270: 15, 253: 15, 241:
15, 384: 14, 374: 14, 348: 14, 334: 14, 324: 14, 313: 14, 291: 14, 288: 14, 285: 14, 283: 14, 278: 1
4, 266: 14, 264: 14, 259: 14, 362: 13, 354: 13, 330: 13, 328: 13, 302: 13, 299: 13, 275: 13, 269: 13
, 261: 13, 473: 12, 403: 12, 396: 12, 368: 12, 343: 12, 333: 12, 312: 12, 308: 12, 252: 12, 244: 12,
239: 12, 670: 11, 432: 11, 411: 11, 402: 11, 394: 11, 377: 11, 369: 11, 335: 11, 327: 11, 325: 11, 3
22: 11, 320: 11, 315: 11, 307: 11, 301: 11, 294: 11, 280: 11, 262: 11, 258: 11, 254: 11, 480: 10, 44
7: 10, 430: 10, 429: 10, 375: 10, 367: 10, 359: 10, 358: 10, 351: 10, 339: 10, 338: 10, 329: 10, 318
: 10, 314: 10, 311: 10, 306: 10, 298: 10, 295: 10, 281: 10, 263: 10, 465: 9, 463: 9, 444: 9, 442: 9,
427: 9, 426: 9, 418: 9, 390: 9, 387: 9, 386: 9, 370: 9, 365: 9, 363: 9, 360: 9, 355: 9, 353: 9, 349:
9, 332: 9, 310: 9, 292: 9, 760: 8, 693: 8, 677: 8, 652: 8, 600: 8, 580: 8, 554: 8, 506: 8, 500: 8, 4
97: 8, 468: 8, 456: 8, 446: 8, 439: 8, 431: 8, 421: 8, 420: 8, 415: 8, 413: 8, 393: 8, 376: 8, 366:
8, 352: 8, 345: 8, 337: 8, 305: 8, 296: 8, 286: 8, 658: 7, 655: 7, 577: 7, 569: 7, 548: 7, 545: 7, 5
35: 7, 499: 7, 490: 7, 487: 7, 484: 7, 476: 7, 471: 7, 470: 7, 466: 7, 462: 7, 440: 7, 438: 7, 437:
7, 436: 7, 419: 7, 408: 7, 398: 7, 389: 7, 388: 7, 383: 7, 382: 7, 381: 7, 361: 7, 347: 7, 331: 7, 2
93: 7, 268: 7, 798: 6, 687: 6, 642: 6, 628: 6, 611: 6, 591: 6, 578: 6, 551: 6, 536: 6, 532: 6, 528:
6, 527: 6, 514: 6, 510: 6, 508: 6, 504: 6, 502: 6, 496: 6, 494: 6, 488: 6, 474: 6, 469: 6, 461: 6, 4
58: 6, 457: 6, 452: 6, 433: 6, 425: 6, 423: 6, 416: 6, 412: 6, 409: 6, 401: 6, 399: 6, 395: 6, 385:
6, 380: 6, 378: 6, 373: 6, 342: 6, 341: 6, 326: 6, 323: 6, 316: 6, 1641: 5, 942: 5, 917: 5, 806: 5,
799: 5, 771: 5, 711: 5, 704: 5, 702: 5, 663: 5, 654: 5, 650: 5, 648: 5, 644: 5, 641: 5, 586: 5, 585:
5, 570: 5, 565: 5, 558: 5, 549: 5, 544: 5, 543: 5, 541: 5, 534: 5, 530: 5, 529: 5, 523: 5, 491: 5, 4
82: 5, 475: 5, 464: 5, 459: 5, 450: 5, 449: 5, 428: 5, 424: 5, 422: 5, 410: 5, 406: 5, 405: 5, 392:
5, 372: 5, 356: 5, 319: 5, 309: 5, 297: 5, 1256: 4, 1086: 4, 1067: 4, 1033: 4, 967: 4, 959: 4, 931:
4, 923: 4, 918: 4, 905: 4, 900: 4, 890: 4, 887: 4, 869: 4, 820: 4, 782: 4, 776: 4, 762: 4, 758: 4, 7
42: 4, 735: 4, 734: 4, 699: 4, 697: 4, 695: 4, 694: 4, 689: 4, 686: 4, 685: 4, 684: 4, 675: 4, 673:
4, 668: 4, 660: 4, 656: 4, 635: 4, 631: 4, 627: 4, 626: 4, 624: 4, 619: 4, 610: 4, 609: 4, 608: 4, 6
07: 4, 598: 4, 592: 4, 590: 4, 588: 4, 587: 4, 574: 4, 573: 4, 566: 4, 562: 4, 561: 4, 559: 4, 557:
4, 556: 4, 539: 4, 537: 4, 526: 4, 525: 4, 522: 4, 521: 4, 519: 4, 515: 4, 512: 4, 493: 4, 483: 4, 4
81: 4, 479: 4, 478: 4, 467: 4, 455: 4, 453: 4, 448: 4, 441: 4, 407: 4, 404: 4, 400: 4, 397: 4, 391:
4, 371: 4, 336: 4, 1973: 3, 1942: 3, 1894: 3, 1815: 3, 1671: 3, 1588: 3, 1524: 3, 1521: 3, 1513: 3,
1473: 3, 1450: 3, 1443: 3, 1382: 3, 1377: 3, 1368: 3, 1346: 3, 1343: 3, 1258: 3, 1240: 3, 1182: 3, 1
173: 3, 1169: 3, 1159: 3, 1150: 3, 1144: 3, 1126: 3, 1109: 3, 1107: 3, 1088: 3, 1082: 3, 1079: 3, 10
70: 3, 1061: 3, 1038: 3, 1029: 3, 1021: 3, 1019: 3, 1016: 3, 1012: 3, 1007: 3, 997: 3, 996: 3, 995:
3, 979: 3, 976: 3, 970: 3, 966: 3, 964: 3, 947: 3, 934: 3, 924: 3, 916: 3, 914: 3, 902: 3, 898: 3, 8
95: 3, 893: 3, 872: 3, 864: 3, 862: 3, 852: 3, 846: 3, 837: 3, 828: 3, 824: 3, 823: 3, 816: 3, 813:
3, 800: 3, 794: 3, 774: 3, 773: 3, 759: 3, 757: 3, 754: 3, 752: 3, 750: 3, 749: 3, 748: 3, 741: 3, 7
40: 3, 737: 3, 732: 3, 731: 3, 730: 3, 727: 3, 725: 3, 721: 3, 717: 3, 716: 3, 714: 3, 706: 3, 701:
3, 683: 3, 680: 3, 671: 3, 667: 3, 665: 3, 661: 3, 659: 3, 657: 3, 653: 3, 649: 3, 646: 3, 645: 3, 6
43: 3, 640: 3, 625: 3, 621: 3, 620: 3, 617: 3, 616: 3, 613: 3, 603: 3, 601: 3, 599: 3, 584: 3, 582:
3, 579: 3, 575: 3, 568: 3, 567: 3, 563: 3, 555: 3, 552: 3, 547: 3, 540: 3, 533: 3, 531: 3, 524: 3, 5
20: 3, 517: 3, 516: 3, 513: 3, 511: 3, 505: 3, 501: 3, 498: 3, 495: 3, 489: 3, 485: 3, 477: 3, 472:
3, 443: 3, 434: 3, 417: 3, 2069: 2, 2052: 2, 2041: 2, 1993: 2, 1982: 2, 1898: 2, 1891: 2, 1868: 2, 1
820: 2, 1783: 2, 1776: 2, 1763: 2, 1752: 2, 1737: 2, 1720: 2, 1715: 2, 1701: 2, 1680: 2, 1640: 2, 16
04: 2, 1597: 2, 1577: 2, 1575: 2, 1574: 2, 1570: 2, 1555: 2, 1553: 2, 1545: 2, 1470: 2, 1457: 2, 144
7: 2, 1410: 2, 1407: 2, 1406: 2, 1393: 2, 1388: 2, 1383: 2, 1366: 2, 1354: 2, 1348: 2, 1336: 2, 1333
: 2, 1330: 2, 1325: 2, 1318: 2, 1317: 2, 1314: 2, 1300: 2, 1299: 2, 1291: 2, 1290: 2, 1284: 2, 1277:
2, 1275: 2, 1260: 2, 1255: 2, 1251: 2, 1249: 2, 1243: 2, 1231: 2, 1230: 2, 1222: 2, 1217: 2, 1189: 2
, 1180: 2, 1177: 2, 1175: 2, 1171: 2, 1165: 2, 1161: 2, 1156: 2, 1153: 2, 1143: 2, 1133: 2, 1122: 2,
```

1120: 2, 1119: 2, 1111: 2, 1108: 2, 1069: 2, 1062: 2, 1059: 2, 1053: 2, 1050: 2, 1048: 2, 1047: 2, 1046: 2, 1041: 2, 1039: 2, 1035: 2, 1028: 2, 1025: 2, 1022: 2, 1014: 2, 1013: 2, 1004: 2, 1001: 2, 1000: 2, 999: 2, 998: 2, 990: 2, 988: 2, 986: 2, 982: 2, 978: 2, 977: 2, 974: 2, 969: 2, 965: 2, 961: 2, 955: 2, 952: 2, 950: 2, 948: 2, 940: 2, 939: 2, 937: 2, 936: 2, 935: 2, 929: 2, 922: 2, 920: 2, 915: 2, 913: 2, 910: 2, 901: 2, 899: 2, 894: 2, 888: 2, 885: 2, 884: 2, 880: 2, 879: 2, 871: 2, 870: 2, 866: 2, 865: 2, 854: 2, 853: 2, 851: 2, 850: 2, 848: 2, 847: 2, 839: 2, 836: 2, 835: 2, 832: 2, 831: 2, 827: 2, 818: 2, 817: 2, 812: 2, 811: 2, 809: 2, 805: 2, 804: 2, 801: 2, 796: 2, 791: 2, 789: 2, 786: 2, 780: 2, 778: 2, 775: 2, 772: 2, 769: 2, 768: 2, 767: 2, 766: 2, 764: 2, 763: 2, 761: 2, 756: 2, 755: 2, 753: 2, 747: 2, 746: 2, 745: 2, 743: 2, 733: 2, 729: 2, 724: 2, 720: 2, 713: 2, 712: 2, 710: 2, 709: 2, 708: 2, 707: 2, 705: 2, 703: 2, 700: 2, 692: 2, 679: 2, 674: 2, 672: 2, 669: 2, 664: 2, 662: 2, 651: 2, 647: 2, 638: 2, 637: 2, 636: 2, 632: 2, 623: 2, 614: 2, 606: 2, 605: 2, 602: 2, 595: 2, 589: 2, 583: 2, 581: 2, 576: 2, 572: 2, 571: 2, 564: 2, 553: 2, 546: 2, 538: 2, 518: 2, 509: 2, 507: 2, 503: 2, 486: 2, 454: 2, 451: 2, 445: 2, 435: 2, 414: 2, 2098: 1, 2097: 1, 2080: 1, 2074: 1, 2073: 1, 2071: 1, 2068: 1, 2058: 1, 2056: 1, 2048: 1, 2034: 1, 2031: 1, 2020: 1, 2019: 1, 2017: 1, 2011: 1, 2010: 1, 2004: 1, 1997: 1, 1990: 1, 1984: 1, 1975: 1, 1974: 1, 1967: 1, 1966: 1, 1963: 1, 1958: 1, 1953: 1, 1949: 1, 1948: 1, 1940: 1, 1931: 1, 1922: 1, 1917: 1, 1916: 1, 1915: 1, 1914: 1, 1912: 1, 1911: 1, 1901: 1, 1899: 1, 1892: 1, 1889: 1, 1884: 1, 1878: 1, 1874: 1, 1871: 1, 1870: 1, 1863: 1, 1861: 1, 1857: 1, 1856: 1, 1854: 1, 1853: 1, 1847: 1, 1845: 1, 1844: 1, 1841: 1, 1840: 1, 1834: 1, 1833: 1, 1832: 1, 1829: 1, 1821: 1, 1817: 1, 1814: 1, 1808: 1, 1799: 1, 1796: 1, 1794: 1, 1792: 1, 1789: 1, 1786: 1, 1778: 1, 1774: 1, 1771: 1, 1770: 1, 1769: 1, 1768: 1, 1766: 1, 1758: 1, 1756: 1, 1751: 1, 1748: 1, 1742: 1, 1739: 1, 1738: 1, 1728: 1, 1724: 1, 1721: 1, 1719: 1, 1708: 1, 1706: 1, 1698: 1, 1697: 1, 1694: 1, 1692: 1, 1685: 1, 1684: 1, 1681: 1, 1678: 1, 1676: 1, 1673: 1, 1670: 1, 1669: 1, 1667: 1, 1661: 1, 1660: 1, 1655: 1, 1653: 1, 1647: 1, 1645: 1, 1644: 1, 1631: 1, 1630: 1, 1629: 1, 1627: 1, 1625: 1, 1622: 1, 1613: 1, 1610: 1, 1607: 1, 1606: 1, 1605: 1, 1600: 1, 1596: 1, 1595: 1, 1590: 1, 1589: 1, 1585: 1, 1584: 1, 1576: 1, 1573: 1, 1563: 1, 1561: 1, 1560: 1, 1559: 1, 1557: 1, 1550: 1, 1542: 1, 1537: 1, 1536: 1, 1534: 1, 1527: 1, 1525: 1, 1523: 1, 1519: 1, 1518: 1, 1517: 1, 1516: 1, 1510: 1, 1509: 1, 1506: 1, 1505: 1, 1504: 1, 1502: 1, 1500: 1, 1499: 1, 1494: 1, 1491: 1, 1489: 1, 1488: 1, 1487: 1, 1486: 1, 1484: 1, 1477: 1, 1472: 1, 1469: 1, 1465: 1, 1461: 1, 1460: 1, 1459: 1, 1453: 1, 1452: 1, 1449: 1, 1445: 1, 1440: 1, 1438: 1, 1436: 1, 1432: 1, 1431: 1, 1430: 1, 1427: 1, 1422: 1, 1420: 1, 1418: 1, 1417: 1, 1412: 1, 1405: 1, 1403: 1, 1400: 1, 1399: 1, 1398: 1, 1397: 1, 1396: 1, 1390: 1, 1389: 1, 1387: 1, 1386: 1, 1384: 1, 1378: 1, 1375: 1, 1374: 1, 1371: 1, 1370: 1, 1369: 1, 1365: 1, 1364: 1, 1360: 1, 1358: 1, 1355: 1, 1353: 1, 1352: 1, 1340: 1, 1339: 1, 1334: 1, 1332: 1, 1323: 1, 1319: 1, 1316: 1, 1315: 1, 1309: 1, 1308: 1, 1305: 1, 1304: 1, 1298: 1, 1297: 1, 1293: 1, 1287: 1, 1281: 1, 1279: 1, 1276: 1, 1274: 1, 1272: 1, 1271: 1, 1269: 1, 1268: 1, 1267: 1, 1263: 1, 1262: 1, 1257: 1, 1254: 1, 1247: 1, 1245: 1, 1242: 1, 1232: 1, 1229: 1, 1227: 1, 1226: 1, 1224: 1, 1218: 1, 1215: 1, 1213: 1, 1211: 1, 1208: 1, 1201: 1, 1200: 1, 1195: 1, 1194: 1, 1190: 1, 1188: 1, 1185: 1, 1183: 1, 1181: 1, 1179: 1, 1176: 1, 1170: 1, 1166: 1, 1163: 1, 1162: 1, 1160: 1, 1157: 1, 1155: 1, 1151: 1, 1148: 1, 1142: 1, 1141: 1, 1140: 1, 1135: 1, 1134: 1, 1128: 1, 1125: 1, 1124: 1, 1123: 1, 1114: 1, 1113: 1, 1110: 1, 1105: 1, 1102: 1, 1098: 1, 1096: 1, 1093: 1, 1091: 1, 1089: 1, 1087: 1, 1078: 1, 1077: 1, 1075: 1, 1074: 1, 1072: 1, 1071: 1, 1068: 1, 1066: 1, 1065: 1, 1064: 1, 1060: 1, 1056: 1, 1055: 1, 1054: 1, 1051: 1, 1049: 1, 1042: 1, 1040: 1, 1037: 1, 1036: 1, 1034: 1, 1032: 1, 1024: 1, 1020: 1, 1017: 1, 1011: 1, 1009: 1, 1006: 1, 1005: 1, 1003: 1, 1002: 1, 994: 1, 993: 1, 992: 1, 991: 1, 989: 1, 987: 1, 983: 1, 981: 1, 980: 1, 973: 1, 972: 1, 968: 1, 963: 1, 960: 1, 957: 1, 956: 1, 954: 1, 953: 1, 951: 1, 946: 1, 945: 1, 944: 1, 938: 1, 933: 1, 932: 1, 930: 1, 927: 1, 925: 1, 921: 1, 912: 1, 909: 1, 908: 1, 907: 1, 904: 1, 903: 1, 896: 1, 891: 1, 886: 1, 883: 1, 882: 1, 881: 1, 877: 1, 874: 1, 873: 1, 868: 1, 867: 1, 858: 1, 857: 1, 856: 1, 855: 1, 849: 1, 844: 1, 843: 1, 842: 1, 841: 1, 840: 1, 834: 1, 833: 1, 830: 1, 829: 1, 826: 1, 819: 1, 814: 1, 810: 1, 808: 1, 803: 1, 802: 1, 797: 1, 792: 1, 790: 1, 787: 1, 781: 1, 779: 1, 777: 1, 770: 1, 765: 1, 751: 1, 738: 1, 728: 1, 726: 1, 723: 1, 722: 1, 719: 1, 715: 1, 698: 1, 696: 1, 690: 1, 681: 1, 678: 1, 676: 1, 666: 1, 639: 1, 634: 1, 633: 1, 630: 1, 629: 1, 615: 1, 604: 1, 597: 1, 596: 1, 594: 1, 593: 1, 560: 1, 550: 1, 542: 1, 492: 1, 460: 1}})

In [47]:

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

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

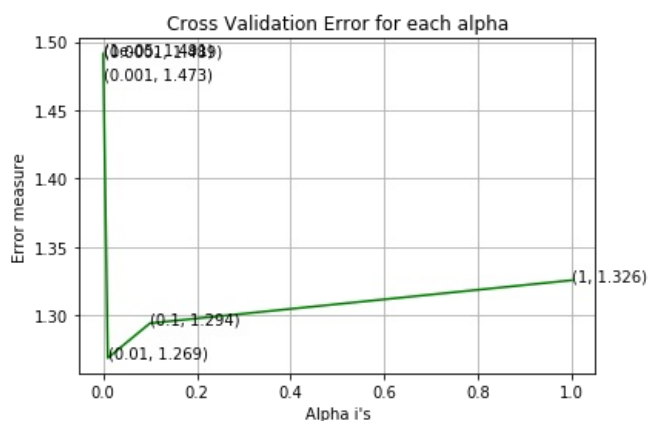
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_text_feature_ngrams, y_train)
    predict_y = sig_clf.predict_proba(cv_text_feature_ngrams)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The 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_ngrams, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_text_feature_ngrams, y_train)

predict_y = sig_clf.predict_proba(train_text_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_text_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_text_feature_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

For values of alpha = 1e-05 The log loss is: 1.4914608845323882  
For values of alpha = 0.0001 The log loss is: 1.4893264737788092  
For values of alpha = 0.001 The log loss is: 1.4726563650207807  
For values of alpha = 0.01 The log loss is: 1.2689359374469482  
For values of alpha = 0.1 The log loss is: 1.294338494805306  
For values of alpha = 1 The log loss is: 1.3257615120116675



For values of best alpha = 0.01 The train log loss is: 0.7598148191815772  
For values of best alpha = 0.01 The cross validation log loss is: 1.2689359374469482  
For values of best alpha = 0.01 The test log loss is: 1.277667004524932

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

**Ans.** Yes, it seems like!

In [48]:

```
def get_intersec_text(df):
    df_text_vec = CountVectorizer(ngram_range=(1,2))
    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 [49]:

```
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")
```

73.432 % of word of test data appeared in train data  
73.459 % of word of Cross Validation appeared in train data

## 4. Machine Learning Models

In [61]:

```
#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("Percentage of mis-classified points :", 100*np.count_nonzero((pred_y - test_y))/test_y.shape[0])
    plot_confusion_matrix(test_y, pred_y)
```

In [51]:

```
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)
```

## Stacking the three types of features

In [53]:

```
# 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_ngrams = hstack((train_gene_feature_ngrams, train_variation_feature_ngrams))
test_gene_var_ngrams = hstack((test_gene_feature_ngrams, test_variation_feature_ngrams))
cv_gene_var_ngrams = hstack((cv_gene_feature_ngrams, cv_variation_feature_ngrams))

train_x_ngrams = hstack((train_gene_var_ngrams, train_text_feature_ngrams)).tocsr()
train_y = np.array(list(train_df['Class']))

test_x_ngrams = hstack((test_gene_var_ngrams, test_text_feature_ngrams)).tocsr()
test_y = np.array(list(test_df['Class']))

cv_x_ngrams = hstack((cv_gene_var_ngrams, cv_text_feature_ngrams)).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 [54]:

```
print("BOW(n grams = 1,2) features representation:\n")
print("(Number of data points * number of features) in train data = ", train_x_ngrams.shape)
print("(Number of data points * number of features) in test data = ", test_x_ngrams.shape)
print("(Number of data points * number of features) in cross validation data = ", cv_x_ngrams.shape)
```

BOW(n\_grams = 1,2) features representation:

```
(Number of data points * number of features) in train data = (2124, 2364918)
(Number of data points * number of features) in test data = (665, 2364918)
(Number of data points * number of features) in cross validation data = (532, 2364918)
```

In [55]:

```
print(" Response encoding features :\n")
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)
```

## 4.3.1 Logistic Regression With Class balancing

### 4.3.1.1. Hyperparameter tuning

In [56]:

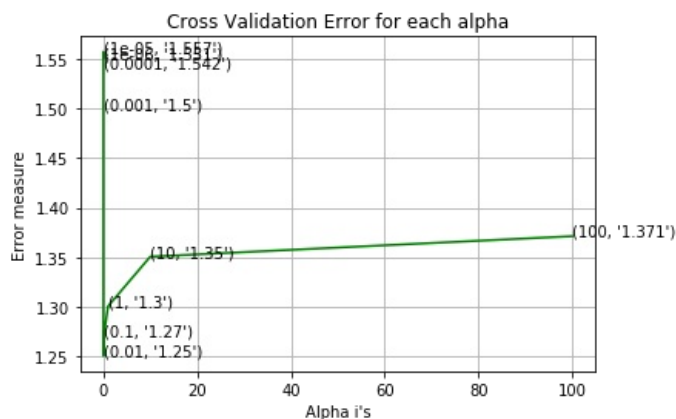
```
alpha = [10 ** x for x in range(-6, 3)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_ngrams, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_ngrams, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_ngrams)
    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_ngrams, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_ngrams, train_y)

predict_y = sig_clf.predict_proba(train_x_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
for alpha = 1e-06
Log Loss : 1.5507962982043606
for alpha = 1e-05
Log Loss : 1.5574692532467123
for alpha = 0.0001
Log Loss : 1.541812402343388
for alpha = 0.001
Log Loss : 1.4995155091305106
for alpha = 0.01
Log Loss : 1.2501309862948178
for alpha = 0.1
Log Loss : 1.2702929185519398
for alpha = 1
Log Loss : 1.3003345704918075
for alpha = 10
Log Loss : 1.3504027199342965
for alpha = 100
Log Loss : 1.3710652999098063
```



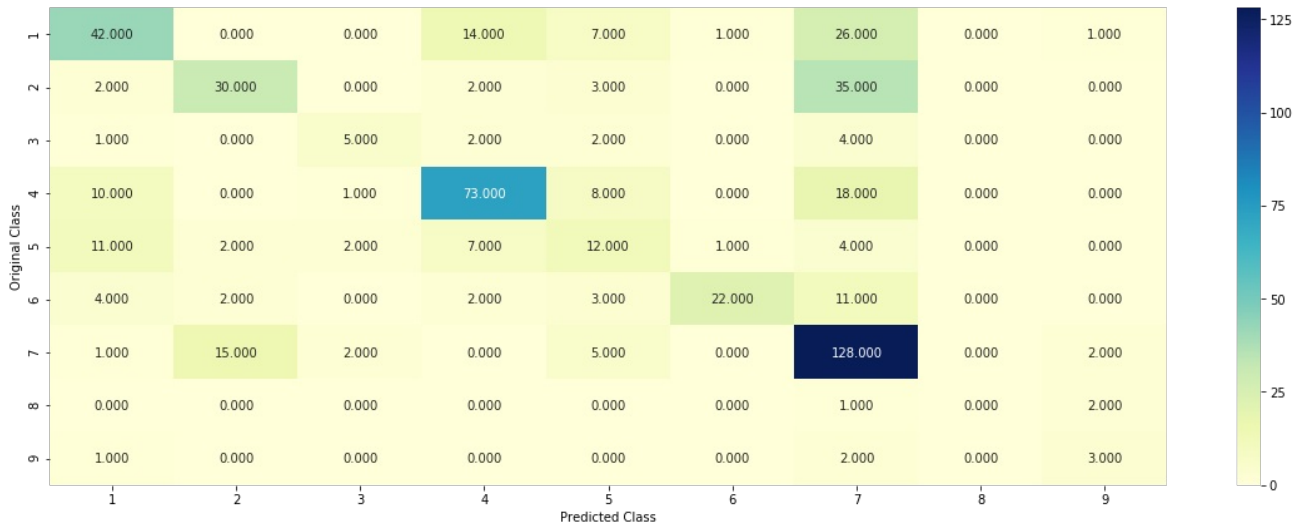
```
For values of best alpha = 0.01 The train log loss is: 0.7462651006935817
For values of best alpha = 0.01 The cross validation log loss is: 1.2501309862948178
For values of best alpha = 0.01 The test log loss is: 1.250034641514856
```

4.3.1.2. Testing the model with best hyper paramters

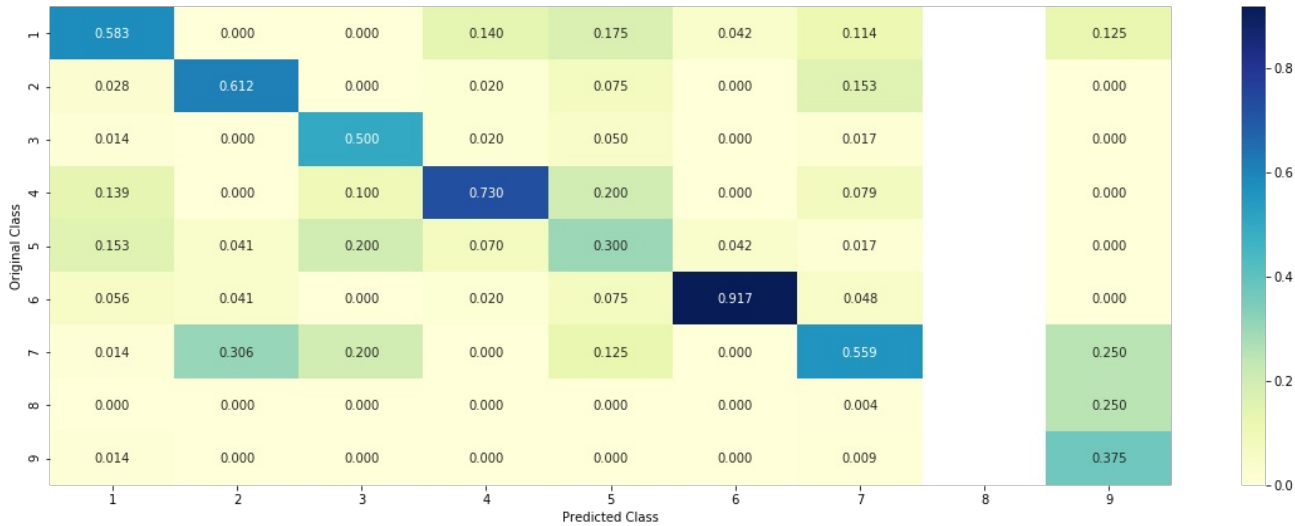
In [66]:

```
#clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_ngrams, train_y, cv_x_ngrams, cv_y, sig_clf)
```

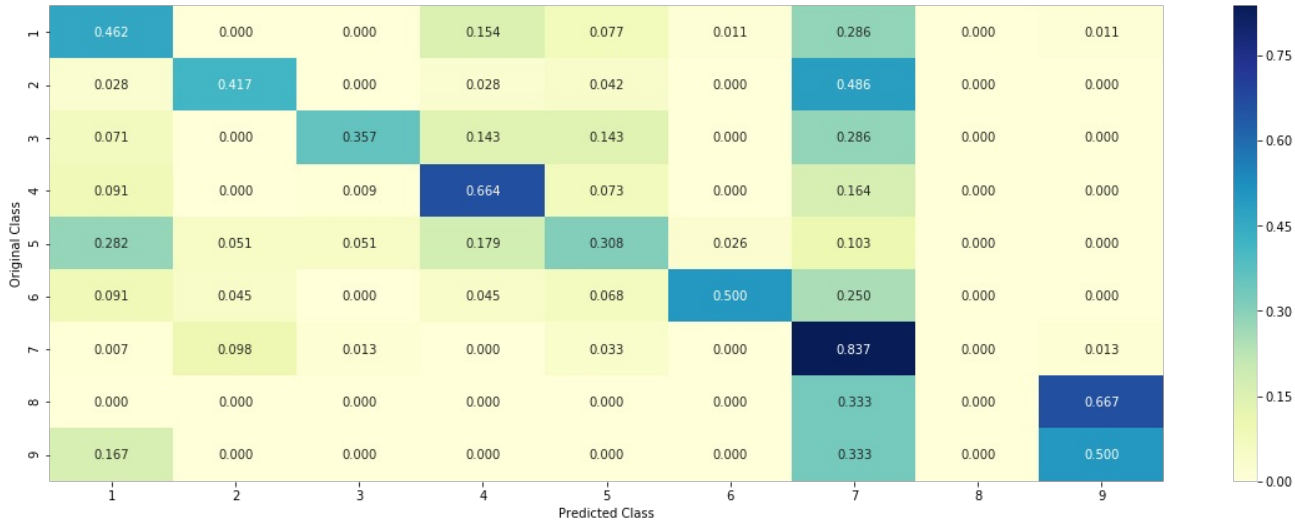
Log loss : 1.2135678532698875  
Percentage of mis-classified points : 40.78947368421053  
----- Confusion matrix -----



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



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



4.3.1.3. Feature Importance



In [117]:

```
gene_count_vec = CountVectorizer(ngram_range=(1,2))
var_count_vec = CountVectorizer(ngram_range=(1,2))
text_count_vec = CountVectorizer(ngram_range=(1,2))

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())
fea1_len
```

Out[117]:

233

In [67]:

```
# 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(ngram_range=(1,2))
    var_count_vec = CountVectorizer(ngram_range=(1,2))
    text_count_vec = CountVectorizer(ngram_range=(1,2))

    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()) #233
    fea2_len = len(var_count_vec.get_feature_names()) #2056

    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")
```

#### 4.3.1.3.1. Correctly Classified point

In [122]:

```
test_point_index = 55
no_feature = 5000
predicted_cls = sig_clf.predict(test_x_ngrams[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_ngrams[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-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['Variation'].iloc[test_point_index], no_feature)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0917 0.0706 0.012  0.0687 0.0368 0.0237 0.6888 0.0036 0.004 ]]
Actual Class : 7
-----
```

#### 4.3.1.3.2. Incorrectly Classified point



In [73]:

```
test_point_index = 1
no_feature = 5000
predicted_cls = sig_clf.predict(test_x_ngrams[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_ngrams[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-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['Variation'].iloc[test_point_index], no_feature)
```

Predicted Class : 7  
Predicted Class Probabilities: [[0.0319 0.0019 0.0062 0.0111 0.0126 0.004 0.9234 0.0046 0.0043]]  
Actual Class : 2  
-----

## 4.3.2 Logistic Regression Without Class balancing

### 4.3.2.1. Hyper paramter tuning

In [74]:

```
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_ngrams, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_ngrams, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_ngrams)
    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()

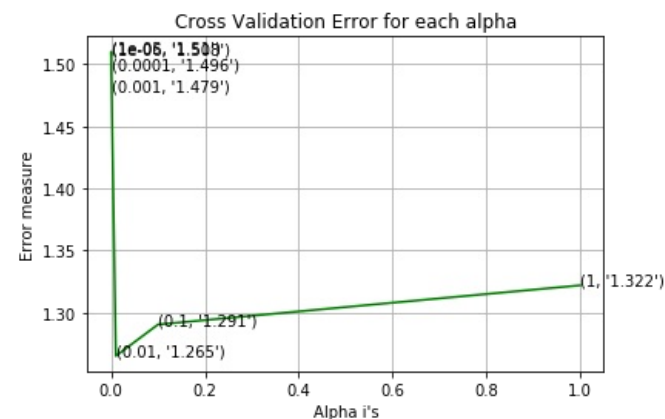
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_ngrams, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_ngrams, train_y)

predict_y = sig_clf.predict_proba(train_x_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_ngrams)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```

for alpha = 1e-06
Log Loss : 1.507948803787281
for alpha = 1e-05
Log Loss : 1.5101009495762914
for alpha = 0.0001
Log Loss : 1.495701978727535
for alpha = 0.001
Log Loss : 1.4790303612474258
for alpha = 0.01
Log Loss : 1.2652506481375811
for alpha = 0.1
Log Loss : 1.2905828012502156
for alpha = 1
Log Loss : 1.3220584787281902

```



```

For values of best alpha = 0.01 The train log loss is: 0.7422312510667094
For values of best alpha = 0.01 The cross validation log loss is: 1.2652506481375811
For values of best alpha = 0.01 The test log loss is: 1.2766745775786867

```

#### 4.3.2.2. Testing model with best hyper parameters

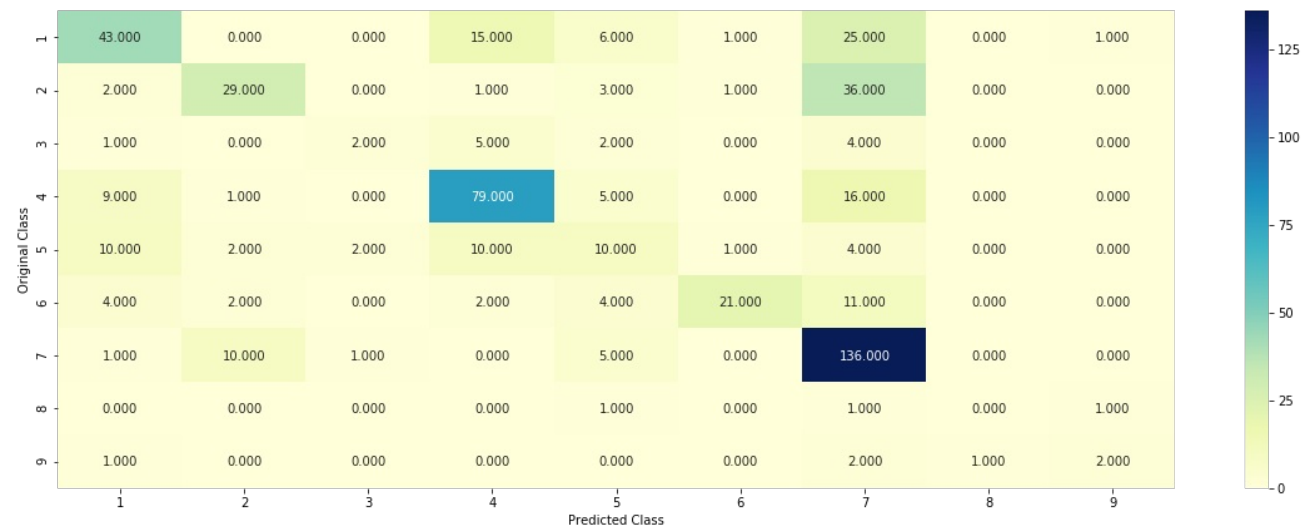
In [75]:

```
predict_and_plot_confusion_matrix(train_x_ngrams, train_y, cv_x_ngrams, cv_y, sig_clf)
```

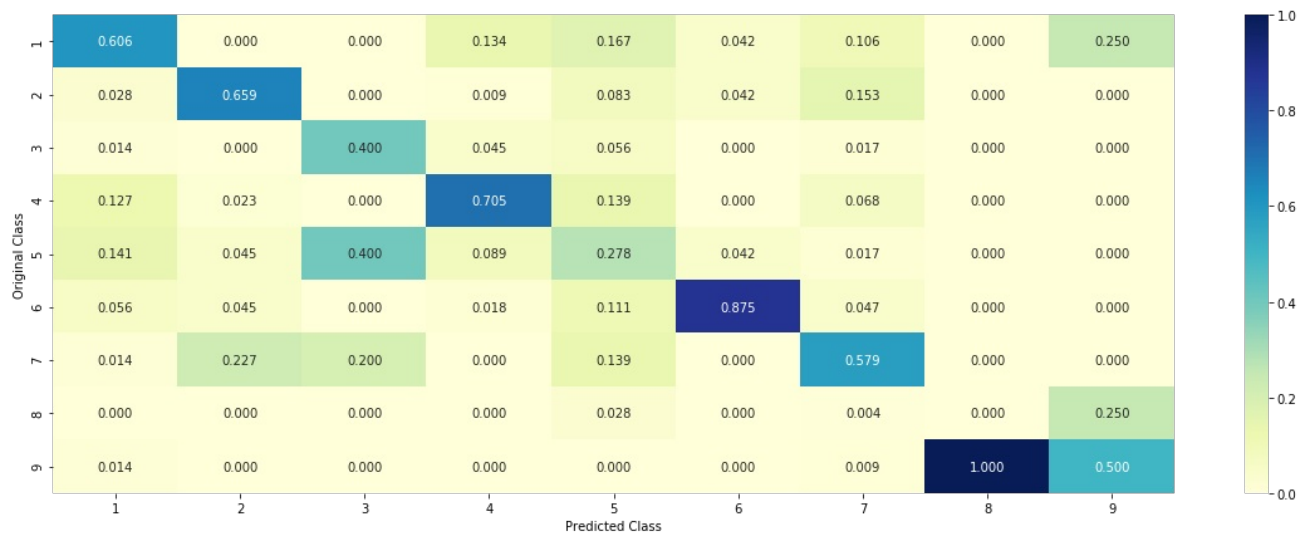
```

Log loss : 1.226452619392822
Percentage of mis-classified points : 39.473684210526315
----- Confusion matrix -----

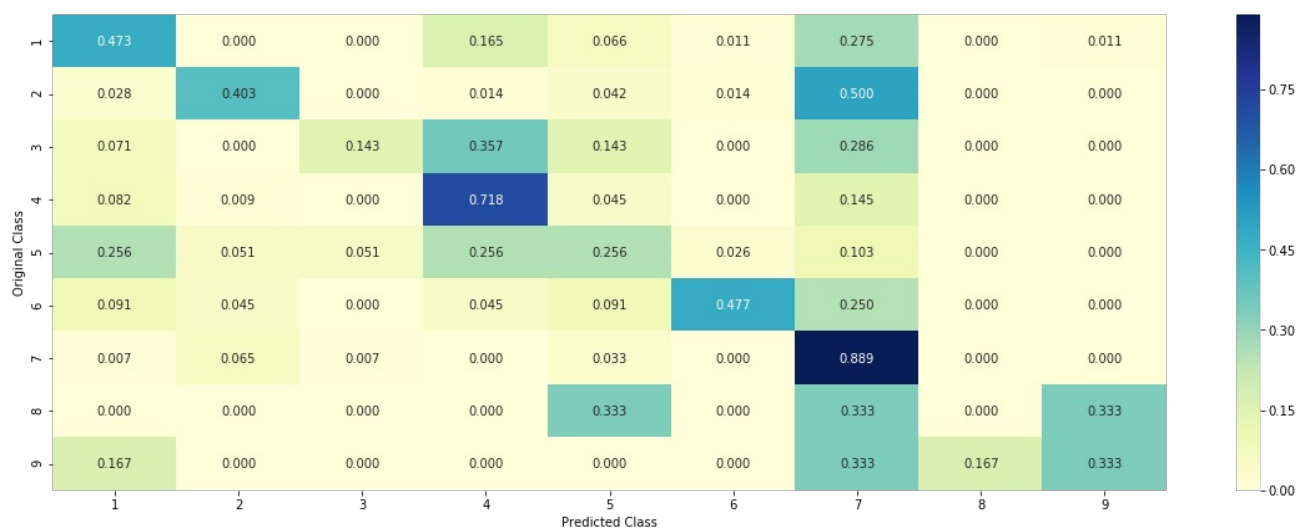
```



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



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



#### 4.3.2.3. Feature Importance, Correctly Classified point

In [76]:

```
test_point_index = 1
no_feature = 500
predicted_cls = sig_clf.predict(test_x_ngrams[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_ngrams[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-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['Variation'].iloc[test_point_index], no_feature)
```

Predicted Class : 7  
 Predicted Class Probabilities: [[3.500e-02 2.800e-03 1.500e-03 1.770e-02 9.800e-03 2.300e-03 9.288e-01  
 1.200e-03 8.000e-04]]  
 Actual Class : 2

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

#### 4.3.2.4. Feature Importance, Incorrectly Classified point

In [77]:

```
test_point_index = 4
no_feature = 500
predicted_cls = sig_clf.predict(test_x_ngrams[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_ngrams[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-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['Variation'].iloc[test_point_index], no_feature)
```

```
Predicted Class : 1
Predicted Class Probabilities: [[0.3183 0.1043 0.0114 0.2137 0.189  0.0243 0.1307 0.0042 0.0042]]
Actual Class : 5
-----
Out of the top 500 features 0 are present in query point
```

## What we did so far?

The task for this experiment is to correctly classify a given genetic variations/mutations into one of the nine classes, based on evidences gathered from text-based clinical literature. There was a no latency requirement, which means it is okay to take some time before arriving at a conclusion. Here, we should keep in mind that making an incorrect assumption can be costly, as it will severely affect a patient's line of treatment. We will build a classifier which will give probabilistic outputs instead of just giving us the class labels.

Our dataset for this task consists of three features - Gene, Variants, Text. We will use these three features to build a cancer detection model.

In the EDA sections there are some very useful information about the structure of the given data. We have 3321 data points unequally distributed among 9 classes. The distribution of the number of sentences tells us that there are almost 50% data points which has 427 sentences at max. Almost 75% data points have 744 number of sentences. There's however a sudden increase beyond the 75th percentile mark. We will keep this in mind as we move along.

As a general rule of thumb, the pre-processing is done by removing unwanted words, characters and converting the entire text to lower case. Converting the text to lower case is important because, the model will treat the capital and small letters as different words even if they are the same!

I have some new features. These are, MergedText - Each of the variation data point merges with the the corresponding text data into one single string. TopWordsInText - This feature will give us the top 75 most frequently occurring words in a given text data. We will use these features to get tfidf representations to train our model. In feature engineering part 1, we will use the gene and variations data to build a tfidf-vectorizer. With this tfidf-vectorizer we will encode the text data, and use this as a feature. This is very relevant because, we have already seen that gene and variation itself brought down the log loss to almost half of a random model.

We need to be careful while splitting the data into train, test and cross validation datasets. This is because we want the three datasets to have almost equal distributions of classes. We can draw a simple plot to check these distributions.

Now, it's time for use to see what are the top most frequently occurring words for a given class. I have used word clouds for this purpose. The word clouds can straight away give us relevant keywords which helps us understand what type of words are mots common for every class. For example we see keywords like "breast cancer", "lung cancer", "ovarian cancer", "gene mutation", "tyrosie kinase", "egfr mutation","amino acids" and so on, belonging to each of the classes. Here I have given just a few examples. We also see that there are some classes which talks about some particular variation/mutant type more frequently than others.

Why have we chosen our metric to be 'log loss' ?

As we have discussed above, we will use our Key Performance Indicator to be 'log' loss. Minimising the Log Loss is basically equivalent to maximising the accuracy of the classifier. In order to calculate Log loss, the classifier must actually assign a probability value for each of the class labels. Ideally as the predicted class probabilities improve, the Log loss keeps on reducing. Log loss penalises the classifier very heavily if it classifies a Class 1 to be Class 4 and vice versa. For example, if for a particular observation, the classifier assigns a very small probability to the correct class then the corresponding contribution to the Log Loss will be very large. Naturally this is going to have a significant impact on the overall Log Loss for the classifier, which will then become higher. But, in other scenario if the classifier assigns a high probability value to the correct class, then the Log loss will reduce significantly. Now, imagine a scenario where the model doesn't predict probability scores. It only predicts class labels. In this case, class labels can be either 0 or 1. So we won't get a deep understanding or any interpretability about why a particular pair of question has been labeled as either 0 or 1. Chosing the KPI as Log loss gives us this freedom. We can interpret the models. We can say this two questions are 95% similar or 80% similar, instead of just bluntly classifying them as duplicates.

For deep understanding of log loss please visit: <https://datawookie.netlify.com/blog/2015/12/making-sense-of-logarithmic-loss/> (<https://datawookie.netlify.com/blog/2015/12/making-sense-of-logarithmic-loss/>) (<https://datawookie.netlify.com/blog/2015/12/making-sense-of-logarithmic-loss/>) (<https://datawookie.netlify.com/blog/2015/12/making-sense-of-logarithmic-loss/>)

Now, the minimum value of log loss in the ideal scenario would be 0. But what about the maximum? For this we will build a random model. This random model, which is like the worst model that can be made for this problem, gave us a log-loss of almost 2.5. This gives us an upper limit. We now know what is the worst log loss that our model can give? With this in mind we will try to bring the log loss as

close to zero as possible.

Typically for high dimensional data, logistic regression and linear SVMs work very well. We can use KNN, Naive Bayes for lower dimensional data representations. In a separate experiment we can also try models like XGBoost and LightGBM to extract features. Anyway, as the first step we need to encode the text feature into numerical data. How? We have two ways - a simple one hot encoding of the text corpus using Bag of Words approach and a response coding approach. 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.

### **Very Important note about response coding: Avoiding response leakage.**

1. We have to be extremely careful not to use the test and cross validate data for response coding. This is because we don't the issue of data leakage.
2. Suppose we have a variant V2 present in Test / Cross Val dataset. But V2 is not present in Train. So in that case, while building the response bales for V2, we will just assign equal probability values to each of the 9 array values. Proba =  $1/9$  for each og them.
3. We will take the help of laplace smoothing in order to achieve this. Without laplace smoothing, we would get a 0/0 error.
4. We are seeing the data that is present in the cross validation and test data during the time of training. So we are literally leaking the information that is present in test/cv data at the time of training.

This should be strictly avoided, as we do not want a data leakage issue.

After all this is done, we will perform univariate analysis using all the features. As we can see, all the individual features brought down the log loss by a significant amount. Now we will combine these three features in various ways to get the best possible log loss that we could get.

The other metric we have chosen for this problem is the confusion matrix. By suing the confusion matrix, precision matrix and the recall matrix we can actually see what percentage of points are correctly classified in each of the 9 classes.

We will use various models like naivae bayes, k-nearest neighbors, logistic regression, random forest etc. Before we begin training our models we have to make sure that we build our model such that it is interpretable. Not only the model will tell us the class type, but it will also tell us the exact reason why it thinks a given query point belongs to a certain class.

First we will use TFIDF features and run all the models. Here we see that the best log loss we get is from the logistic regression model with balanced class values. Instead of using all the words in the dataset, we will use only the top 1000 words which occur in the text feature.

Now, since we know that logistic regression performs better for high dimensional data, we will repeat our experiment by encoding the text to BOW representations using bigrams and unigrams. We have seen in the word clouds before how unigrams and bigrams are so important in determining the context of any class. The BOW features with bigrams performed fairly well. But, there wasn't any drastic improvement in the log loss.

With the initial first cut solution, we will move onto the feature engineering stage, where we will perform hyperparameter tuning and try and reduce the log loss as much as possibele.

In stage 1 of feature engineering, we will combine the words present in the Gene column with those present in the Variation column. We will build a tfidf vectorizer using this corpus and use this to transform the text column. We will use this as our 4th feature and see a how a logistic regression performs on this model. Using this feature we were able to reduce the log loss to almost 0.95, which is a significant drop.

In feature engineering stage 2, we will add a 5th feature. We will take the 75 most frequently occurring words in each text and use this corpus to build a tfidf vector representation of the top words. Using this method, I was able to bring down the log loss to 0.92 with a 31% misclassification error. This is by far the best model that we have seen.

Given below are the model performances for all the models tried so far.

### **Comparing model performances:**

#### **TF-IDF Feature representations:**

In [19]:

```
from prettytable import PrettyTable
table = PrettyTable()
table.field_names = ["Model", "Vectorizer", "CV Log Loss", "Test Log Loss", "Misclassification Error"]
table.add_row(["Random Model", 'Response Coded', "~2.5", "~2.5", ">65%"])
table.add_row(["Naive Bayes", 'TFIDF', 0.525, 1.146, "36.48%"])
table.add_row(["KNN", 'TFIDF', 0.642, 1.032, "36.27%"])
table.add_row(["LR (Class Balancing)", 'TFIDF', 0.442, 0.979, "34.27%"])
table.add_row(["LR (No Class Balancing)", 'TFIDF', 0.433, 1.005, "35.15%"])
table.add_row(["Linear-SVM", 'TFIDF', 0.504, 1.0420, "35.52%"])
table.add_row(["Random Forest (OneHot)", 'TFIDF', 0.852, 1.175, "41.35%"])
table.add_row(["Random Forest (ResponseCode)", 'TFIDF', 0.059, 1.257, "41.91%"])
table.add_row(["StackingClassfier", 'TFIDF', 0.533, 1.135, "37.29%"])
table.add_row(["MaximumVotingClassfier", 'TFIDF', 0.833, 1.173, "37.44%"])
print(table)
```

Model	Vectorizer	CV Log Loss	Test Log Loss	Misclassification Error
Random Model	Response Coded	~2.5	~2.5	>65%
Naive Bayes	TFIDF	0.525	1.146	36.48%
KNN	TFIDF	0.642	1.032	36.27%
LR (Class Balancing)	TFIDF	0.442	0.979	34.27%
LR (No Class Balancing)	TFIDF	0.433	1.005	35.15%
Linear-SVM	TFIDF	0.504	1.042	35.52%
Random Forest (OneHot)	TFIDF	0.852	1.175	41.35%
Random Forest (ResponseCode)	TFIDF	0.059	1.257	41.91%
StackingClassfier	TFIDF	0.533	1.135	37.29%
MaximumVotingClassfier	TFIDF	0.833	1.173	37.44%

**BOW (bigrams) Feature representations:**

In [18]:

```
from prettytable import PrettyTable
table = PrettyTable()
table.field_names = ["Model", "Vectorizer", "CV Log Loss", "Test Log Loss", "Misclassification Error"]
table.add_row(["Random Model", 'Response Coded', "~2.5", "~2.5", ">65%"])

table.add_row(["LR (Class Balancing)", 'BOW', 0.746, 1.25, "40.78%"])
table.add_row(["LR (No Class Balancing)", 'BOW', 0.742, 1.277, "39.47%"])
table.add_row(["Linear-SVM", 'BOW', 0.504, 1.0420, "35.52%"])

print(table)
```

Model	Vectorizer	CV Log Loss	Test Log Loss	Misclassification Error
Random Model	Response Coded	~2.5	~2.5	>65%
LR (Class Balancing)	BOW	0.746	1.25	40.78%
LR (No Class Balancing)	BOW	0.742	1.277	39.47%
Linear-SVM	BOW	0.504	1.042	35.52%

**Feature Engineering Results:**

In [17]:

```
from prettytable import PrettyTable
table = PrettyTable()
table.field_names = ["Model", "Vectorizer", "CV Log Loss", "Test Log Loss", "Misclassification Error"]
table.add_row(["Random Model", 'Response Coded', "~2.5", "~2.5", ">65%"])

table.add_row(["LR (Class Balancing) (FE1)", 'TFIDF', 0.456, 0.970, "33.23%"])
table.add_row(["LR (No Class Balancing) (FE1)", 'TFIDF', 0.445, 0.963, "33.24%"])
table.add_row(["LR (Class Balancing) (FE2)", 'TFIDF', 0.448, 0.925, "32.48%"])
table.add_row(["LR (No Class Balancing) (FE2)", 'TFIDF', 0.439, 0.92, "32.03%"])

print(table)
```

Model	Vectorizer	CV Log Loss	Test Log Loss	Misclassification Error
Random Model	Response Coded	~2.5	~2.5	>65%
LR (Class Balancing) (FE1)	TFIDF	0.456	0.97	33.23%
LR (No Class Balancing) (FE1)	TFIDF	0.445	0.963	33.24%
LR (Class Balancing) (FE2)	TFIDF	0.448	0.925	32.48%
LR (No Class Balancing) (FE2)	TFIDF	0.439	0.92	32.03%