

# PCD\_BOW\_ngram\_LR

July 21, 2020

Personalized cancer diagnosis

## 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>
2. <https://www.youtube.com/watch?v=UwbuW7oK8rk>
3. <https://www.youtube.com/watch?v=qxXRKVompI8>

### 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>
- 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 are have a common column called ID
- Data file's information:

```
<li>
training_variants (ID , Gene, Variations, Class)
</li>
<li>
training_text (ID, Text)
</li>
```

### 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

### 2.2.2. Performance Metric

Source: <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

```
[ ]: from google.colab import drive
drive.mount('/content/drive')
```

Go to this URL in a browser: [https://accounts.google.com/o/oauth2/auth?client\\_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect\\_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aoob&response\\_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly](https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aoob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly)

Enter your authorization code:

ûûûûûûûûûûûû

Mounted at /content/drive

```
[ ]: import os
os.chdir("/content/drive/My Drive/PCD")
!ls -l
```

total 207220

-rw----- 1 root root 212125752 Jun 20 2018 training\_text

-rw----- 1 root root 66688 Jun 23 2017 training\_variants

## 3. Exploratory Data Analysis

```
[ ]: import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import nltk
nltk.download("stopwords")
import warnings
import numpy as np
```

```

from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import SGDClassifier
from imblearn.over_sampling import SMOTE
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from sklearn.model_selection import StratifiedKFold
from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier
warnings.filterwarnings("ignore")

from mlxtend.classifier import StackingClassifier

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression

```

[nltk\_data] Downloading package stopwords to /root/nltk\_data...

[nltk\_data] Unzipping corpora/stopwords.zip.

### 3.1. Reading Data

#### 3.1.1. Reading Gene and Variation Data

```

[ ]: data = pd.read_csv('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']

```
[ ]:  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 Fields are

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

### 3.1.2. Reading Text Data

```
[ ]: # note the separator in this file
data_text = pd.
    ↳ read_csv("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']
```

```
[ ]:  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

```
[ ]: # 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

```

```

[:]: #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 : 28.683154000000002 seconds

```

```

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

```

```

[:]:
  ID    Gene  ... Class                                TEXT
0  0  FAM58A  ...    1  cyclin dependent kinases cdks regulate variety...
1  1    CBL  ...    2  abstract background non small cell lung cancer...
2  2    CBL  ...    2  abstract background non small cell lung cancer...
3  3    CBL  ...    3  recent evidence demonstrated acquired uniparen...
4  4    CBL  ...    4  oncogenic mutations monomeric casitas b lineag...

[5 rows x 5 columns]

```

```

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

```

```

[:]:
  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

```

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

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

```
[ ]:      ID  Gene Variation  Class      TEXT
      1109  1109  FANCA     S1088F      1  FANCA S1088F
```

### 3.1.4. Test, Train and Cross Validation Split

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

```
[ ]: 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

```
[ ]: 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

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

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

```

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

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

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

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

# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:

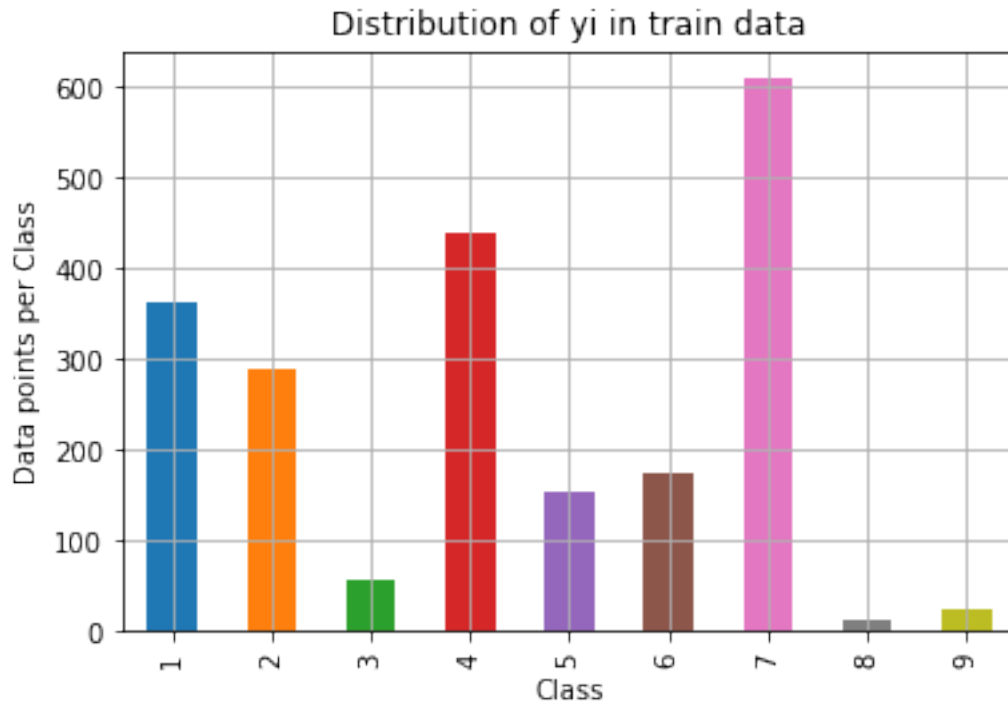
```



```

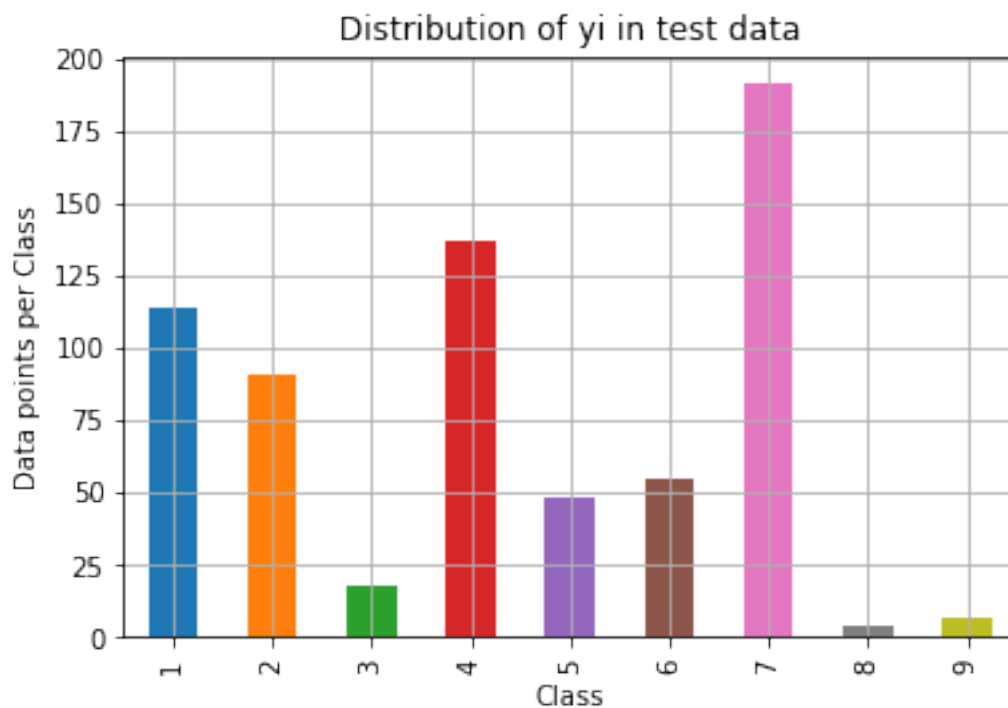
print('Number of data points in class', i+1, ':', cv_class_distribution.
→values[i], '(', np.round((cv_class_distribution.values[i]/cv_df.
→shape[0]*100), 3), '%)')

```



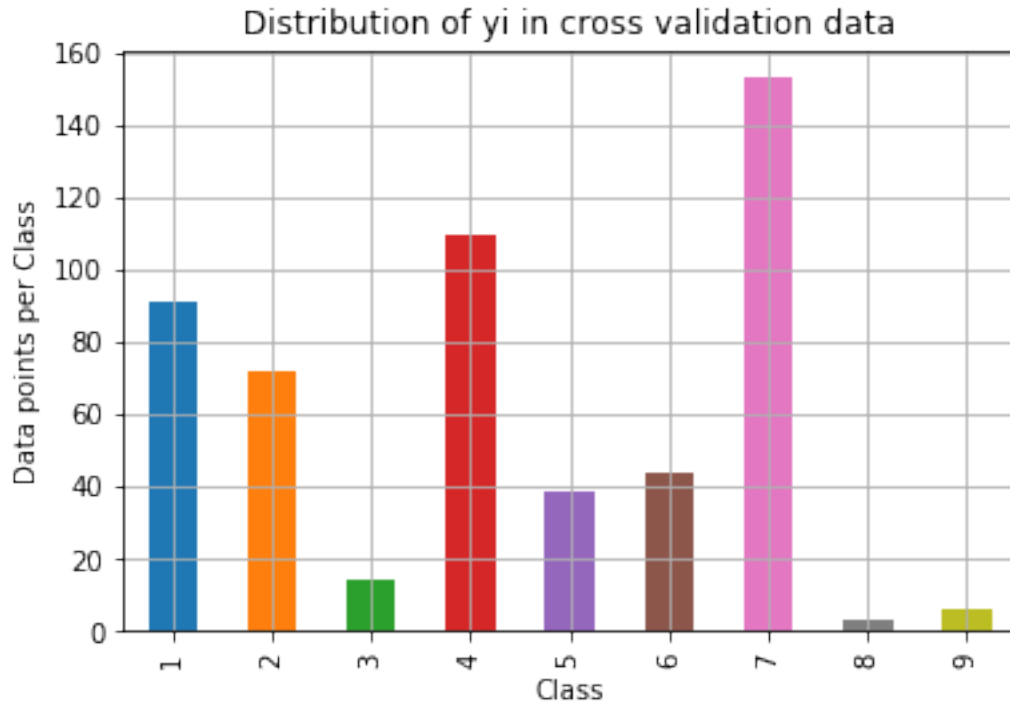
Number of data points in class 7 : 609 ( 28.672 %)  
 Number of data points in class 4 : 439 ( 20.669 %)  
 Number of data points in class 1 : 363 ( 17.09 %)  
 Number of data points in class 2 : 289 ( 13.606 %)  
 Number of data points in class 6 : 176 ( 8.286 %)  
 Number of data points in class 5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

### 3.2 Prediction using a 'Random' Model

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

```

[ ]: # 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)
    # divide 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]]

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()
```

```
[ ]: # we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to generate 9 numbers and divide each of the numbers by their
    → sum
# ref: https://stackoverflow.com/a/18662466/4084039
test_data_len = test_df.shape[0]
cv_data_len = cv_df.shape[0]

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

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

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

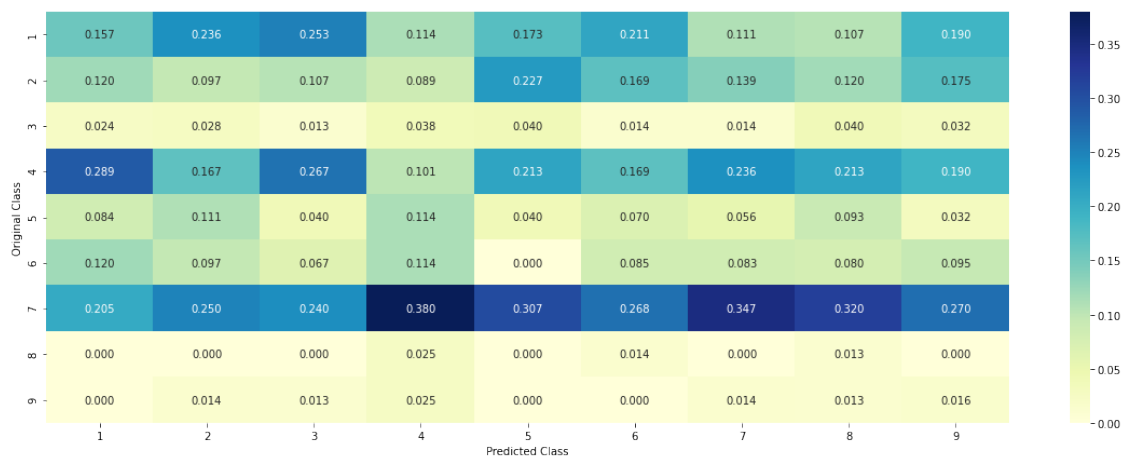
Log loss on Cross Validation Data using Random Model 2.5088038472695926

Log loss on Test Data using Random Model 2.487897745132674

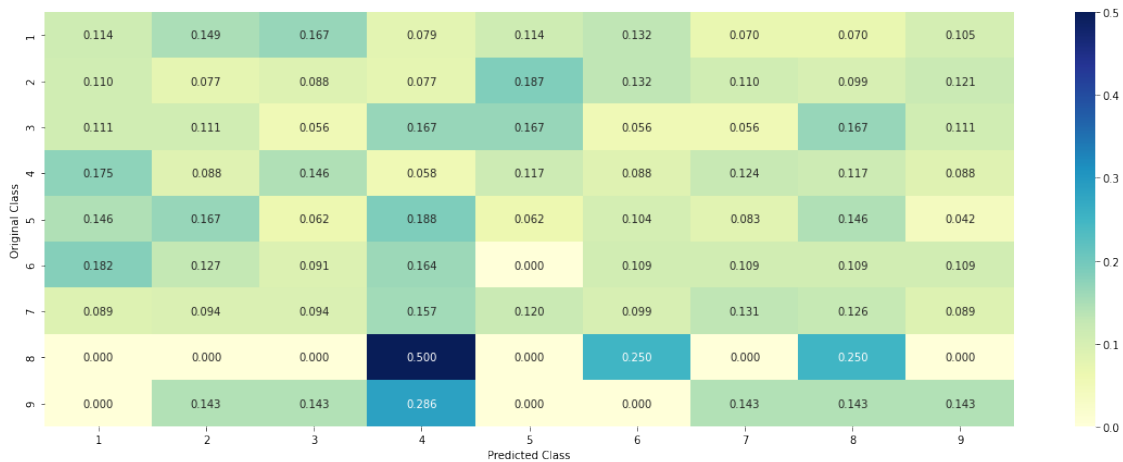
----- Confusion matrix -----



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



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



### 3.3 Univariate Analysis

```
[ ]: # code for response coding with Laplace smoothing.
# alpha : used for laplace smoothing
# feature: ['gene', 'variation']
# df: ['train_df', 'test_df', 'cv_df']
# algorithm
# -----
# Consider all unique values and the number of occurrences of given feature in
# → train data dataframe
# build a vector (1*9) , the first element = (number of times it occurred in
# → class1 + 10*alpha / number of 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
→occured in whole data
for i, denominator in value_count.items():
    # vec will contain (p(yi==1/Gi) probability of gene/variation belongs
→to perticular class
    # vec is 9 diamensional vector
    vec = []
    for k in range(1,10):
        # print(train_df.loc[(train_df['Class']==1) &
→(train_df['Gene']=='BRCA1')])
        #
        # ID   Gene          Variation   Class
        # 2470  2470  BRCA1          S1715C      1
        # 2486  2486  BRCA1          S1841R      1
        # 2614  2614  BRCA1           M1R        1
        # 2432  2432  BRCA1          L1657P      1
        # 2567  2567  BRCA1          T1685A      1
        # 2583  2583  BRCA1          E1660G      1
        # 2634  2634  BRCA1          W1718L      1
        # cls_cnt.shape[0] will return the number of rows

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

```



```

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    #      {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.
    → 0681818181818177, 0.13636363636363635, 0.25, 0.19318181818181818, 0.
    → 03787878787878788, 0.03787878787878788, 0.03787878787878788],
    #      'TP53': [0.32142857142857145, 0.061224489795918366, 0.
    → 061224489795918366, 0.27040816326530615, 0.061224489795918366, 0.
    → 066326530612244902, 0.051020408163265307, 0.051020408163265307, 0.
    → 056122448979591837],
    #      'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.
    → 0681818181818177, 0.0681818181818177, 0.0625, 0.346590909090912, 0.
    → 0625, 0.056818181818181816],
    #      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.
    → 060606060606060608, 0.078787878787878782, 0.1393939393939394, 0.
    → 34545454545454546, 0.060606060606060608, 0.060606060606060608, 0.
    → 060606060606060608],
    #      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.
    → 069182389937106917, 0.46540880503144655, 0.075471698113207544, 0.
    → 062893081761006289, 0.069182389937106917, 0.062893081761006289, 0.
    → 062893081761006289],
    #      'KIT': [0.066225165562913912, 0.25165562913907286, 0.
    → 072847682119205295, 0.072847682119205295, 0.066225165562913912, 0.
    → 066225165562913912, 0.27152317880794702, 0.066225165562913912, 0.
    → 066225165562913912],
    #      'BRAF': [0.066666666666666666, 0.17999999999999999, 0.
    → 073333333333333334, 0.073333333333333334, 0.093333333333333338, 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 will 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] to gv_fea
for index, row in df.iterrows():
    if row[feature] in dict(value_count).keys():
        gv_fea.append(gv_dict[row[feature]])
    else:
        gv_fea.append([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
#
    gv_fea.append([-1,-1,-1,-1,-1,-1,-1,-1,-1])
return gv_fea

```

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

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

### 3.2.1 Univariate Analysis on Gene Feature

Q1. Gene, What type of feature it is ?

Ans. Gene is a categorical variable

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

```

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

```

Number of Unique Genes : 240

|        |     |
|--------|-----|
| BRCA1  | 168 |
| TP53   | 103 |
| EGFR   | 98  |
| BRCA2  | 84  |
| PTEN   | 76  |
| BRAF   | 63  |
| KIT    | 57  |
| ALK    | 45  |
| PIK3CA | 41  |
| ERBB2  | 39  |

Name: Gene, dtype: int64

```

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

```

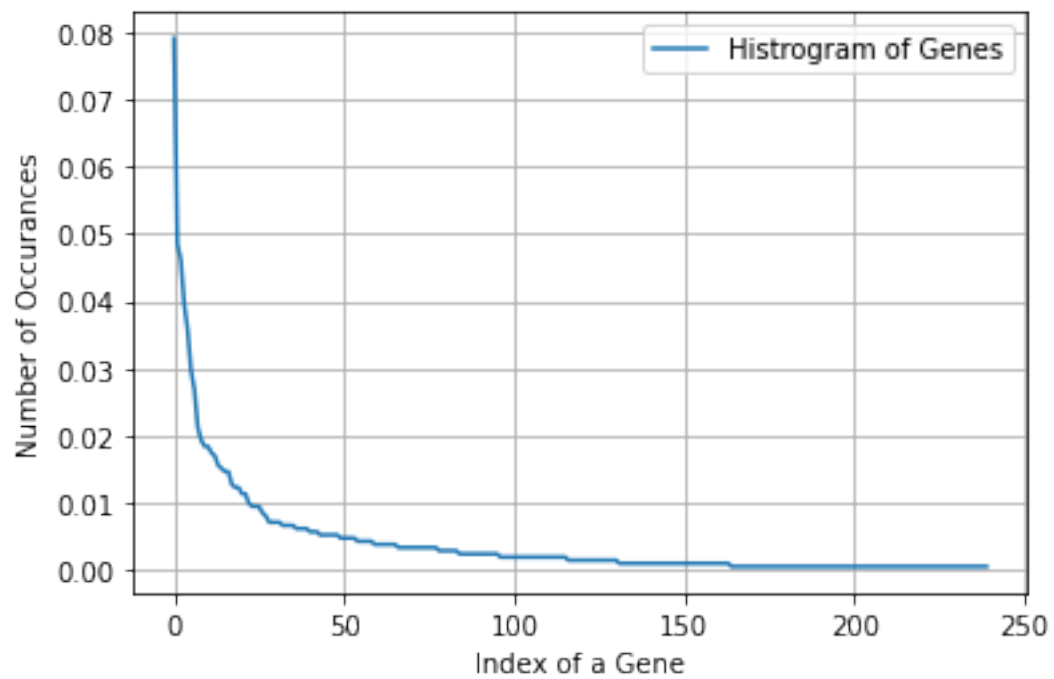
Ans: There are 240 different categories of genes in the train data, and they are distributed as follows

```

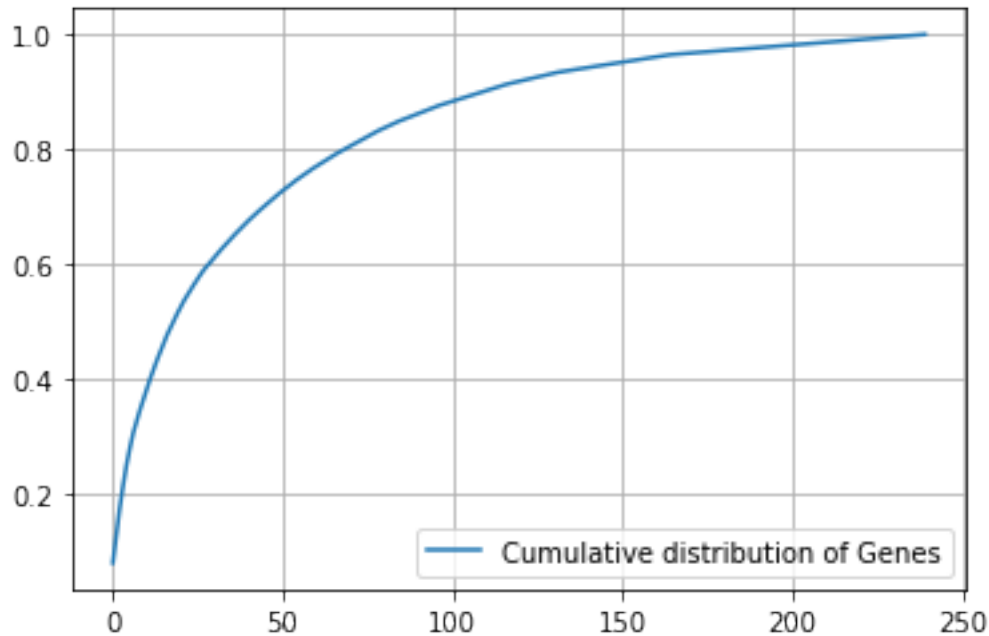
[:]: s = sum(unique_genes.values);
h = unique_genes.values/s;
plt.plot(h, label="Histogram of Genes")
plt.xlabel('Index of a Gene')

```

```
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```



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



Q3. How to featurize this Gene feature ?

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

One hot Encoding

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.

```
[ ]: #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))

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

train\_gene\_feature\_responseCoding is converted feature using response coding method. The shape of gene feature: (2124, 9)

```
[ ]: # one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer(ngram_range=(1,3))
train_gene_feature_onehotCoding = gene_vectorizer.
    ↳fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

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

```
[ ]: 2537    BRCA1
      2449    BRCA1
      2846    BRCA2
      1755    IDH1
      991     TSC1
      Name: Gene, dtype: object
```

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

```
[ ]: ['abl1',
      'acvr1',
      'ago2',
      'akt1',
      'akt2',
      'akt3',
      'alk',
      'apc',
      'ar',
      'araf',
      'arid1b',
      'arid2',
      'arid5b',
      'asxl1',
      'asxl2',
      'atm',
      'atrx',
      'aurkb',
      'axin1',
      'axl',
      'b2m',
      'bap1',
      'bard1',
      'bcl10',
      'bcl2',
      'bcl2l11',
      'bcor',
      'braf',
      'brca1',
      'brca2',
      'brd4',
```

'brip1',  
'btk',  
'card11',  
'carm1',  
'casp8',  
'cbl',  
'ccnd1',  
'ccnd2',  
'ccnd3',  
'cdh1',  
'cdk12',  
'cdk4',  
'cdk6',  
'cdk8',  
'cdkn1a',  
'cdkn1b',  
'cdkn2a',  
'cdkn2b',  
'cdkn2c',  
'cebpa',  
'chek2',  
'cic',  
'crebbp',  
'ctcf',  
'ctla4',  
'ctnnb1',  
'ddr2',  
'dicer1',  
'dnmt3a',  
'dnmt3b',  
'egfr',  
'elf3',  
'ep300',  
'epas1',  
'epcam',  
'erbb2',  
'erbb3',  
'erbb4',  
'ercc2',  
'ercc4',  
'erg',  
'errfi1',  
'esr1',  
'etv1',  
'etv6',  
'ewsr1',  
'ezh2',

'fam58a',  
'fanca',  
'fat1',  
'fbxw7',  
'fgf19',  
'fgf3',  
'fgfr1',  
'fgfr2',  
'fgfr3',  
'fgfr4',  
'flt1',  
'flt3',  
'foxa1',  
'foxl2',  
'foxp1',  
'fubp1',  
'gata3',  
'gli1',  
'gna11',  
'gnaq',  
'gnas',  
'h3f3a',  
'hist1h1c',  
'hla',  
'hnf1a',  
'hras',  
'idh1',  
'idh2',  
'igf1r',  
'ikzf1',  
'il7r',  
'jak1',  
'jak2',  
'jun',  
'kdm5a',  
'kdm5c',  
'kdm6a',  
'kdr',  
'keap1',  
'kit',  
'klf4',  
'kmt2a',  
'kmt2b',  
'kmt2c',  
'kmt2d',  
'knstrn',  
'kras',

'lats2',  
'map2k1',  
'map2k2',  
'map2k4',  
'map3k1',  
'mapk1',  
'mdm2',  
'med12',  
'mef2b',  
'men1',  
'met',  
'mlh1',  
'mpl',  
'msh2',  
'msh6',  
'mtor',  
'myc',  
'mycn',  
'myd88',  
'myod1',  
'ncor1',  
'nf1',  
'nf2',  
'nfe2l2',  
'nfkb1a',  
'nkx2',  
'notch1',  
'notch2',  
'npm1',  
'nras',  
'nsd1',  
'ntrk1',  
'ntrk2',  
'ntrk3',  
'nup93',  
'pax8',  
'pbrm1',  
'pdgfra',  
'pdgfrb',  
'pik3ca',  
'pik3cb',  
'pik3cd',  
'pik3r1',  
'pik3r2',  
'pim1',  
'pms2',  
'pole',



'ppm1d',  
'ppp2r1a',  
'ppp6c',  
'prdm1',  
'pten',  
'ptpn11',  
'ptprd',  
'ptprt',  
'rab35',  
'rac1',  
'rad21',  
'rad50',  
'rad51b',  
'rad51c',  
'rad54l',  
'raf1',  
'rara',  
'rasa1',  
'rb1',  
'rbm10',  
'ret',  
'rheb',  
'rhoa',  
'rictor',  
'rit1',  
'rnf43',  
'ros1',  
'runx1',  
'rxra',  
'rybp',  
'setd2',  
'sf3b1',  
'shq1',  
'smad2',  
'smad3',  
'smad4',  
'smarca4',  
'smarcb1',  
'smo',  
'sos1',  
'sox9',  
'spop',  
'src',  
'srsf2',  
'stag2',  
'stat3',  
'stk11',

```
'tcf3',
'tcf7l2',
'tert',
'tet1',
'tet2',
'tgfbr1',
'tgfbr2',
'tmprss2',
'tp53',
'tp53bp1',
'tsc1',
'tsc2',
'u2af1',
'vegfa',
'vhl',
'whsc1',
'whsc1l1',
'xpo1',
'xrcc2',
'yap1']
```

```
[ ]: print("train_gene_feature_onehotCoding is converted feature using one-hot_
→encoding method. The shape of gene feature:",
→train_gene_feature_onehotCoding.shape)
```

train\_gene\_feature\_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 239)

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  $y_i$ .

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

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
→generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
→fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
→learning_rate=optimal, eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ])          Fit linear model with
→Stochastic Gradient Descent.
```

```

# predict(X)          Predict class labels for samples in X.

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

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

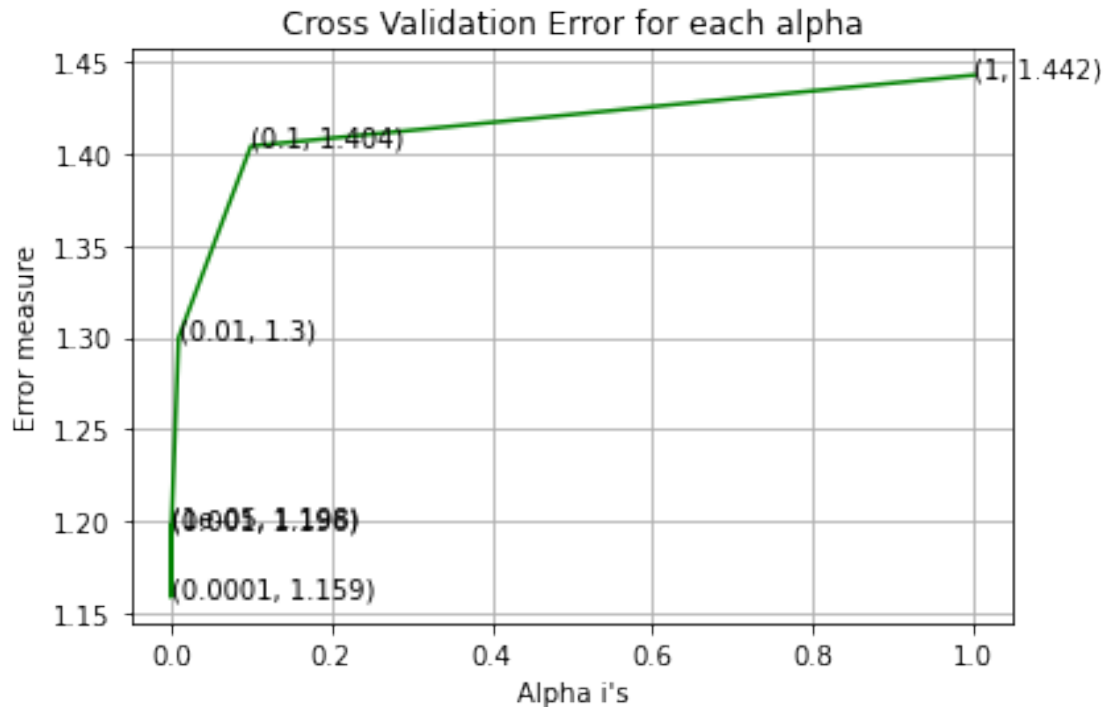
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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

```

For values of alpha = 1e-05 The log loss is: 1.1975800795058558  
 For values of alpha = 0.0001 The log loss is: 1.1590426179417324  
 For values of alpha = 0.001 The log loss is: 1.1959898640450148  
 For values of alpha = 0.01 The log loss is: 1.2999352042715278  
 For values of alpha = 0.1 The log loss is: 1.4038858401096033  
 For values of alpha = 1 The log loss is: 1.4423897723525352



For values of best alpha = 0.0001 The train log loss is: 1.0061824975433085  
 For values of best alpha = 0.0001 The cross validation log loss is:  
 1.1590426179417324  
 For values of best alpha = 0.0001 The test log loss is: 1.1500749249331368

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.

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

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 240 genes in train dataset?

Ans

1. In test data 653 out of 665 : 98.19548872180451
2. In cross validation data 514 out of 532 : 96.61654135338345

### 3.2.2 Univariate Analysis on Variation Feature

Q7. Variation, What type of feature is it ?

Ans. Variation is a categorical variable

Q8. How many categories are there?

```
[ ]: 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 : 1918

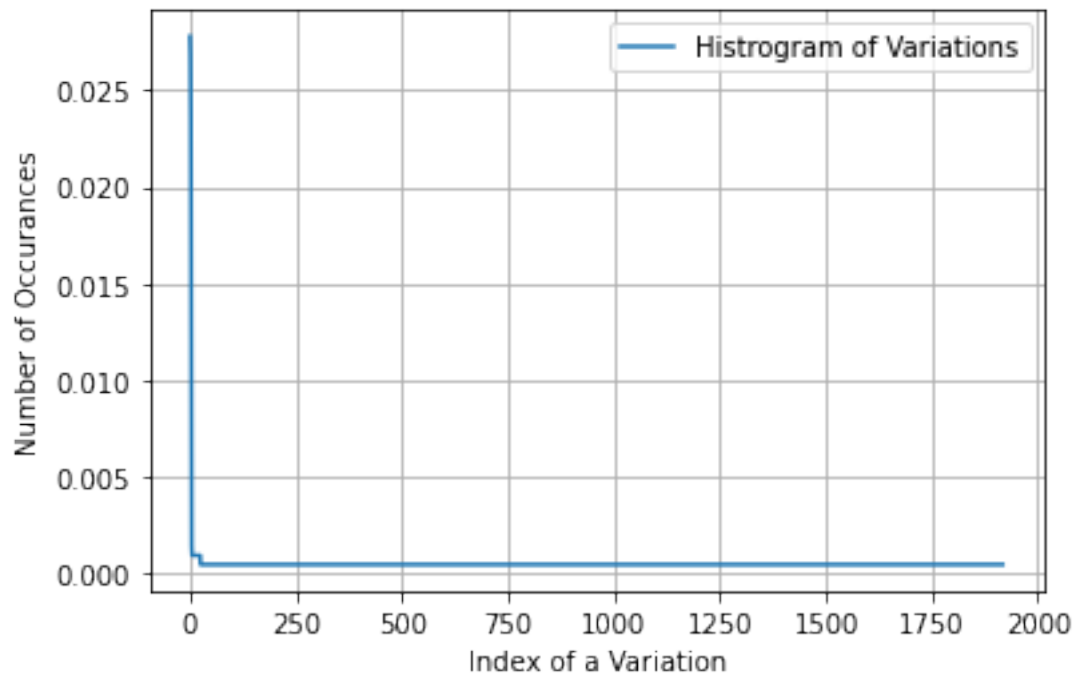
|                      |    |
|----------------------|----|
| Truncating_Mutations | 59 |
| Deletion             | 57 |
| Amplification        | 50 |
| Fusions              | 22 |
| Overexpression       | 3  |
| T58I                 | 2  |
| G67R                 | 2  |
| S308A                | 2  |
| E17K                 | 2  |
| E542K                | 2  |

Name: Variation, dtype: int64

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

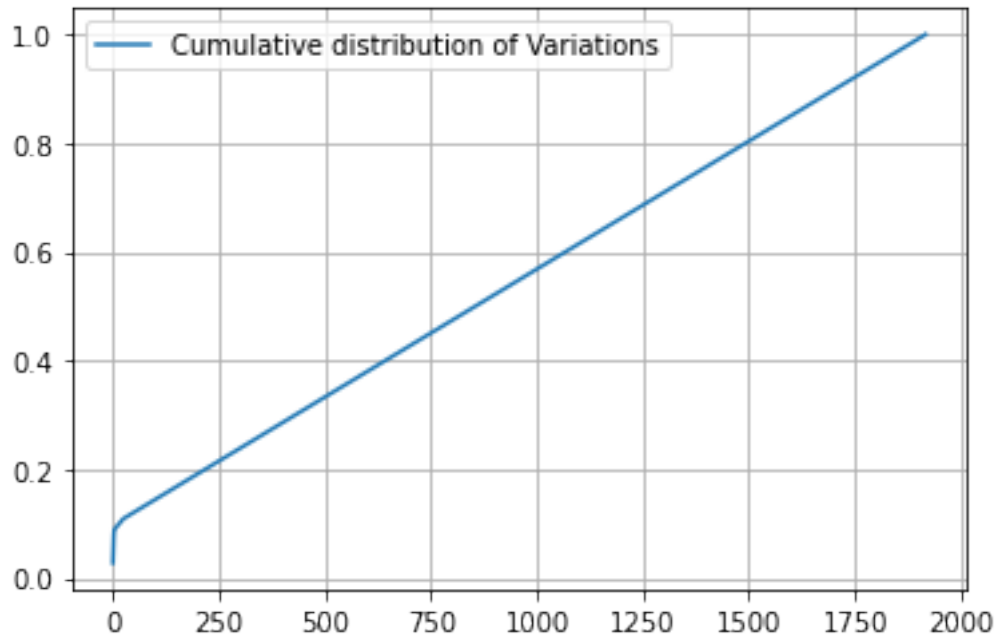
Ans: There are 1918 different categories of variations in the train data, and they are distributed as follows

```
[ ]: 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()
```



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

```
[0.02777778 0.05461394 0.07815443 ... 0.99905838 0.99952919 1.          ]
```



Q9. How to featurize this Variation feature ?

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

One hot Encoding

Response coding

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

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

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

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

```
[ ]: # one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer(ngram_range=(1,3))
train_variation_feature_onehotCoding = variation_vectorizer.
    ↳fit_transform(train_df['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.
    ↳transform(test_df['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.
    ↳transform(cv_df['Variation'])

[ ]: print("train_variation_feature_onehotEncoded is converted feature using the
    ↳onne-hot encoding method. The shape of Variation feature:",
    ↳train_variation_feature_onehotCoding.shape)
```

train\_variation\_feature\_onehotEncoded is converted feature using the onne-hot encoding method. The shape of Variation feature: (2124, 2053)

Q10. How good is this Variation feature in predicting y\_i?

Let's build a model just like the earlier!

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

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
    ↳generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
    ↳fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
    ↳learning_rate=optimal, eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

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

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

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



```

predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)

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

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

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

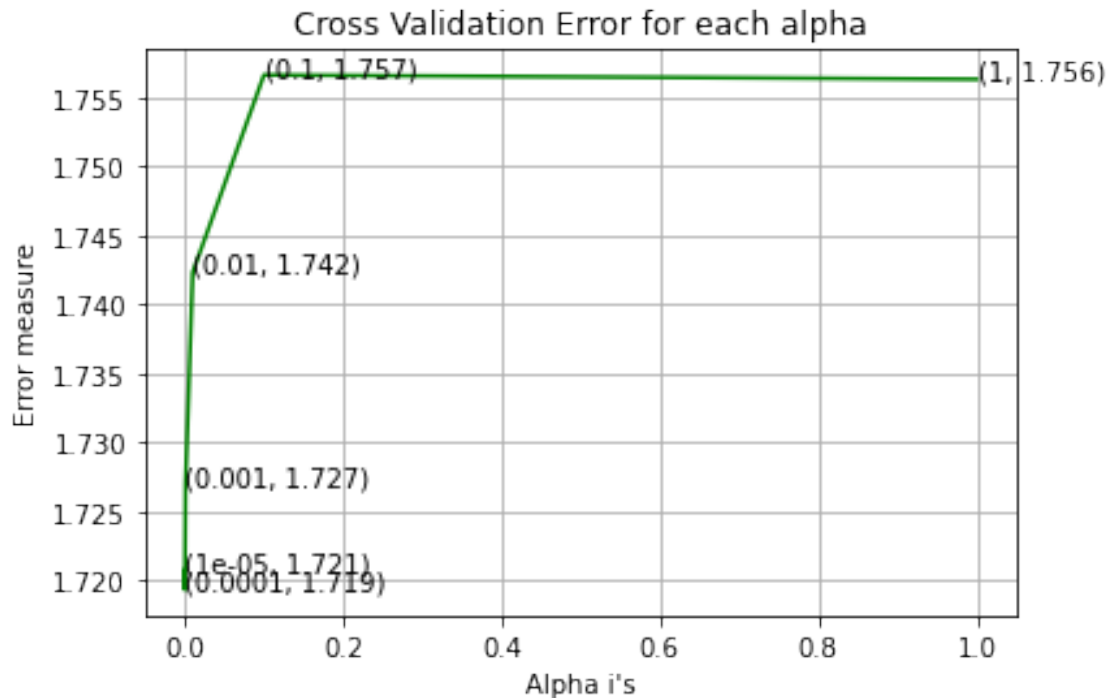
predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:
→",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
print('For values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

```

```

For values of alpha = 1e-05 The log loss is: 1.7207476487572677
For values of alpha = 0.0001 The log loss is: 1.7193230273243498
For values of alpha = 0.001 The log loss is: 1.726995530183674
For values of alpha = 0.01 The log loss is: 1.742324818179992
For values of alpha = 0.1 The log loss is: 1.7566376272022663
For values of alpha = 1 The log loss is: 1.7563276745470453

```



For values of best alpha = 0.0001 The train log loss is: 0.710578246017075  
 For values of best alpha = 0.0001 The cross validation log loss is:  
 1.7193230273243498  
 For values of best alpha = 0.0001 The test log loss is: 1.7099233127881008

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.

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

Q12. How many data points are covered by total 1918 genes in test and cross validation data sets?

Ans

1. In test data 57 out of 665 : 8.571428571428571
2. In cross validation data 58 out of 532 : 10.902255639097744

### 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?

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

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

```
[ ]: 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
```

```
[ ]: # building a CountVectorizer with all the words that occurred minimum 3 times in
→train data
text_vectorizer = CountVectorizer(min_df=3,ngram_range=(1,3))
train_text_feature_onehotCoding = text_vectorizer.
→fit_transform(train_df['TEXT'])
# getting all the feature names (words)
train_text_features= text_vectorizer.get_feature_names()

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

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

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

Total number of unique words in train data : 1861236

```
[ ]: 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)
```

```
[ ]: #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)
```

```
[ ]: # 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
```

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

```

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

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

```

```

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

```

```

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

```

```

Counter({3: 445398, 4: 280222, 5: 193407, 6: 173703, 7: 106059, 10: 84860, 8:
79705, 9: 76345, 12: 37095, 13: 36233, 14: 35178, 11: 29310, 16: 24627, 15:
21830, 20: 16947, 21: 15285, 17: 12165, 18: 11922, 24: 11196, 23: 8827, 41:
7981, 19: 7459, 25: 7433, 31: 7244, 22: 6353, 36: 6278, 45: 5986, 26: 5365, 28:
4766, 27: 4344, 66: 3968, 30: 3850, 29: 3359, 32: 3343, 58: 2775, 33: 2619, 42:
2395, 34: 2372, 37: 2347, 40: 2267, 35: 2267, 39: 1914, 38: 1901, 46: 1680, 48:
1567, 43: 1564, 44: 1509, 47: 1359, 50: 1330, 72: 1318, 49: 1113, 52: 1104, 51:
1098, 54: 992, 60: 972, 55: 960, 56: 891, 53: 889, 62: 862, 67: 855, 57: 821,
59: 767, 61: 751, 63: 714, 70: 709, 64: 686, 65: 659, 69: 644, 68: 636, 73: 625,
75: 549, 71: 541, 82: 536, 90: 527, 74: 508, 80: 490, 76: 478, 78: 472, 77: 452,
81: 428, 79: 426, 83: 421, 84: 412, 85: 377, 91: 344, 86: 334, 92: 329, 87: 328,
96: 319, 88: 318, 93: 305, 98: 302, 89: 301, 100: 286, 95: 279, 94: 270, 97:
264, 99: 252, 116: 249, 101: 247, 108: 242, 105: 241, 102: 236, 132: 234, 103:
234, 112: 226, 110: 225, 111: 224, 104: 221, 114: 208, 107: 199, 120: 198, 113:
196, 109: 189, 106: 185, 117: 180, 123: 179, 144: 175, 135: 172, 122: 170, 125:
169, 115: 167, 133: 163, 118: 163, 126: 162, 121: 161, 124: 156, 119: 156, 127:
152, 130: 148, 137: 146, 128: 144, 134: 140, 140: 138, 131: 138, 129: 138, 141:
133, 139: 133, 147: 130, 138: 124, 151: 123, 136: 123, 150: 122, 145: 120, 153:
119, 146: 119, 142: 116, 152: 114, 148: 113, 156: 112, 155: 109, 154: 103, 143:
99, 180: 98, 161: 98, 165: 96, 168: 95, 164: 95, 157: 91, 149: 91, 170: 90, 158:
90, 163: 88, 162: 88, 169: 87, 159: 87, 160: 85, 183: 83, 174: 83, 225: 80, 172:
80, 166: 80, 198: 77, 179: 75, 173: 75, 167: 74, 184: 73, 175: 73, 200: 70, 191:
70, 190: 69, 203: 67, 178: 67, 207: 66, 205: 66, 187: 65, 185: 65, 182: 65, 189:
64, 204: 63, 202: 63, 193: 63, 199: 62, 196: 62, 177: 62, 171: 62, 186: 61, 192:
60, 181: 60, 216: 59, 201: 59, 176: 58, 230: 57, 195: 56, 223: 51, 208: 51, 231:
50, 213: 50, 188: 50, 228: 49, 217: 49, 211: 49, 197: 49, 276: 48, 218: 48, 246:
47, 212: 47, 194: 47, 250: 46, 234: 46, 209: 46, 256: 45, 210: 45, 232: 44, 226:
44, 206: 44, 265: 43, 244: 43, 242: 43, 229: 43, 241: 42, 238: 42, 270: 41, 248:

```

41, 263: 40, 233: 40, 221: 40, 220: 40, 215: 40, 222: 39, 269: 38, 237: 38, 258:  
 37, 257: 37, 254: 37, 236: 37, 224: 37, 304: 36, 264: 36, 247: 36, 245: 36, 243:  
 36, 298: 35, 278: 35, 252: 35, 251: 35, 235: 35, 227: 35, 214: 35, 300: 34, 288:  
 34, 280: 34, 261: 34, 255: 34, 219: 34, 274: 33, 260: 33, 259: 33, 249: 33, 239:  
 33, 287: 32, 262: 32, 240: 32, 279: 31, 324: 30, 314: 30, 294: 30, 277: 30, 271:  
 30, 268: 30, 303: 29, 297: 29, 289: 29, 286: 29, 330: 28, 317: 28, 290: 28, 285:  
 28, 284: 28, 283: 28, 282: 28, 275: 28, 319: 27, 315: 27, 312: 27, 310: 27, 272:  
 27, 364: 26, 331: 26, 328: 26, 305: 26, 302: 26, 301: 26, 273: 26, 267: 26, 336:  
 25, 335: 25, 306: 25, 295: 25, 292: 25, 266: 25, 253: 25, 381: 24, 365: 24, 337:  
 24, 329: 24, 325: 24, 309: 24, 308: 24, 296: 24, 353: 23, 343: 23, 293: 23, 327:  
 22, 391: 21, 374: 21, 362: 21, 355: 21, 347: 21, 340: 21, 334: 21, 299: 21, 450:  
 20, 371: 20, 369: 20, 358: 20, 341: 20, 320: 20, 281: 20, 378: 19, 360: 19, 352:  
 19, 348: 19, 333: 19, 326: 19, 316: 19, 313: 19, 291: 19, 407: 18, 361: 18, 357:  
 18, 344: 18, 342: 18, 323: 18, 321: 18, 311: 18, 307: 18, 1019: 17, 460: 17,  
 436: 17, 423: 17, 406: 17, 385: 17, 376: 17, 368: 17, 350: 17, 688: 16, 551: 16,  
 451: 16, 430: 16, 424: 16, 415: 16, 409: 16, 405: 16, 390: 16, 384: 16, 383: 16,  
 370: 16, 359: 16, 339: 16, 332: 16, 474: 15, 470: 15, 417: 15, 413: 15, 399: 15,  
 394: 15, 382: 15, 380: 15, 354: 15, 345: 15, 585: 14, 464: 14, 462: 14, 455: 14,  
 448: 14, 445: 14, 440: 14, 438: 14, 418: 14, 410: 14, 396: 14, 387: 14, 372: 14,  
 363: 14, 346: 14, 318: 14, 517: 13, 510: 13, 509: 13, 485: 13, 468: 13, 452: 13,  
 416: 13, 411: 13, 397: 13, 395: 13, 393: 13, 392: 13, 377: 13, 349: 13, 338: 13,  
 720: 12, 557: 12, 539: 12, 534: 12, 511: 12, 498: 12, 496: 12, 495: 12, 459: 12,  
 454: 12, 453: 12, 439: 12, 437: 12, 435: 12, 432: 12, 421: 12, 419: 12, 408: 12,  
 403: 12, 402: 12, 401: 12, 389: 12, 388: 12, 375: 12, 356: 12, 322: 12, 626: 11,  
 589: 11, 576: 11, 574: 11, 549: 11, 547: 11, 501: 11, 473: 11, 461: 11, 456: 11,  
 446: 11, 444: 11, 442: 11, 431: 11, 429: 11, 425: 11, 420: 11, 414: 11, 404: 11,  
 400: 11, 373: 11, 367: 11, 366: 11, 713: 10, 702: 10, 663: 10, 634: 10, 616: 10,  
 613: 10, 575: 10, 564: 10, 563: 10, 560: 10, 558: 10, 550: 10, 544: 10, 532: 10,  
 527: 10, 483: 10, 472: 10, 463: 10, 441: 10, 433: 10, 428: 10, 398: 10, 386: 10,  
 351: 10, 841: 9, 763: 9, 735: 9, 695: 9, 643: 9, 625: 9, 624: 9, 604: 9, 598: 9,  
 553: 9, 543: 9, 540: 9, 537: 9, 536: 9, 535: 9, 533: 9, 515: 9, 503: 9, 499: 9,  
 497: 9, 484: 9, 481: 9, 480: 9, 476: 9, 475: 9, 469: 9, 443: 9, 427: 9, 422: 9,  
 1137: 8, 929: 8, 809: 8, 716: 8, 673: 8, 590: 8, 587: 8, 586: 8, 583: 8, 582: 8,  
 573: 8, 565: 8, 556: 8, 555: 8, 552: 8, 531: 8, 529: 8, 528: 8, 523: 8, 522: 8,  
 508: 8, 506: 8, 502: 8, 500: 8, 494: 8, 492: 8, 489: 8, 487: 8, 478: 8, 467: 8,  
 465: 8, 426: 8, 412: 8, 379: 8, 1163: 7, 1082: 7, 1028: 7, 986: 7, 953: 7, 907:  
 7, 902: 7, 854: 7, 816: 7, 795: 7, 786: 7, 779: 7, 770: 7, 767: 7, 766: 7, 755:  
 7, 753: 7, 715: 7, 701: 7, 665: 7, 654: 7, 646: 7, 621: 7, 620: 7, 606: 7, 602:  
 7, 588: 7, 580: 7, 578: 7, 572: 7, 569: 7, 562: 7, 545: 7, 541: 7, 530: 7, 520:  
 7, 519: 7, 514: 7, 513: 7, 512: 7, 507: 7, 504: 7, 493: 7, 490: 7, 488: 7, 477:  
 7, 458: 7, 457: 7, 449: 7, 434: 7, 1211: 6, 1159: 6, 1088: 6, 985: 6, 981: 6,  
 931: 6, 921: 6, 920: 6, 909: 6, 830: 6, 826: 6, 815: 6, 799: 6, 751: 6, 747: 6,  
 744: 6, 731: 6, 729: 6, 728: 6, 727: 6, 719: 6, 718: 6, 711: 6, 710: 6, 709: 6,  
 708: 6, 691: 6, 685: 6, 684: 6, 677: 6, 671: 6, 667: 6, 657: 6, 642: 6, 638: 6,  
 631: 6, 628: 6, 617: 6, 615: 6, 614: 6, 612: 6, 609: 6, 607: 6, 600: 6, 571: 6,  
 570: 6, 566: 6, 546: 6, 542: 6, 526: 6, 525: 6, 521: 6, 518: 6, 479: 6, 1606: 5,  
 1573: 5, 1474: 5, 1410: 5, 1383: 5, 1302: 5, 1269: 5, 1250: 5, 1192: 5, 1089: 5,  
 1080: 5, 1078: 5, 1010: 5, 992: 5, 980: 5, 976: 5, 972: 5, 966: 5, 961: 5, 960:

5, 949: 5, 948: 5, 936: 5, 932: 5, 901: 5, 899: 5, 887: 5, 869: 5, 865: 5, 862:  
 5, 859: 5, 853: 5, 850: 5, 846: 5, 824: 5, 822: 5, 811: 5, 807: 5, 803: 5, 798:  
 5, 781: 5, 778: 5, 771: 5, 768: 5, 764: 5, 759: 5, 756: 5, 754: 5, 750: 5, 746:  
 5, 740: 5, 726: 5, 723: 5, 721: 5, 696: 5, 693: 5, 686: 5, 675: 5, 669: 5, 668:  
 5, 662: 5, 656: 5, 655: 5, 648: 5, 647: 5, 645: 5, 644: 5, 641: 5, 640: 5, 639:  
 5, 636: 5, 632: 5, 630: 5, 619: 5, 610: 5, 608: 5, 597: 5, 593: 5, 591: 5, 579:  
 5, 577: 5, 568: 5, 554: 5, 538: 5, 516: 5, 505: 5, 486: 5, 447: 5, 2457: 4,  
 2211: 4, 2208: 4, 2056: 4, 1784: 4, 1680: 4, 1601: 4, 1585: 4, 1531: 4, 1430: 4,  
 1401: 4, 1348: 4, 1319: 4, 1316: 4, 1308: 4, 1299: 4, 1268: 4, 1265: 4, 1256: 4,  
 1244: 4, 1243: 4, 1239: 4, 1231: 4, 1222: 4, 1213: 4, 1210: 4, 1174: 4, 1166: 4,  
 1161: 4, 1123: 4, 1115: 4, 1108: 4, 1102: 4, 1074: 4, 1052: 4, 1051: 4, 1045: 4,  
 1044: 4, 1042: 4, 1032: 4, 1025: 4, 1022: 4, 1000: 4, 998: 4, 995: 4, 967: 4,  
 954: 4, 943: 4, 930: 4, 924: 4, 914: 4, 913: 4, 906: 4, 903: 4, 896: 4, 889: 4,  
 886: 4, 885: 4, 882: 4, 880: 4, 871: 4, 858: 4, 848: 4, 842: 4, 840: 4, 839: 4,  
 838: 4, 837: 4, 834: 4, 833: 4, 831: 4, 829: 4, 828: 4, 827: 4, 819: 4, 813: 4,  
 801: 4, 797: 4, 796: 4, 792: 4, 790: 4, 785: 4, 780: 4, 777: 4, 773: 4, 772: 4,  
 769: 4, 765: 4, 762: 4, 758: 4, 745: 4, 742: 4, 738: 4, 734: 4, 733: 4, 712: 4,  
 704: 4, 703: 4, 700: 4, 699: 4, 694: 4, 689: 4, 687: 4, 683: 4, 682: 4, 680: 4,  
 676: 4, 674: 4, 672: 4, 666: 4, 660: 4, 659: 4, 658: 4, 653: 4, 651: 4, 650: 4,  
 635: 4, 633: 4, 629: 4, 627: 4, 623: 4, 611: 4, 605: 4, 603: 4, 601: 4, 599: 4,  
 595: 4, 594: 4, 584: 4, 581: 4, 567: 4, 524: 4, 482: 4, 5160: 3, 3217: 3, 2996:  
 3, 2791: 3, 2605: 3, 2522: 3, 2505: 3, 2446: 3, 2380: 3, 2365: 3, 2364: 3, 2248:  
 3, 2149: 3, 2053: 3, 2028: 3, 2005: 3, 1940: 3, 1920: 3, 1914: 3, 1909: 3, 1906:  
 3, 1901: 3, 1865: 3, 1846: 3, 1812: 3, 1786: 3, 1781: 3, 1754: 3, 1742: 3, 1734:  
 3, 1729: 3, 1710: 3, 1702: 3, 1694: 3, 1630: 3, 1627: 3, 1579: 3, 1550: 3, 1525:  
 3, 1495: 3, 1494: 3, 1490: 3, 1486: 3, 1452: 3, 1441: 3, 1421: 3, 1416: 3, 1413:  
 3, 1400: 3, 1394: 3, 1367: 3, 1365: 3, 1364: 3, 1357: 3, 1352: 3, 1349: 3, 1335:  
 3, 1332: 3, 1330: 3, 1325: 3, 1324: 3, 1292: 3, 1285: 3, 1272: 3, 1271: 3, 1259:  
 3, 1248: 3, 1241: 3, 1236: 3, 1226: 3, 1225: 3, 1212: 3, 1207: 3, 1203: 3, 1187:  
 3, 1186: 3, 1184: 3, 1182: 3, 1181: 3, 1175: 3, 1172: 3, 1167: 3, 1157: 3, 1155:  
 3, 1153: 3, 1150: 3, 1139: 3, 1131: 3, 1127: 3, 1125: 3, 1124: 3, 1117: 3, 1112:  
 3, 1107: 3, 1106: 3, 1103: 3, 1100: 3, 1092: 3, 1087: 3, 1084: 3, 1081: 3, 1071:  
 3, 1066: 3, 1064: 3, 1060: 3, 1058: 3, 1043: 3, 1029: 3, 1027: 3, 1014: 3, 1009:  
 3, 1004: 3, 984: 3, 982: 3, 978: 3, 975: 3, 973: 3, 971: 3, 962: 3, 957: 3, 951:  
 3, 947: 3, 942: 3, 938: 3, 937: 3, 935: 3, 934: 3, 933: 3, 928: 3, 923: 3, 922:  
 3, 915: 3, 910: 3, 894: 3, 879: 3, 878: 3, 874: 3, 867: 3, 861: 3, 857: 3, 856:  
 3, 852: 3, 851: 3, 844: 3, 843: 3, 836: 3, 835: 3, 823: 3, 821: 3, 814: 3, 806:  
 3, 804: 3, 789: 3, 788: 3, 787: 3, 774: 3, 760: 3, 757: 3, 749: 3, 741: 3, 739:  
 3, 737: 3, 725: 3, 722: 3, 717: 3, 706: 3, 697: 3, 692: 3, 681: 3, 678: 3, 670:  
 3, 652: 3, 649: 3, 622: 3, 618: 3, 561: 3, 559: 3, 548: 3, 491: 3, 471: 3, 466:  
 3, 12233: 2, 6429: 2, 6401: 2, 5919: 2, 5716: 2, 5595: 2, 5217: 2, 4967: 2,  
 4936: 2, 4821: 2, 4784: 2, 4772: 2, 4702: 2, 4479: 2, 4371: 2, 4279: 2, 4243: 2,  
 4097: 2, 4095: 2, 4038: 2, 3982: 2, 3971: 2, 3848: 2, 3769: 2, 3637: 2, 3592: 2,  
 3576: 2, 3575: 2, 3552: 2, 3519: 2, 3449: 2, 3404: 2, 3400: 2, 3348: 2, 3338: 2,  
 3265: 2, 3165: 2, 3162: 2, 3158: 2, 3143: 2, 3125: 2, 3122: 2, 3113: 2, 3100: 2,  
 3094: 2, 3081: 2, 3080: 2, 3044: 2, 2933: 2, 2910: 2, 2898: 2, 2893: 2, 2889: 2,  
 2870: 2, 2855: 2, 2850: 2, 2810: 2, 2765: 2, 2739: 2, 2657: 2, 2656: 2, 2645: 2,  
 2639: 2, 2632: 2, 2627: 2, 2590: 2, 2580: 2, 2576: 2, 2557: 2, 2553: 2, 2524: 2,

2503: 2, 2499: 2, 2487: 2, 2482: 2, 2451: 2, 2449: 2, 2447: 2, 2445: 2, 2434: 2,  
 2414: 2, 2408: 2, 2395: 2, 2389: 2, 2388: 2, 2369: 2, 2340: 2, 2336: 2, 2332: 2,  
 2329: 2, 2289: 2, 2271: 2, 2249: 2, 2237: 2, 2230: 2, 2196: 2, 2192: 2, 2183: 2,  
 2177: 2, 2162: 2, 2145: 2, 2134: 2, 2130: 2, 2127: 2, 2119: 2, 2118: 2, 2116: 2,  
 2104: 2, 2096: 2, 2091: 2, 2086: 2, 2079: 2, 2070: 2, 2059: 2, 2055: 2, 2052: 2,  
 2051: 2, 2042: 2, 2037: 2, 2033: 2, 2030: 2, 2019: 2, 2007: 2, 2004: 2, 1992: 2,  
 1978: 2, 1976: 2, 1971: 2, 1959: 2, 1952: 2, 1934: 2, 1932: 2, 1912: 2, 1910: 2,  
 1907: 2, 1903: 2, 1888: 2, 1885: 2, 1878: 2, 1870: 2, 1866: 2, 1862: 2, 1859: 2,  
 1851: 2, 1844: 2, 1837: 2, 1834: 2, 1828: 2, 1814: 2, 1813: 2, 1808: 2, 1799: 2,  
 1794: 2, 1792: 2, 1785: 2, 1782: 2, 1773: 2, 1772: 2, 1770: 2, 1766: 2, 1764: 2,  
 1761: 2, 1760: 2, 1743: 2, 1741: 2, 1739: 2, 1736: 2, 1725: 2, 1720: 2, 1717: 2,  
 1716: 2, 1712: 2, 1705: 2, 1704: 2, 1699: 2, 1698: 2, 1693: 2, 1691: 2, 1686: 2,  
 1681: 2, 1676: 2, 1675: 2, 1668: 2, 1667: 2, 1664: 2, 1661: 2, 1653: 2, 1651: 2,  
 1649: 2, 1645: 2, 1638: 2, 1637: 2, 1629: 2, 1622: 2, 1609: 2, 1608: 2, 1600: 2,  
 1599: 2, 1596: 2, 1595: 2, 1594: 2, 1588: 2, 1569: 2, 1568: 2, 1567: 2, 1563: 2,  
 1554: 2, 1552: 2, 1548: 2, 1547: 2, 1543: 2, 1538: 2, 1537: 2, 1534: 2, 1533: 2,  
 1529: 2, 1527: 2, 1522: 2, 1517: 2, 1514: 2, 1506: 2, 1505: 2, 1499: 2, 1498: 2,  
 1493: 2, 1489: 2, 1482: 2, 1481: 2, 1480: 2, 1477: 2, 1471: 2, 1467: 2, 1465: 2,  
 1461: 2, 1456: 2, 1444: 2, 1440: 2, 1437: 2, 1435: 2, 1424: 2, 1423: 2, 1422: 2,  
 1417: 2, 1402: 2, 1393: 2, 1391: 2, 1389: 2, 1386: 2, 1385: 2, 1380: 2, 1379: 2,  
 1376: 2, 1370: 2, 1369: 2, 1368: 2, 1353: 2, 1346: 2, 1343: 2, 1339: 2, 1337: 2,  
 1333: 2, 1328: 2, 1321: 2, 1313: 2, 1310: 2, 1309: 2, 1304: 2, 1303: 2, 1298: 2,  
 1296: 2, 1294: 2, 1288: 2, 1287: 2, 1281: 2, 1277: 2, 1276: 2, 1275: 2, 1274: 2,  
 1262: 2, 1260: 2, 1255: 2, 1254: 2, 1247: 2, 1242: 2, 1230: 2, 1229: 2, 1224: 2,  
 1221: 2, 1219: 2, 1216: 2, 1209: 2, 1200: 2, 1199: 2, 1197: 2, 1196: 2, 1193: 2,  
 1185: 2, 1179: 2, 1176: 2, 1168: 2, 1158: 2, 1156: 2, 1154: 2, 1152: 2, 1151: 2,  
 1147: 2, 1146: 2, 1144: 2, 1138: 2, 1133: 2, 1120: 2, 1116: 2, 1114: 2, 1113: 2,  
 1094: 2, 1091: 2, 1079: 2, 1077: 2, 1072: 2, 1057: 2, 1056: 2, 1055: 2, 1054: 2,  
 1050: 2, 1048: 2, 1047: 2, 1039: 2, 1037: 2, 1036: 2, 1035: 2, 1023: 2, 1021: 2,  
 1020: 2, 1008: 2, 1005: 2, 1002: 2, 997: 2, 994: 2, 991: 2, 988: 2, 974: 2, 970:  
 2, 969: 2, 964: 2, 958: 2, 956: 2, 955: 2, 945: 2, 940: 2, 939: 2, 926: 2, 925:  
 2, 918: 2, 917: 2, 916: 2, 911: 2, 905: 2, 898: 2, 892: 2, 891: 2, 888: 2, 884:  
 2, 883: 2, 876: 2, 873: 2, 870: 2, 863: 2, 847: 2, 845: 2, 832: 2, 825: 2, 818:  
 2, 817: 2, 810: 2, 808: 2, 805: 2, 802: 2, 800: 2, 794: 2, 793: 2, 791: 2, 784:  
 2, 783: 2, 775: 2, 761: 2, 752: 2, 743: 2, 730: 2, 724: 2, 714: 2, 707: 2, 679:  
 2, 664: 2, 637: 2, 592: 2, 152286: 1, 116758: 1, 80687: 1, 68299: 1, 67007: 1,  
 65890: 1, 65751: 1, 65181: 1, 63579: 1, 62634: 1, 56389: 1, 53998: 1, 48504: 1,  
 48255: 1, 46138: 1, 45818: 1, 43770: 1, 42838: 1, 42492: 1, 41678: 1, 41399: 1,  
 40855: 1, 40732: 1, 39299: 1, 39251: 1, 38967: 1, 38316: 1, 38278: 1, 36371: 1,  
 36103: 1, 35945: 1, 35570: 1, 34410: 1, 33569: 1, 33311: 1, 32927: 1, 31211: 1,  
 30800: 1, 29269: 1, 27723: 1, 25924: 1, 25811: 1, 25715: 1, 25680: 1, 25453: 1,  
 25424: 1, 25382: 1, 24829: 1, 24445: 1, 24053: 1, 23989: 1, 23943: 1, 23901: 1,  
 23898: 1, 22931: 1, 22442: 1, 22139: 1, 22031: 1, 21752: 1, 21193: 1, 21030: 1,  
 20805: 1, 20734: 1, 20184: 1, 20082: 1, 19856: 1, 19624: 1, 19576: 1, 19385: 1,  
 19093: 1, 19075: 1, 19070: 1, 19008: 1, 18912: 1, 18764: 1, 18644: 1, 18515: 1,  
 18370: 1, 18326: 1, 18114: 1, 18059: 1, 18056: 1, 18025: 1, 17770: 1, 17728: 1,  
 17657: 1, 17655: 1, 17515: 1, 17342: 1, 17292: 1, 17284: 1, 17135: 1, 16963: 1,  
 16921: 1, 16865: 1, 16846: 1, 16760: 1, 16622: 1, 16573: 1, 16497: 1, 16176: 1,



16073: 1, 15965: 1, 15877: 1, 15874: 1, 15616: 1, 15412: 1, 15404: 1, 15321: 1,  
15310: 1, 15296: 1, 15203: 1, 15200: 1, 15048: 1, 14955: 1, 14952: 1, 14737: 1,  
14631: 1, 14546: 1, 14507: 1, 14417: 1, 14340: 1, 14322: 1, 14262: 1, 14208: 1,  
14147: 1, 14101: 1, 13808: 1, 13675: 1, 13633: 1, 13572: 1, 13350: 1, 13338: 1,  
13254: 1, 13119: 1, 13089: 1, 12950: 1, 12902: 1, 12899: 1, 12878: 1, 12801: 1,  
12753: 1, 12642: 1, 12580: 1, 12562: 1, 12523: 1, 12512: 1, 12472: 1, 12469: 1,  
12391: 1, 12310: 1, 12300: 1, 12259: 1, 12250: 1, 12235: 1, 12229: 1, 12211: 1,  
12129: 1, 12057: 1, 12010: 1, 11967: 1, 11934: 1, 11922: 1, 11890: 1, 11775: 1,  
11731: 1, 11700: 1, 11693: 1, 11691: 1, 11683: 1, 11648: 1, 11598: 1, 11529: 1,  
11491: 1, 11438: 1, 11423: 1, 11357: 1, 11335: 1, 11280: 1, 11266: 1, 11187: 1,  
11141: 1, 11061: 1, 11059: 1, 11039: 1, 11025: 1, 10922: 1, 10859: 1, 10814: 1,  
10583: 1, 10570: 1, 10470: 1, 10432: 1, 10390: 1, 10359: 1, 10325: 1, 10318: 1,  
10251: 1, 10199: 1, 10193: 1, 10159: 1, 10156: 1, 10120: 1, 10047: 1, 10042: 1,  
9996: 1, 9898: 1, 9894: 1, 9838: 1, 9837: 1, 9810: 1, 9796: 1, 9645: 1, 9619: 1,  
9589: 1, 9530: 1, 9517: 1, 9506: 1, 9496: 1, 9408: 1, 9405: 1, 9389: 1, 9364: 1,  
9279: 1, 9270: 1, 9261: 1, 9255: 1, 9215: 1, 9199: 1, 9177: 1, 9125: 1, 9112: 1,  
9082: 1, 9079: 1, 9052: 1, 9021: 1, 9003: 1, 8999: 1, 8932: 1, 8926: 1, 8921: 1,  
8853: 1, 8799: 1, 8792: 1, 8779: 1, 8773: 1, 8720: 1, 8678: 1, 8623: 1, 8591: 1,  
8513: 1, 8511: 1, 8431: 1, 8425: 1, 8414: 1, 8362: 1, 8323: 1, 8276: 1, 8262: 1,  
8251: 1, 8250: 1, 8248: 1, 8226: 1, 8225: 1, 8214: 1, 8200: 1, 8180: 1, 8178: 1,  
8175: 1, 8130: 1, 8127: 1, 8102: 1, 8098: 1, 8084: 1, 8063: 1, 8045: 1, 8042: 1,  
8038: 1, 8026: 1, 8025: 1, 8024: 1, 7967: 1, 7960: 1, 7928: 1, 7911: 1, 7889: 1,  
7885: 1, 7837: 1, 7822: 1, 7821: 1, 7796: 1, 7789: 1, 7756: 1, 7747: 1, 7745: 1,  
7738: 1, 7732: 1, 7725: 1, 7716: 1, 7714: 1, 7672: 1, 7659: 1, 7581: 1, 7490: 1,  
7489: 1, 7469: 1, 7466: 1, 7432: 1, 7318: 1, 7317: 1, 7302: 1, 7278: 1, 7268: 1,  
7231: 1, 7224: 1, 7215: 1, 7212: 1, 7209: 1, 7199: 1, 7176: 1, 7161: 1, 7153: 1,  
7129: 1, 7122: 1, 7105: 1, 7104: 1, 7052: 1, 7042: 1, 7039: 1, 7034: 1, 7011: 1,  
7009: 1, 6989: 1, 6961: 1, 6954: 1, 6932: 1, 6931: 1, 6906: 1, 6899: 1, 6886: 1,  
6880: 1, 6876: 1, 6837: 1, 6835: 1, 6834: 1, 6827: 1, 6810: 1, 6803: 1, 6795: 1,  
6733: 1, 6717: 1, 6704: 1, 6698: 1, 6692: 1, 6690: 1, 6665: 1, 6641: 1, 6631: 1,  
6627: 1, 6624: 1, 6607: 1, 6598: 1, 6595: 1, 6592: 1, 6588: 1, 6582: 1, 6558: 1,  
6543: 1, 6531: 1, 6525: 1, 6523: 1, 6517: 1, 6516: 1, 6508: 1, 6501: 1, 6486: 1,  
6472: 1, 6459: 1, 6448: 1, 6436: 1, 6421: 1, 6410: 1, 6378: 1, 6366: 1, 6318: 1,  
6278: 1, 6277: 1, 6255: 1, 6254: 1, 6224: 1, 6218: 1, 6213: 1, 6145: 1, 6085: 1,  
6082: 1, 6074: 1, 6073: 1, 6061: 1, 6048: 1, 6045: 1, 6044: 1, 6042: 1, 6033: 1,  
6018: 1, 5997: 1, 5976: 1, 5974: 1, 5960: 1, 5952: 1, 5948: 1, 5942: 1, 5920: 1,  
5886: 1, 5871: 1, 5850: 1, 5846: 1, 5819: 1, 5800: 1, 5797: 1, 5787: 1, 5761: 1,  
5760: 1, 5757: 1, 5742: 1, 5741: 1, 5729: 1, 5727: 1, 5699: 1, 5691: 1, 5676: 1,  
5667: 1, 5659: 1, 5641: 1, 5630: 1, 5626: 1, 5610: 1, 5594: 1, 5584: 1, 5574: 1,  
5547: 1, 5520: 1, 5516: 1, 5503: 1, 5486: 1, 5469: 1, 5460: 1, 5418: 1, 5400: 1,  
5379: 1, 5372: 1, 5364: 1, 5352: 1, 5351: 1, 5345: 1, 5341: 1, 5340: 1, 5333: 1,  
5331: 1, 5304: 1, 5301: 1, 5278: 1, 5277: 1, 5257: 1, 5248: 1, 5238: 1, 5227: 1,  
5212: 1, 5208: 1, 5180: 1, 5159: 1, 5157: 1, 5156: 1, 5142: 1, 5127: 1, 5125: 1,  
5119: 1, 5100: 1, 5081: 1, 5077: 1, 5073: 1, 5069: 1, 5050: 1, 5046: 1, 5043: 1,  
5035: 1, 5029: 1, 5028: 1, 5014: 1, 5012: 1, 5008: 1, 5001: 1, 4996: 1, 4971: 1,  
4966: 1, 4932: 1, 4929: 1, 4912: 1, 4909: 1, 4877: 1, 4869: 1, 4863: 1, 4859: 1,  
4841: 1, 4827: 1, 4815: 1, 4812: 1, 4808: 1, 4807: 1, 4805: 1, 4800: 1, 4795: 1,  
4788: 1, 4775: 1, 4769: 1, 4768: 1, 4743: 1, 4733: 1, 4732: 1, 4726: 1, 4694: 1,

4690: 1, 4679: 1, 4676: 1, 4669: 1, 4658: 1, 4635: 1, 4627: 1, 4622: 1, 4591: 1,  
4568: 1, 4553: 1, 4541: 1, 4533: 1, 4532: 1, 4527: 1, 4524: 1, 4520: 1, 4519: 1,  
4498: 1, 4495: 1, 4481: 1, 4477: 1, 4471: 1, 4466: 1, 4460: 1, 4458: 1, 4454: 1,  
4453: 1, 4445: 1, 4444: 1, 4436: 1, 4427: 1, 4424: 1, 4405: 1, 4400: 1, 4398: 1,  
4390: 1, 4388: 1, 4386: 1, 4383: 1, 4353: 1, 4338: 1, 4335: 1, 4334: 1, 4320: 1,  
4315: 1, 4312: 1, 4301: 1, 4299: 1, 4286: 1, 4280: 1, 4278: 1, 4277: 1, 4267: 1,  
4257: 1, 4252: 1, 4251: 1, 4244: 1, 4218: 1, 4214: 1, 4192: 1, 4181: 1, 4179: 1,  
4173: 1, 4158: 1, 4156: 1, 4150: 1, 4148: 1, 4145: 1, 4143: 1, 4141: 1, 4122: 1,  
4121: 1, 4120: 1, 4117: 1, 4116: 1, 4111: 1, 4109: 1, 4098: 1, 4091: 1, 4089: 1,  
4085: 1, 4071: 1, 4063: 1, 4048: 1, 4036: 1, 4034: 1, 4026: 1, 4022: 1, 4020: 1,  
4019: 1, 4013: 1, 4009: 1, 4008: 1, 3999: 1, 3997: 1, 3991: 1, 3986: 1, 3984: 1,  
3973: 1, 3968: 1, 3967: 1, 3951: 1, 3939: 1, 3932: 1, 3928: 1, 3915: 1, 3913: 1,  
3911: 1, 3902: 1, 3901: 1, 3897: 1, 3893: 1, 3884: 1, 3875: 1, 3870: 1, 3869: 1,  
3868: 1, 3866: 1, 3864: 1, 3853: 1, 3846: 1, 3842: 1, 3839: 1, 3838: 1, 3834: 1,  
3826: 1, 3823: 1, 3809: 1, 3805: 1, 3800: 1, 3794: 1, 3789: 1, 3782: 1, 3781: 1,  
3772: 1, 3764: 1, 3748: 1, 3746: 1, 3744: 1, 3735: 1, 3720: 1, 3715: 1, 3713: 1,  
3709: 1, 3708: 1, 3707: 1, 3704: 1, 3697: 1, 3696: 1, 3688: 1, 3683: 1, 3677: 1,  
3675: 1, 3670: 1, 3657: 1, 3656: 1, 3653: 1, 3644: 1, 3640: 1, 3636: 1, 3634: 1,  
3631: 1, 3627: 1, 3624: 1, 3617: 1, 3611: 1, 3608: 1, 3601: 1, 3600: 1, 3593: 1,  
3590: 1, 3572: 1, 3569: 1, 3568: 1, 3566: 1, 3553: 1, 3547: 1, 3545: 1, 3538: 1,  
3536: 1, 3535: 1, 3533: 1, 3529: 1, 3528: 1, 3523: 1, 3512: 1, 3510: 1, 3509: 1,  
3508: 1, 3498: 1, 3493: 1, 3492: 1, 3484: 1, 3482: 1, 3474: 1, 3473: 1, 3472: 1,  
3470: 1, 3469: 1, 3468: 1, 3460: 1, 3459: 1, 3457: 1, 3456: 1, 3455: 1, 3443: 1,  
3439: 1, 3433: 1, 3429: 1, 3426: 1, 3425: 1, 3423: 1, 3422: 1, 3415: 1, 3412: 1,  
3407: 1, 3399: 1, 3390: 1, 3383: 1, 3369: 1, 3368: 1, 3366: 1, 3357: 1, 3356: 1,  
3350: 1, 3347: 1, 3344: 1, 3336: 1, 3333: 1, 3331: 1, 3325: 1, 3320: 1, 3319: 1,  
3316: 1, 3314: 1, 3313: 1, 3309: 1, 3302: 1, 3297: 1, 3294: 1, 3293: 1, 3291: 1,  
3289: 1, 3278: 1, 3269: 1, 3268: 1, 3267: 1, 3264: 1, 3259: 1, 3253: 1, 3250: 1,  
3248: 1, 3244: 1, 3243: 1, 3237: 1, 3234: 1, 3233: 1, 3229: 1, 3228: 1, 3223: 1,  
3222: 1, 3213: 1, 3207: 1, 3205: 1, 3203: 1, 3195: 1, 3193: 1, 3188: 1, 3186: 1,  
3183: 1, 3180: 1, 3177: 1, 3175: 1, 3170: 1, 3169: 1, 3167: 1, 3166: 1, 3164: 1,  
3159: 1, 3156: 1, 3154: 1, 3138: 1, 3131: 1, 3118: 1, 3107: 1, 3105: 1, 3093: 1,  
3091: 1, 3075: 1, 3065: 1, 3059: 1, 3053: 1, 3051: 1, 3049: 1, 3048: 1, 3047: 1,  
3046: 1, 3042: 1, 3038: 1, 3033: 1, 3014: 1, 3012: 1, 3009: 1, 3008: 1, 3007: 1,  
3000: 1, 2998: 1, 2991: 1, 2988: 1, 2984: 1, 2983: 1, 2977: 1, 2976: 1, 2973: 1,  
2968: 1, 2966: 1, 2959: 1, 2957: 1, 2945: 1, 2941: 1, 2934: 1, 2929: 1, 2927: 1,  
2919: 1, 2918: 1, 2908: 1, 2906: 1, 2901: 1, 2896: 1, 2895: 1, 2884: 1, 2878: 1,  
2877: 1, 2873: 1, 2872: 1, 2864: 1, 2863: 1, 2862: 1, 2861: 1, 2858: 1, 2848: 1,  
2844: 1, 2842: 1, 2836: 1, 2832: 1, 2825: 1, 2824: 1, 2822: 1, 2821: 1, 2818: 1,  
2815: 1, 2811: 1, 2808: 1, 2797: 1, 2795: 1, 2794: 1, 2789: 1, 2781: 1, 2775: 1,  
2772: 1, 2758: 1, 2757: 1, 2752: 1, 2746: 1, 2737: 1, 2735: 1, 2732: 1, 2728: 1,  
2727: 1, 2725: 1, 2718: 1, 2710: 1, 2708: 1, 2707: 1, 2704: 1, 2702: 1, 2701: 1,  
2698: 1, 2690: 1, 2682: 1, 2679: 1, 2678: 1, 2671: 1, 2662: 1, 2661: 1, 2660: 1,  
2658: 1, 2654: 1, 2653: 1, 2648: 1, 2644: 1, 2643: 1, 2637: 1, 2636: 1, 2635: 1,  
2634: 1, 2630: 1, 2622: 1, 2619: 1, 2618: 1, 2615: 1, 2614: 1, 2611: 1, 2604: 1,  
2601: 1, 2599: 1, 2597: 1, 2595: 1, 2588: 1, 2581: 1, 2578: 1, 2573: 1, 2572: 1,  
2571: 1, 2568: 1, 2564: 1, 2560: 1, 2556: 1, 2555: 1, 2554: 1, 2552: 1, 2547: 1,  
2545: 1, 2543: 1, 2542: 1, 2541: 1, 2538: 1, 2536: 1, 2535: 1, 2526: 1, 2521: 1,

2520: 1, 2519: 1, 2517: 1, 2515: 1, 2514: 1, 2512: 1, 2510: 1, 2509: 1, 2508: 1,  
2504: 1, 2498: 1, 2493: 1, 2488: 1, 2484: 1, 2475: 1, 2474: 1, 2473: 1, 2466: 1,  
2464: 1, 2454: 1, 2452: 1, 2450: 1, 2442: 1, 2441: 1, 2440: 1, 2438: 1, 2436: 1,  
2435: 1, 2433: 1, 2427: 1, 2426: 1, 2422: 1, 2421: 1, 2420: 1, 2419: 1, 2416: 1,  
2413: 1, 2412: 1, 2410: 1, 2407: 1, 2406: 1, 2405: 1, 2394: 1, 2393: 1, 2392: 1,  
2391: 1, 2390: 1, 2386: 1, 2385: 1, 2382: 1, 2373: 1, 2371: 1, 2370: 1, 2368: 1,  
2367: 1, 2366: 1, 2362: 1, 2359: 1, 2356: 1, 2355: 1, 2351: 1, 2348: 1, 2346: 1,  
2345: 1, 2344: 1, 2343: 1, 2341: 1, 2323: 1, 2322: 1, 2311: 1, 2309: 1, 2305: 1,  
2303: 1, 2301: 1, 2299: 1, 2297: 1, 2295: 1, 2290: 1, 2288: 1, 2286: 1, 2285: 1,  
2283: 1, 2278: 1, 2277: 1, 2275: 1, 2273: 1, 2272: 1, 2258: 1, 2254: 1, 2251: 1,  
2246: 1, 2245: 1, 2244: 1, 2243: 1, 2238: 1, 2236: 1, 2232: 1, 2231: 1, 2223: 1,  
2214: 1, 2213: 1, 2212: 1, 2210: 1, 2200: 1, 2199: 1, 2198: 1, 2197: 1, 2193: 1,  
2191: 1, 2188: 1, 2186: 1, 2185: 1, 2184: 1, 2179: 1, 2176: 1, 2175: 1, 2170: 1,  
2169: 1, 2164: 1, 2161: 1, 2157: 1, 2156: 1, 2154: 1, 2152: 1, 2151: 1, 2147: 1,  
2146: 1, 2144: 1, 2137: 1, 2136: 1, 2135: 1, 2132: 1, 2124: 1, 2123: 1, 2122: 1,  
2110: 1, 2107: 1, 2105: 1, 2102: 1, 2098: 1, 2090: 1, 2087: 1, 2084: 1, 2083: 1,  
2075: 1, 2074: 1, 2071: 1, 2067: 1, 2066: 1, 2064: 1, 2063: 1, 2058: 1, 2054: 1,  
2049: 1, 2047: 1, 2043: 1, 2041: 1, 2031: 1, 2029: 1, 2027: 1, 2025: 1, 2023: 1,  
2022: 1, 2021: 1, 2018: 1, 2017: 1, 2015: 1, 2012: 1, 2011: 1, 2006: 1, 2002: 1,  
1998: 1, 1997: 1, 1990: 1, 1984: 1, 1980: 1, 1973: 1, 1970: 1, 1969: 1, 1967: 1,  
1966: 1, 1965: 1, 1964: 1, 1962: 1, 1961: 1, 1957: 1, 1954: 1, 1950: 1, 1947: 1,  
1943: 1, 1941: 1, 1939: 1, 1938: 1, 1935: 1, 1933: 1, 1929: 1, 1928: 1, 1927: 1,  
1925: 1, 1922: 1, 1917: 1, 1916: 1, 1915: 1, 1913: 1, 1911: 1, 1908: 1, 1904: 1,  
1900: 1, 1899: 1, 1894: 1, 1893: 1, 1892: 1, 1891: 1, 1887: 1, 1882: 1, 1881: 1,  
1880: 1, 1879: 1, 1876: 1, 1875: 1, 1873: 1, 1872: 1, 1869: 1, 1860: 1, 1858: 1,  
1856: 1, 1853: 1, 1848: 1, 1847: 1, 1843: 1, 1842: 1, 1839: 1, 1838: 1, 1836: 1,  
1832: 1, 1831: 1, 1830: 1, 1825: 1, 1824: 1, 1822: 1, 1821: 1, 1819: 1, 1817: 1,  
1811: 1, 1810: 1, 1809: 1, 1805: 1, 1804: 1, 1791: 1, 1787: 1, 1780: 1, 1778: 1,  
1776: 1, 1775: 1, 1774: 1, 1771: 1, 1768: 1, 1763: 1, 1759: 1, 1758: 1, 1757: 1,  
1753: 1, 1752: 1, 1751: 1, 1750: 1, 1748: 1, 1746: 1, 1745: 1, 1744: 1, 1740: 1,  
1738: 1, 1737: 1, 1724: 1, 1723: 1, 1721: 1, 1719: 1, 1709: 1, 1706: 1, 1695: 1,  
1688: 1, 1683: 1, 1679: 1, 1674: 1, 1673: 1, 1672: 1, 1670: 1, 1669: 1, 1665: 1,  
1660: 1, 1658: 1, 1655: 1, 1652: 1, 1646: 1, 1643: 1, 1642: 1, 1641: 1, 1640: 1,  
1634: 1, 1633: 1, 1631: 1, 1623: 1, 1621: 1, 1618: 1, 1616: 1, 1615: 1, 1614: 1,  
1613: 1, 1611: 1, 1610: 1, 1603: 1, 1602: 1, 1597: 1, 1592: 1, 1591: 1, 1589: 1,  
1587: 1, 1583: 1, 1581: 1, 1580: 1, 1578: 1, 1577: 1, 1571: 1, 1570: 1, 1566: 1,  
1565: 1, 1564: 1, 1561: 1, 1558: 1, 1557: 1, 1555: 1, 1553: 1, 1551: 1, 1549: 1,  
1545: 1, 1544: 1, 1542: 1, 1540: 1, 1536: 1, 1535: 1, 1532: 1, 1528: 1, 1526: 1,  
1523: 1, 1521: 1, 1520: 1, 1519: 1, 1518: 1, 1516: 1, 1515: 1, 1511: 1, 1508: 1,  
1507: 1, 1504: 1, 1502: 1, 1501: 1, 1496: 1, 1492: 1, 1487: 1, 1483: 1, 1479: 1,  
1475: 1, 1473: 1, 1469: 1, 1468: 1, 1464: 1, 1463: 1, 1458: 1, 1457: 1, 1455: 1,  
1454: 1, 1453: 1, 1451: 1, 1450: 1, 1447: 1, 1445: 1, 1443: 1, 1442: 1, 1439: 1,  
1436: 1, 1433: 1, 1432: 1, 1431: 1, 1428: 1, 1427: 1, 1425: 1, 1420: 1, 1419: 1,  
1418: 1, 1415: 1, 1409: 1, 1408: 1, 1407: 1, 1406: 1, 1405: 1, 1404: 1, 1399: 1,  
1398: 1, 1396: 1, 1387: 1, 1384: 1, 1378: 1, 1375: 1, 1373: 1, 1366: 1, 1362: 1,  
1361: 1, 1359: 1, 1351: 1, 1350: 1, 1347: 1, 1345: 1, 1344: 1, 1341: 1, 1340: 1,  
1338: 1, 1336: 1, 1331: 1, 1329: 1, 1327: 1, 1326: 1, 1318: 1, 1315: 1, 1314: 1,  
1311: 1, 1307: 1, 1306: 1, 1301: 1, 1295: 1, 1291: 1, 1290: 1, 1289: 1, 1284: 1,

```

1283: 1, 1282: 1, 1280: 1, 1279: 1, 1278: 1, 1273: 1, 1267: 1, 1266: 1, 1263: 1,
1252: 1, 1251: 1, 1249: 1, 1240: 1, 1238: 1, 1234: 1, 1233: 1, 1232: 1, 1228: 1,
1223: 1, 1218: 1, 1215: 1, 1214: 1, 1208: 1, 1206: 1, 1205: 1, 1204: 1, 1201: 1,
1198: 1, 1195: 1, 1190: 1, 1189: 1, 1188: 1, 1180: 1, 1178: 1, 1171: 1, 1165: 1,
1164: 1, 1160: 1, 1149: 1, 1148: 1, 1145: 1, 1142: 1, 1140: 1, 1136: 1, 1135: 1,
1130: 1, 1128: 1, 1122: 1, 1119: 1, 1110: 1, 1109: 1, 1104: 1, 1101: 1, 1099: 1,
1097: 1, 1096: 1, 1095: 1, 1093: 1, 1086: 1, 1085: 1, 1083: 1, 1075: 1, 1069: 1,
1067: 1, 1065: 1, 1063: 1, 1062: 1, 1061: 1, 1053: 1, 1046: 1, 1041: 1, 1038: 1,
1034: 1, 1031: 1, 1030: 1, 1026: 1, 1024: 1, 1018: 1, 1017: 1, 1016: 1, 1015: 1,
1013: 1, 1011: 1, 1007: 1, 1006: 1, 1003: 1, 999: 1, 996: 1, 990: 1, 989: 1,
987: 1, 983: 1, 979: 1, 977: 1, 968: 1, 965: 1, 959: 1, 952: 1, 950: 1, 946: 1,
944: 1, 927: 1, 912: 1, 908: 1, 900: 1, 897: 1, 895: 1, 893: 1, 890: 1, 877: 1,
875: 1, 872: 1, 868: 1, 866: 1, 864: 1, 860: 1, 855: 1, 849: 1, 820: 1, 812: 1,
782: 1, 748: 1, 736: 1, 732: 1, 705: 1, 698: 1, 690: 1, 661: 1})

```

```

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

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
↳ generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
↳ fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
↳ learning_rate=optimal, eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

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

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

    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_text_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)

```

```

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

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

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

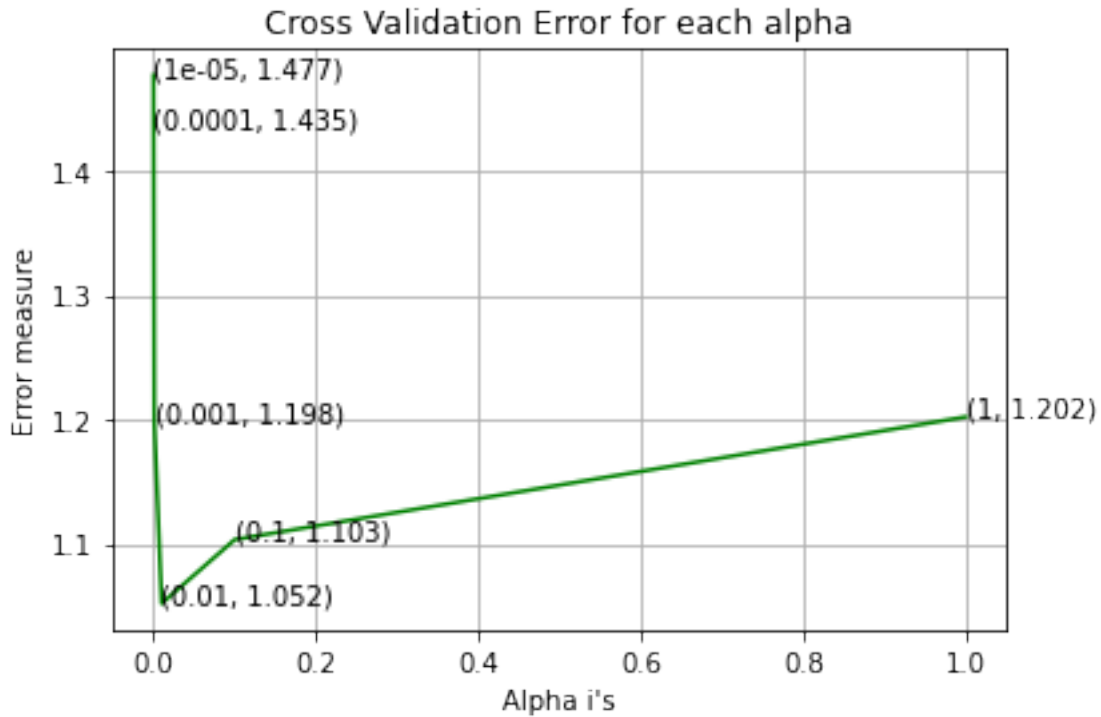
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:
→",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best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

```

```

For values of alpha = 1e-05 The log loss is: 1.4768185342682982
For values of alpha = 0.0001 The log loss is: 1.4346906310450065
For values of alpha = 0.001 The log loss is: 1.1978989975277825
For values of alpha = 0.01 The log loss is: 1.0520080204659137
For values of alpha = 0.1 The log loss is: 1.1033830631722568
For values of alpha = 1 The log loss is: 1.2020149247829088

```



For values of best alpha = 0.01 The train log loss is: 0.7513084613362182

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

1.0520080204659137

For values of best alpha = 0.01 The test log loss is: 1.0962167451622822

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

Ans. Yes, it seems like!

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

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

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

97.097 % of word of test data appeared in train data  
97.421 % of word of Cross Validation appeared in train data

#### 4. Machine Learning Models

```
[ ]: #Data preparation for ML models.

#Misc. functions for ML models

def predict_and_plot_confusion_matrix(train_x, train_y, test_x, test_y, clf):
    clf.fit(train_x, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x, train_y)
    pred_y = sig_clf.predict(test_x)

    # for calculating log_loss we will provide the array of probabilities
    # belongs to each class
    print("Log loss :", log_loss(test_y, sig_clf.predict_proba(test_x)))
    # calculating the number of data points that are misclassified
    print("Number of mis-classified points :", np.count_nonzero((pred_y -
    test_y))/test_y.shape[0])
    plot_confusion_matrix(test_y, pred_y)

[ ]: 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)

[ ]: # this function will be used just for naive bayes
# for the given indices, we will print the name of the features
# and we will check whether the feature present in the test point text or not
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = CountVectorizer(ngram_range=(1,3))
    var_count_vec = CountVectorizer(ngram_range=(1,3))
    text_count_vec = CountVectorizer(min_df=3, ngram_range=(1,3))

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

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

    word_present = 0
    for i, v in enumerate(indices):
        if (v < fea1_len):
```

```

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

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

```

Stacking the three types of BOW features

```

[:]: # merging gene, variance and text features

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

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

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

```



```

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

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

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

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

```

```

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

```

One hot encoding features :

```

(number of data points * number of features) in train data = (2124, 1863528)
(number of data points * number of features) in test data = (665, 1863528)
(number of data points * number of features) in cross validation data = (532,
1863528)

```

```

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

```

Response encoding features :

```

(number of data points * number of features) in train data = (2124, 27)

```

(number of data points \* number of features) in test data = (665, 27)  
(number of data points \* number of features) in cross validation data = (532, 27)

#### 4.1. Base Line Model

##### Logistic Regression with BOW ngram =(1,3) and Class balancing

```
[ ]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
      ↳ generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
      ↳ fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
      ↳ learning_rate=optimal, eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/
      ↳ lessons/geometric-intuition-1/
#-----

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

alpha = [10 ** x for x in range(-6, 3)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
```

```

    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2',
→loss='log', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.
→classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilities we use
→log-probability estimates
    print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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

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

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

```

```

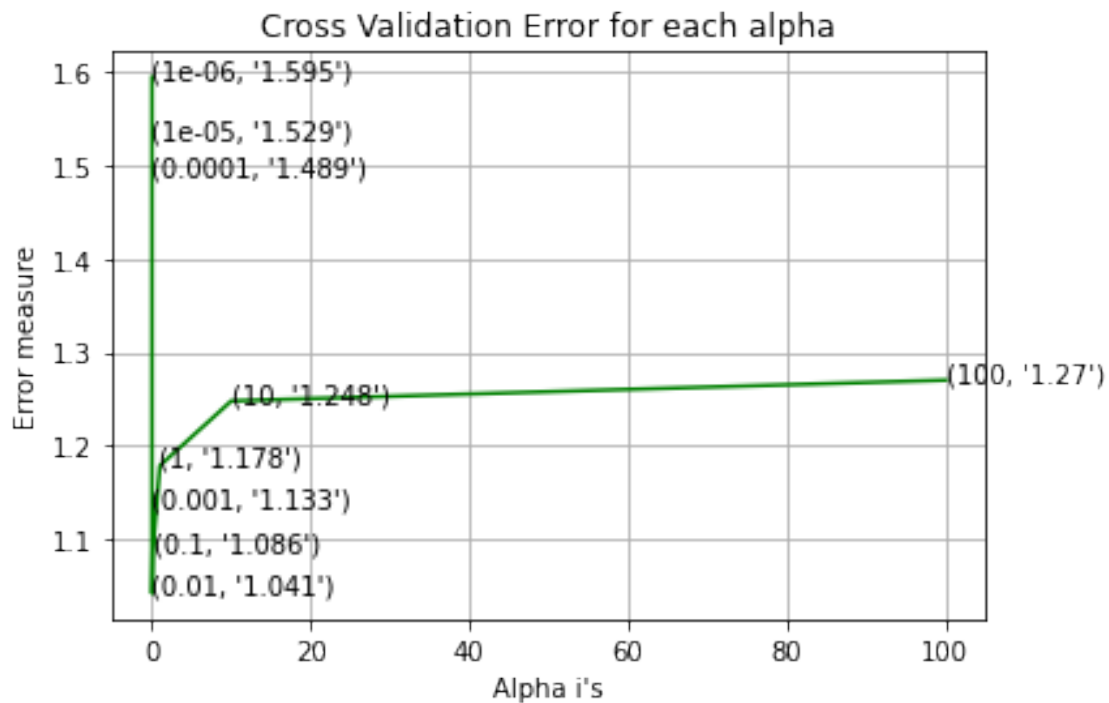
for alpha = 1e-06
Log Loss : 1.5950238895067985
for alpha = 1e-05
Log Loss : 1.528519025799228
for alpha = 0.0001
Log Loss : 1.4889089243498643
for alpha = 0.001

```

```

Log Loss : 1.133179264821804
for alpha = 0.01
Log Loss : 1.041066336881157
for alpha = 0.1
Log Loss : 1.085719611561209
for alpha = 1
Log Loss : 1.1781582970759825
for alpha = 10
Log Loss : 1.247630102990052
for alpha = 100
Log Loss : 1.269920522662356

```



```

For values of best alpha = 0.01 The train log loss is: 0.7335814397454019
For values of best alpha = 0.01 The cross validation log loss is:
1.041066336881157
For values of best alpha = 0.01 The test log loss is: 1.0984763733995617

```

Testing the model with best hyper paramters

```

[ ]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
      ↳ generated/sklearn.linear_model.SGDClassifier.html
      # -----
      # default parameters
      # SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
      ↳ fit_intercept=True, max_iter=None, tol=None,

```

```
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
→ learning_rate=optimal, eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/
→ lessons/geometric-intuition-1/
#-----

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

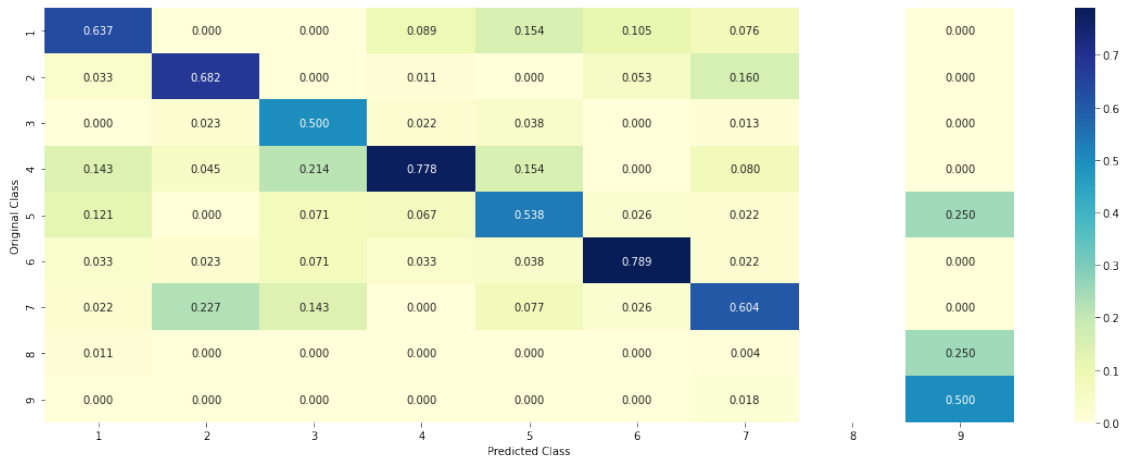
Log loss : 1.041066336881157

Number of mis-classified points : 0.34774436090225563

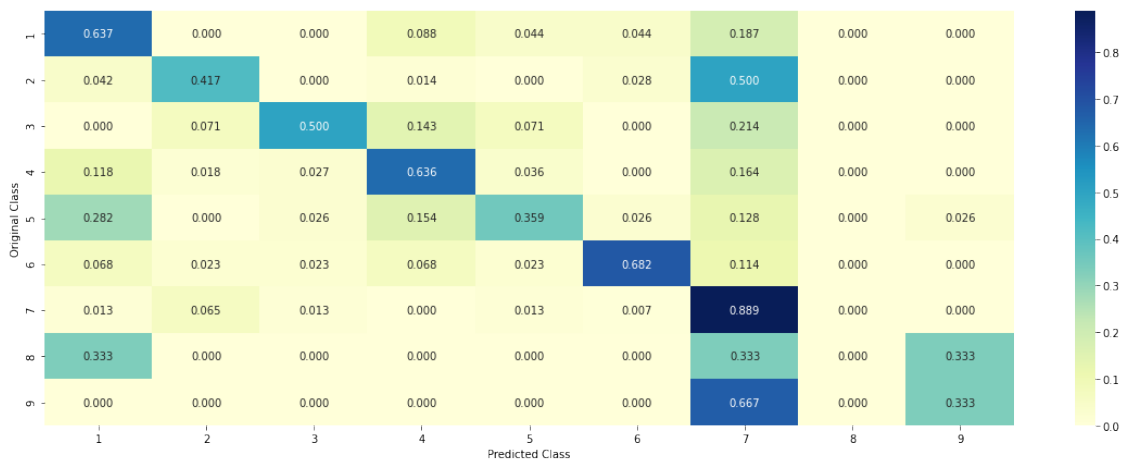
----- Confusion matrix -----



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



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



## Feature Importance

```
[ ]: def get_imp_feature_names(text, indices, removed_ind = []):
    word_present = 0
    tabulte_list = []
    increasingorder_ind = 0
    for i in indices:
        if i < train_gene_feature_onehotCoding.shape[1]:
            tabulte_list.append([increasingorder_ind, "Gene", "Yes"])
        elif i < 18:
            tabulte_list.append([increasingorder_ind, "Variation", "Yes"])
        if ((i > 17) & (i not in removed_ind)) :
            word = train_text_features[i]
            yes_no = True if word in text.split() else False
```

```

        if yes_no:
            word_present += 1
            tabulte_list.append([increasingorder_ind, train_text_features[i],
→yes_no])
            increasingorder_ind += 1
        print(word_present, "most important features are present in our query
→point")
        print("-"*50)
        print("The features that are most important of the ", predicted_cls[0], "
→class:")
        print(tabulate(tabulte_list, headers=["Index", 'Feature name', 'Present or
→Not']))

```

Correctly Classified point

```

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

```

Predicted Class : 6

Predicted Class Probabilities: [[0.0064 0.0248 0.0037 0.0081 0.1144 0.8148  
0.0186 0.0045 0.0047]]

Actual Class : 6

```

-----
389 Text feature [previously] present in test data point [True]
397 Text feature [performed] present in test data point [True]
412 Text feature [found] present in test data point [True]
433 Text feature [presence] present in test data point [True]
435 Text feature [described] present in test data point [True]
498 Text feature [results] present in test data point [True]
Out of the top 500 features 6 are present in query point

```

Incorrectly Classified point

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

Predicted Class : 7

Predicted Class Probabilities: [[0.1994 0.1386 0.0222 0.1884 0.0651 0.0585  
0.3109 0.0062 0.0106]]

Actual Class : 7

```
-----
151 Text feature [missense] present in test data point [True]
190 Text feature [loss] present in test data point [True]
223 Text feature [function] present in test data point [True]
224 Text feature [mim] present in test data point [True]
226 Text feature [individuals] present in test data point [True]
289 Text feature [suppressor] present in test data point [True]
321 Text feature [protein] present in test data point [True]
324 Text feature [dna] present in test data point [True]
325 Text feature [affected] present in test data point [True]
Out of the top 500 features 9 are present in query point
```

Without Class balancing

```
[ ]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
    ↳generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
    ↳fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
    ↳learning_rate=optimal, eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

#-----
```



```

# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/
# → lessons/geometric-intuition-1/
#-----

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

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

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

best_alpha = np.argmin(cv_log_error_array)

```

```

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

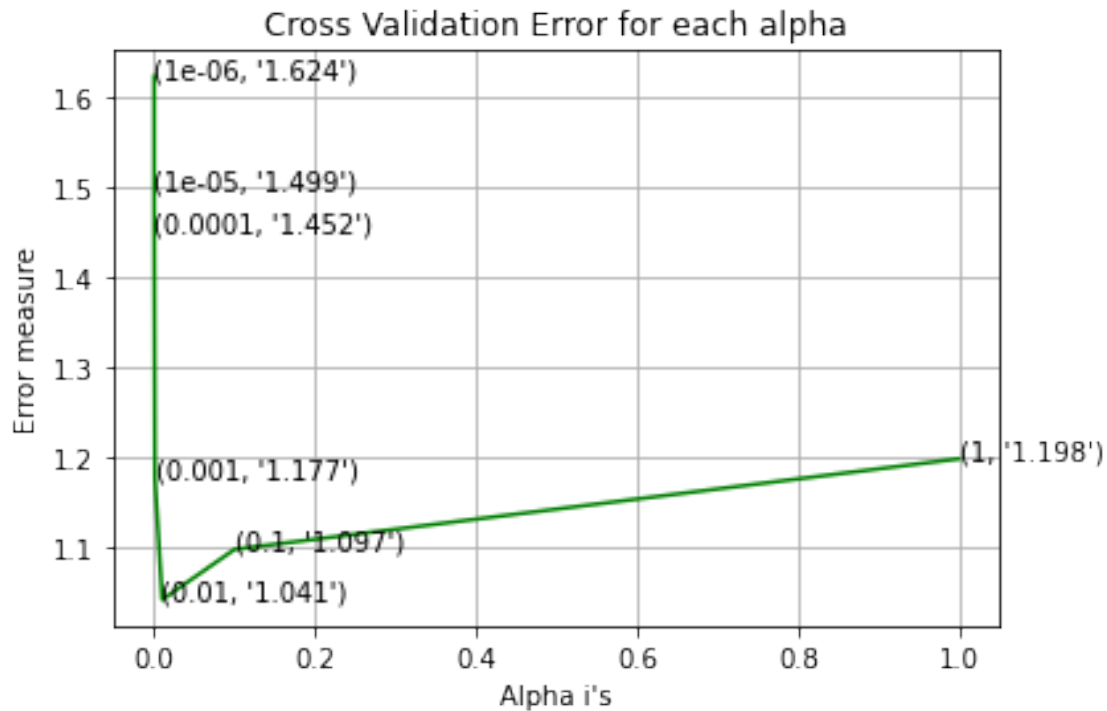
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:
    ", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 1e-06
Log Loss : 1.6238007160702748
for alpha = 1e-05
Log Loss : 1.4993642321529355
for alpha = 0.0001
Log Loss : 1.4516479814262566
for alpha = 0.001
Log Loss : 1.1767979507818411
for alpha = 0.01
Log Loss : 1.041276211781357
for alpha = 0.1
Log Loss : 1.0973385152220392
for alpha = 1
Log Loss : 1.1981450667217117

```



For values of best alpha = 0.01 The train log loss is: 0.7272041854788291  
 For values of best alpha = 0.01 The cross validation log loss is:  
 1.041276211781357  
 For values of best alpha = 0.01 The test log loss is: 1.0895054100865529

Testing model with best hyper parameters

```
[ ]: # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/
      ↳ generated/sklearn.linear_model.SGDClassifier.html
      # -----
      # default parameters
      # SGDClassifier(loss=hinge, penalty=l2, alpha=0.0001, l1_ratio=0.15,
      ↳ fit_intercept=True, max_iter=None, tol=None,
      # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None,
      ↳ learning_rate=optimal, eta0=0.0, power_t=0.5,
      # class_weight=None, warm_start=False, average=False, n_iter=None)

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

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

```
#-----

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

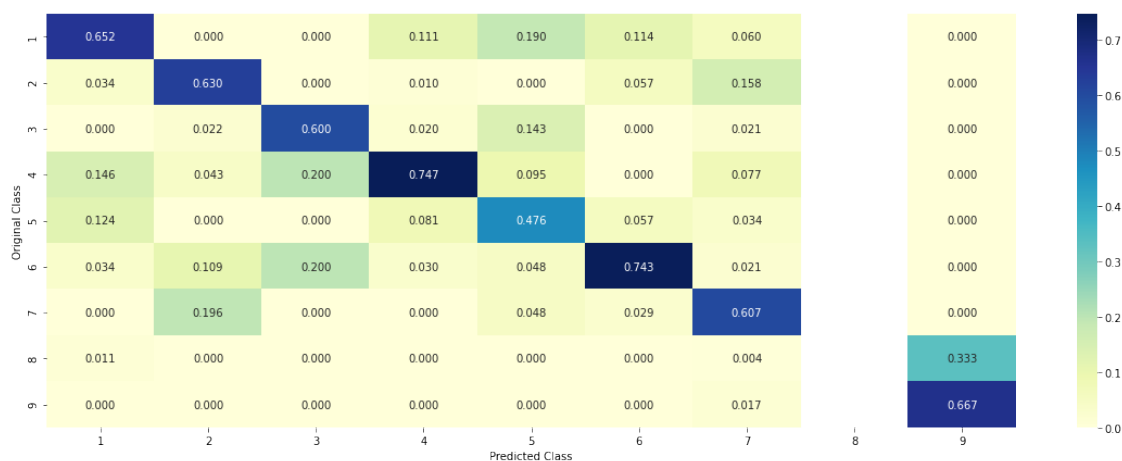
Log loss : 1.041276211781357

Number of mis-classified points : 0.3533834586466165

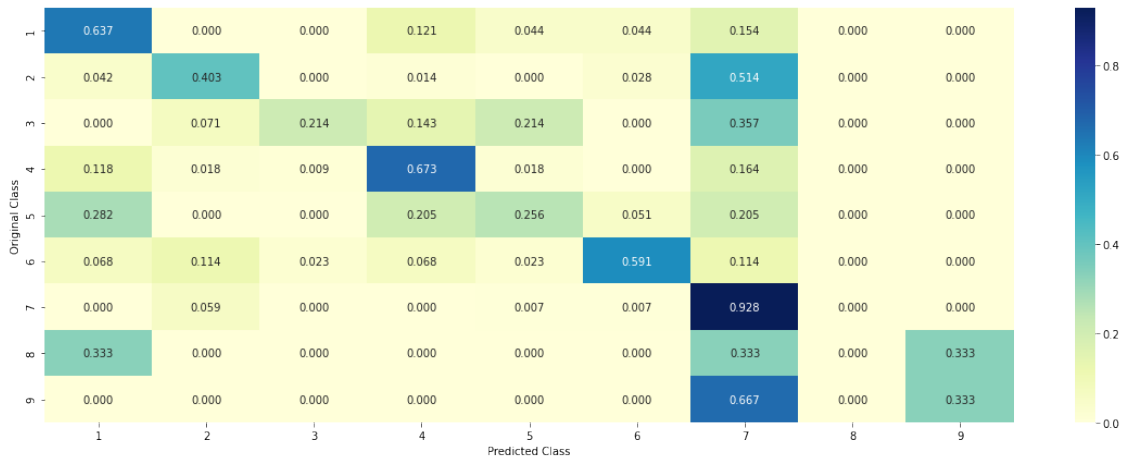
----- Confusion matrix -----



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



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



### Feature Importance, Correctly Classified point

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

Predicted Class : 6

Predicted Class Probabilities: [[0.0066 0.0266 0.0039 0.0088 0.0949 0.8343 0.02  
0.0034 0.0015]]

Actual Class : 6

```
-----
302 Text feature [previously] present in test data point [True]
310 Text feature [performed] present in test data point [True]
315 Text feature [found] present in test data point [True]
332 Text feature [described] present in test data point [True]
333 Text feature [presence] present in test data point [True]
352 Text feature [additional] present in test data point [True]
353 Text feature [results] present in test data point [True]
360 Text feature [detected] present in test data point [True]
```

```

402 Text feature [resulting] present in test data point [True]
408 Text feature [two] present in test data point [True]
411 Text feature [cells] present in test data point [True]
414 Text feature [shown] present in test data point [True]
417 Text feature [identified] present in test data point [True]
426 Text feature [determine] present in test data point [True]
441 Text feature [pcr] present in test data point [True]
445 Text feature [protein] present in test data point [True]
453 Text feature [using] present in test data point [True]
456 Text feature [suggested] present in test data point [True]
462 Text feature [showed] present in test data point [True]
469 Text feature [mutation] present in test data point [True]
475 Text feature [show] present in test data point [True]
476 Text feature [addition] present in test data point [True]
484 Text feature [containing] present in test data point [True]
492 Text feature [mutations] present in test data point [True]
496 Text feature [type] present in test data point [True]
497 Text feature [sequenced] present in test data point [True]
498 Text feature [characterized] present in test data point [True]
Out of the top 500 features 27 are present in query point

```

```

[ ]: test_point_index = 14
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.
    →predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].
    →iloc[test_point_index], test_df['Gene'].
    →iloc[test_point_index], test_df['Variation'].iloc[test_point_index],
    →no_feature)

```

```

Predicted Class : 2
Predicted Class Probabilities: [[0.0043 0.3969 0.0026 0.0021 0.2726 0.0007
0.3156 0.0043 0.0008]]
Actual Class : 7

```

```

-----
316 Text feature [function] present in test data point [True]
330 Text feature [type] present in test data point [True]
454 Text feature [whether] present in test data point [True]
467 Text feature [dominant] present in test data point [True]
487 Text feature [indicated] present in test data point [True]
Out of the top 500 features 5 are present in query point

```

<li>Apply Logistic regression with CountVectorizer Features, including both unigrams and bigrams