

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> (<https://www.forbes.com/sites/matthewherper/2017/06/03/a-new-cancer-drug-helped-almost-everyone-who-took-it-almost-heres-what-it-teaches-us/#2a44ee2f6b25>)
2. <https://www.youtube.com/watch?v=UwbuW7oK8rk> (<https://www.youtube.com/watch?v=UwbuW7oK8rk>)
3. <https://www.youtube.com/watch?v=qxXRKVompl8> (<https://www.youtube.com/watch?v=qxXRKVompl8>)

1.3. Real-world/Business objectives and constraints.

- No low-latency requirement.
- Interpretability is important.
- Errors can be very costly.

- Probability of a data-point belonging to each class is needed.

2. Machine Learning Problem Formulation

2.1. Data

2.1.1. Data Overview

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

2.1.2. Example Data Point

training_variants

ID, Gene, Variation, Class
0, FAM58A, Truncating Mutations, 1
1, CBL, W802*, 2
2, CBL, Q249E, 2
...

training_text

ID, Text
0|Cyclin-dependent kinases (CDKs) regulate a variety of fundamental cellular processes. CDK10 stands out as one of the last orphan CDKs for which no activating cyclin has been identified and no kinase activity revealed. Previous work has shown that CDK10 silencing increases ETS2 (v-ets erythroblastosis virus E26 oncogene homolog 2)-driven activation of the MAPK pathway, which confers tamoxifen resistance to breast cancer cells. The precise mechanisms by which CDK10 modulates ETS2 activity, and more generally the functions of CDK10, remain elusive. Here we

demonstrate that CDK10 is a cyclin-dependent kinase by identifying cyclin M as an activating cyclin. Cyclin M, an orphan cyclin, is the product of FAM58A, whose mutations cause STAR syndrome, a human developmental anomaly whose features include toe syndactyly, telecanthus, and anogenital and renal malformations. We show that STAR syndrome-associated cyclin M mutants are unable to interact with CDK10. Cyclin M silencing phenocopies CDK10 silencing in increasing c-Raf and in conferring tamoxifen resistance to breast cancer cells. CDK10/cyclin M phosphorylates ETS2 in vitro, and in cells it positively controls ETS2 degradation by the proteasome. ETS2 protein levels are increased in cells derived from a STAR patient, and this increase is attributable to decreased cyclin M levels. Altogether, our results reveal an additional regulatory mechanism for ETS2, which plays key roles in cancer and development. They also shed light on the molecular mechanisms underlying STAR syndrome. Cyclin-dependent kinases (CDKs) play a pivotal role in the control of a number of fundamental cellular processes (1). The human genome contains 21 genes encoding proteins that can be considered as members of the CDK family owing to their sequence similarity with bona fide CDKs, those known to be activated by cyclins (2). Although discovered almost 20 y ago (3, 4), CDK10 remains one of the two CDKs without an identified cyclin partner. This knowledge gap has largely impeded the exploration of its biological functions. CDK10 can act as a positive cell cycle regulator in some cells (5, 6) or as a tumor suppressor in others (7, 8). CDK10 interacts with the ETS2 (v-ets erythroblastosis virus E26 oncogene homolog 2) transcription factor and inhibits its transcriptional activity through an unknown mechanism (9). CDK10 knockdown derepresses ETS2, which increases the expression of the c-Raf protein kinase, activates the MAPK pathway, and induces resistance of MCF7 cells to tamoxifen (6). ...

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

2.2.1. Type of Machine Learning Problem

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

2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

2.3. Train, CV and Test Datasets

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

3. Exploratory Data Analysis

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import re
4 import time
5 import warnings
6 import numpy as np
7 from nltk.corpus import stopwords
8 from sklearn.decomposition import TruncatedSVD
9 from sklearn.preprocessing import normalize
10 from sklearn.feature_extraction.text import CountVectorizer
11 from sklearn.manifold import TSNE
12 import seaborn as sns
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.metrics import confusion_matrix
15 from sklearn.metrics.classification import accuracy_score, log_loss
16 from sklearn.feature_extraction.text import TfidfVectorizer
17 from sklearn.linear_model import SGDClassifier
18 from imblearn.over_sampling import SMOTE
19 from collections import Counter
20 from scipy.sparse import hstack
21 from sklearn.multiclass import OneVsRestClassifier
22 from sklearn.svm import SVC
23 from sklearn.model_selection import StratifiedKFold
24 from collections import Counter, defaultdict
25 from sklearn.calibration import CalibratedClassifierCV
26 from sklearn.naive_bayes import MultinomialNB
27 from sklearn.naive_bayes import GaussianNB
28 from sklearn.model_selection import train_test_split
29 from sklearn.model_selection import GridSearchCV
30 import math
31 from sklearn.metrics import normalized_mutual_info_score
32 from sklearn.ensemble import RandomForestClassifier
33 warnings.filterwarnings("ignore")
34
```

```

35 from mlxtend.classifier import StackingClassifier
36
37 from sklearn import model_selection
38 from sklearn.linear_model import LogisticRegression
39

```

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

```

1 data = pd.read_csv('training/training_variants')
2 print('Number of data points : ', data.shape[0])
3 print('Number of features : ', data.shape[1])
4 print('Features : ', data.columns.values)
5 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 mutations used for training.

Fields are

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

3.1.2. Reading Text Data

```

1 # note the separator in this file
2 data_text = pd.read_csv("training/training_text", sep="\|", engine="python", names=["ID", "TEXT"], skiprows:
3 print('Number of data points : ', data_text.shape[0])
4 print('Number of features : ', data_text.shape[1])
5 print('Features : ', data_text.columns.values)
6 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...

```

1 #merging both gene_variations and text data based on ID
2 result = pd.merge(data, data_text, on='ID', how='left')
3 result.head()

```

	ID	Gene	Variation	Class	TEXT
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var...
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc...
2	2	CBL	Q249E	2	Abstract Background Non-small cell lung canc...
3	3	CBL	N454D	3	Recent evidence has demonstrated that acquired...
4	4	CBL	L399V	4	Oncogenic mutations in the monomeric Casitas B...

```
1 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

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

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

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

Feature Engineering

Feature engineering reference - <https://www.analyticsvidhya.com/blog/2018/02/the-different-methods-deal-text-data-predictive-python/> (<https://www.analyticsvidhya.com/blog/2018/02/the-different-methods-deal-text-data-predictive-python/>)


```
1  #Variation share
2  result['Variation_Share'] = result.apply(lambda r: sum([1 for w in r['Variation'].split(' ') if w in r[
3
4  #Text word count
5  result["Text_Word_Count"] = result["TEXT"].apply(lambda x: len(x.split()))
6
7  #Text Character count
8  result['Text_Character_Count'] = result['TEXT'].apply(lambda x: len(str(x)))
9
10 #Text average Length
11 result['Text_Avg_length'] = result['Text_Character_Count'] / result['Text_Word_Count']
12
13 #Gene Character count
14 result['Gene_Character_Count'] = result['Gene'].apply(lambda x: len(str(x)))
15
16 #Variation_Character_Count
17 result['Variation_Character_Count'] = result['Variation'].apply(lambda x: len(str(x)))
18
19 #Var word count
20 result["Var_Word_Count"] = result["Variation"].apply(lambda x: len(x.split()))
21
22 #Var average Length
23 result['var_Avg_length'] = result['Variation_Character_Count'] / result['Var_Word_Count']
24
25 #Stopword text count
26 from nltk.corpus import stopwords
27 stop = stopwords.words('english')
28
29 result['stopwordst'] = result['TEXT'].apply(lambda x: len([x for x in x.split() if x in stop]))
30
31 #Speacial character count
32 result['hastags'] = result['TEXT'].apply(lambda x: len([x for x in x.split() if x.startswith('#'))])
33
34 #Numerics count
```

```

35 result['numerics'] = result['TEXT'].apply(lambda x: len([x for x in x.split() if x.isdigit()]))
36
37 #Uppercase word count
38 result['upper'] = result['TEXT'].apply(lambda x: len([x for x in x.split() if x.isupper()]))
39
40 # Label encoder
41 from sklearn import preprocessing
42 for c in result.columns:
43     if result[c].dtype == 'object':
44         if c == 'Gene':
45             lbl = preprocessing.LabelEncoder()
46             result[c+'_lbl_enc'] = lbl.fit_transform(result[c].values)
47         elif c == 'Variation':
48             lbl = preprocessing.LabelEncoder()
49             result[c+'_lbl_enc'] = lbl.fit_transform(result[c].values)

```

```
1 result.head(2)
```

	ID	Gene	Variation	Class	TEXT	Variation_Share	Text_Word_Count	Text_Character_Count	Text_A
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var...	1	6089	39765	6.53062
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc...	1	5722	36831	6.43673

3.1.3. Preprocessing of text

```
1  # loading stop words from nltk library
2  stop_words = set(stopwords.words('english'))
3
4
5  def nlp_preprocessing(total_text, index, column):
6      if type(total_text) is not int:
7          string = ""
8          # replace every special char with space
9          total_text = re.sub('[^a-zA-Z0-9\n]', ' ', total_text)
10         # replace multiple spaces with single space
11         total_text = re.sub('\s+', ' ', total_text)
12         # converting all the chars into lower-case.
13         total_text = total_text.lower()
14
15         for word in total_text.split():
16             # if the word is a not a stop word then retain that word from the data
17             if not word in stop_words:
18                 string += word + " "
19
20         data_text[column][index] = string
```

```
1 #text processing stage.
2 start_time = time.clock()
3 for index, row in data_text.iterrows():
4     if type(row['TEXT']) is str:
5         nlp_preprocessing(row['TEXT'], index, 'TEXT')
6     else:
7         print("there is no text description for id:",index)
8 print('Time took for preprocessing the text :',time.clock() - start_time, "seconds")
```

```
there is no text description for id: 1109
there is no text description for id: 1277
there is no text description for id: 1407
there is no text description for id: 1639
there is no text description for id: 2755
Time took for preprocessing the text : 121.0215206174142 seconds
```

```
1 #Creating dummies for Gene_lbl_enc & Variation_lbl_enc
2 result = pd.concat([result, pd.get_dummies(result['Gene_lbl_enc']).\
3                     rename(columns=lambda x: 'Gene_lbl_enc_' + str(x))], axis=1)
4
5 result = pd.concat([result, pd.get_dummies(result['Variation_lbl_enc']).\
6                     rename(columns=lambda x: 'Variation_lbl_enc_' + str(x))], axis=1)
```

```
1 result.head(2)
```

	ID	Gene	Variation	Class	TEXT	Variation_Share	Text_Word_Count	Text_Character_Count	Text_A
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var...	1	6089	39765	6.53062
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc...	1	5722	36831	6.43673

2 rows x 3279 columns

Checking feature importance of only newly created variables using Random Forest

```
1 temp_y=result[['Class']]
2 temp_train=result.drop(['ID', 'Gene', 'Variation', 'Class', 'TEXT', 'Gene_lbl_enc', 'Variation_lbl_enc'],axis=
```

```
1 ## Import the random forest model.
2 from sklearn.ensemble import RandomForestClassifier
3 ## This line instantiates the model.
4 rf = RandomForestClassifier()
5 ## Fit the model on your training data.
6 rf.fit(temp_train, temp_y)
```

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

```

1 import pandas as pd
2 feature_importances = pd.DataFrame(rf.feature_importances_,
3                                   index = temp_train.columns,
4                                   columns=['importance']).sort_values('importance',ascending=False)

```

```

1 #Feature importance of newly created variables
2 feature_importances[:15]

```

	importance
stopwordst	0.053990
numerics	0.053875
Text_Word_Count	0.052921
Text_Avg_length	0.051516
upper	0.051413
Text_Character_Count	0.049281
var_Avg_length	0.030022
Variation_Character_Count	0.028296
Gene_Character_Count	0.023299
Variation_Share	0.018827
Gene_lbl_enc_31	0.015835
hastags	0.015086
Gene_lbl_enc_252	0.013824
Gene_lbl_enc_196	0.013398
Variation_lbl_enc_2629	0.012533

3.1.4. Test, Train and Cross Validation Split

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

```
1 y_true = result['Class'].values
2 result.Gene      = result.Gene.str.replace('\s+', '_')
3 result.Variation = result.Variation.str.replace('\s+', '_')
4
5 # split the data into test and train by maintaining same distribution of output variable 'y_true' [stratified]
6 X_train, test_df, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_size=0.2)
7 # split the train data into train and cross validation by maintaining same distribution of output variable
8 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

```
1 print('Number of data points in train data:', train_df.shape[0])
2 print('Number of data points in test data:', test_df.shape[0])
3 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

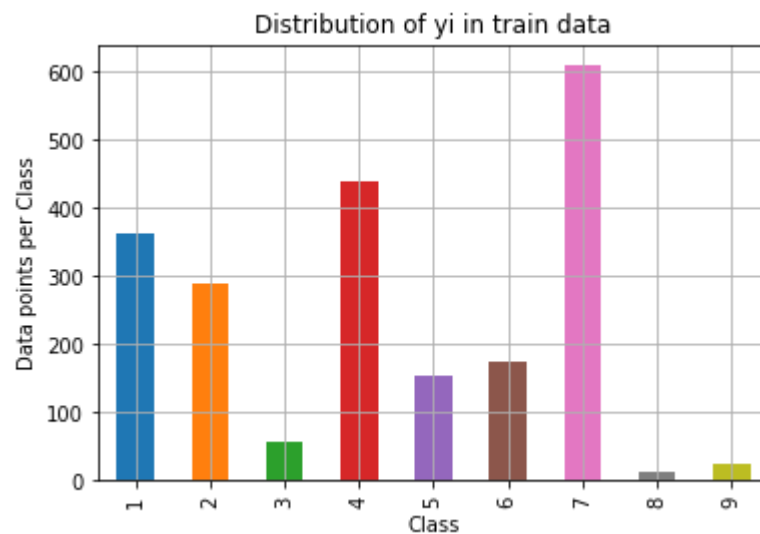
```
1  # it returns a dict, keys as class labels and values as the number of data points in that class
2  train_class_distribution = train_df['Class'].value_counts().sortlevel()
3  test_class_distribution = test_df['Class'].value_counts().sortlevel()
4  cv_class_distribution = cv_df['Class'].value_counts().sortlevel()
5
6  my_colors = 'rgbkymc'
7  train_class_distribution.plot(kind='bar')
8  plt.xlabel('Class')
9  plt.ylabel('Data points per Class')
10 plt.title('Distribution of yi in train data')
11 plt.grid()
12 plt.show()
13
14 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
15 # -(train_class_distribution.values): the minus sign will give us in decreasing order
16 sorted_yi = np.argsort(-train_class_distribution.values)
17 for i in sorted_yi:
18     print('Number of data points in class', i+1, ':', train_class_distribution.values[i], '(', np.round(
19
20
21 print('-'*80)
22 my_colors = 'rgbkymc'
23 test_class_distribution.plot(kind='bar')
24 plt.xlabel('Class')
25 plt.ylabel('Data points per Class')
26 plt.title('Distribution of yi in test data')
27 plt.grid()
28 plt.show()
29
30 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
31 # -(train_class_distribution.values): the minus sign will give us in decreasing order
32 sorted_yi = np.argsort(-test_class_distribution.values)
33 for i in sorted_yi:
34     print('Number of data points in class', i+1, ':', test_class_distribution.values[i], '(', np.round(
```



```

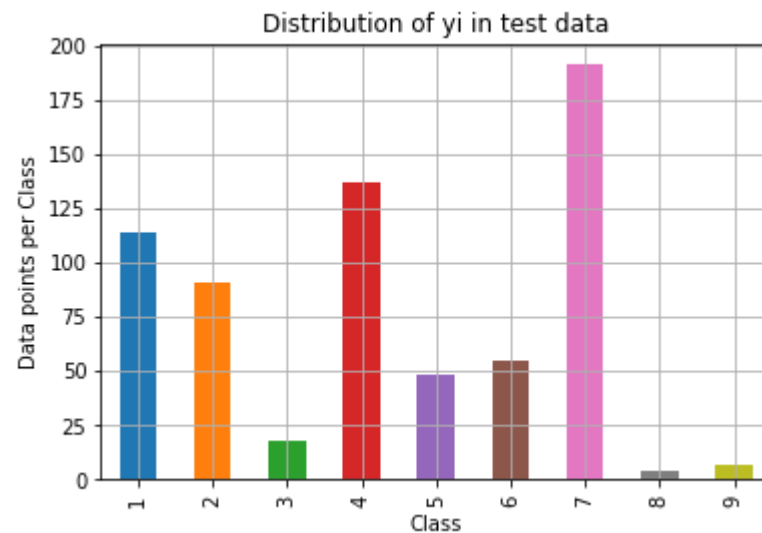
35
36 print('-'*80)
37 my_colors = 'rgbkymc'
38 cv_class_distribution.plot(kind='bar')
39 plt.xlabel('Class')
40 plt.ylabel('Data points per Class')
41 plt.title('Distribution of yi in cross validation data')
42 plt.grid()
43 plt.show()
44
45 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
46 # -(train_class_distribution.values): the minus sign will give us in decreasing order
47 sorted_yi = np.argsort(-train_class_distribution.values)
48 for i in sorted_yi:
49     print('Number of data points in class', i+1, ':', cv_class_distribution.values[i], '(', np.round((cv_
50

```

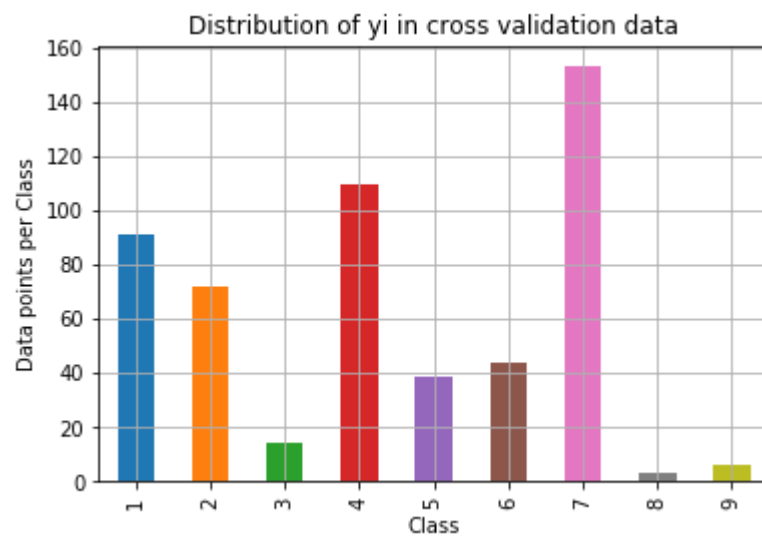


Number of data points in class 7 : 609 (28.672 %)
 Number of data points in class 4 : 439 (20.669 %)
 Number of data points in class 1 : 363 (17.09 %)
 Number of data points in class 2 : 289 (13.606 %)

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.

```

1  # This function plots the confusion matrices given y_i, y_i_hat.
2  def plot_confusion_matrix(test_y, predict_y):
3      C = confusion_matrix(test_y, predict_y)
4      # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j
5
6      A = (((C.T)/(C.sum(axis=1))).T)
7      #divid each element of the confusion matrix with the sum of elements in that column
8
9      # C = [[1, 2],
10     #      [3, 4]]
11     # C.T = [[1, 3],
12     #        [2, 4]]
13     # C.sum(axis = 1)  axis=0 corresonds to columns and axis=1 corresponds to rows in two dimensional
14     # C.sum(axix =1) = [[3, 7]]
15     # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
16     #                             [2/3, 4/7]]
17
18     # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
19     #                             [3/7, 4/7]]
20     # sum of row elements = 1
21
22     B = (C/C.sum(axis=0))
23     #divid each element of the confusion matrix with the sum of elements in that row
24     # C = [[1, 2],
25     #      [3, 4]]
26     # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to rows in two dimensional
27     # C.sum(axix =0) = [[4, 6]]
28     # (C/C.sum(axis=0)) = [[1/4, 2/6],
29     #                       [3/4, 4/6]]
30
31     labels = [1,2,3,4,5,6,7,8,9]
32     # representing A in heatmap format
33     print("-"*20, "Confusion matrix", "-"*20)
34     plt.figure(figsize=(20,7))

```

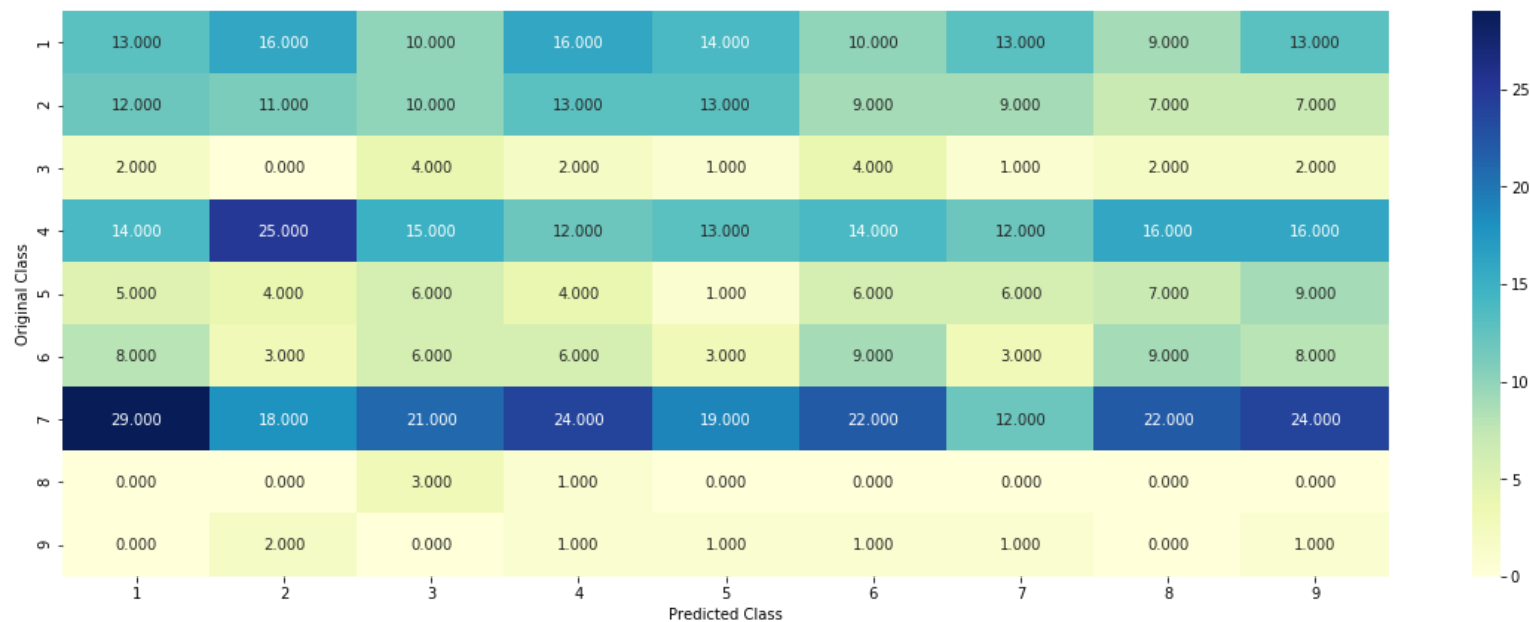
```
35 sns.heatmap(C, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
36 plt.xlabel('Predicted Class')
37 plt.ylabel('Original Class')
38 plt.show()
39
40 print("-"*20, "Precision matrix (Column Sum=1)", "-"*20)
41 plt.figure(figsize=(20,7))
42 sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
43 plt.xlabel('Predicted Class')
44 plt.ylabel('Original Class')
45 plt.show()
46
47 # representing B in heatmap format
48 print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
49 plt.figure(figsize=(20,7))
50 sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
51 plt.xlabel('Predicted Class')
52 plt.ylabel('Original Class')
53 plt.show()
```

```
1  # we need to generate 9 numbers and the sum of numbers should be 1
2  # one solution is to generate 9 numbers and divide each of the numbers by their sum
3  # ref: https://stackoverflow.com/a/18662466/4084039
4  test_data_len = test_df.shape[0]
5  cv_data_len = cv_df.shape[0]
6
7  # we create a output array that has exactly same size as the CV data
8  cv_predicted_y = np.zeros((cv_data_len,9))
9  for i in range(cv_data_len):
10     rand_probs = np.random.rand(1,9)
11     cv_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
12  print("Log loss on Cross Validation Data using Random Model",log_loss(y_cv,cv_predicted_y, eps=1e-15))
13
14
15  # Test-Set error.
16  #we create a output array that has exactly same as the test data
17  test_predicted_y = np.zeros((test_data_len,9))
18  for i in range(test_data_len):
19     rand_probs = np.random.rand(1,9)
20     test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
21  print("Log loss on Test Data using Random Model",log_loss(y_test,test_predicted_y, eps=1e-15))
22
23  predicted_y =np.argmax(test_predicted_y, axis=1)
24  plot_confusion_matrix(y_test, predicted_y+1)
```

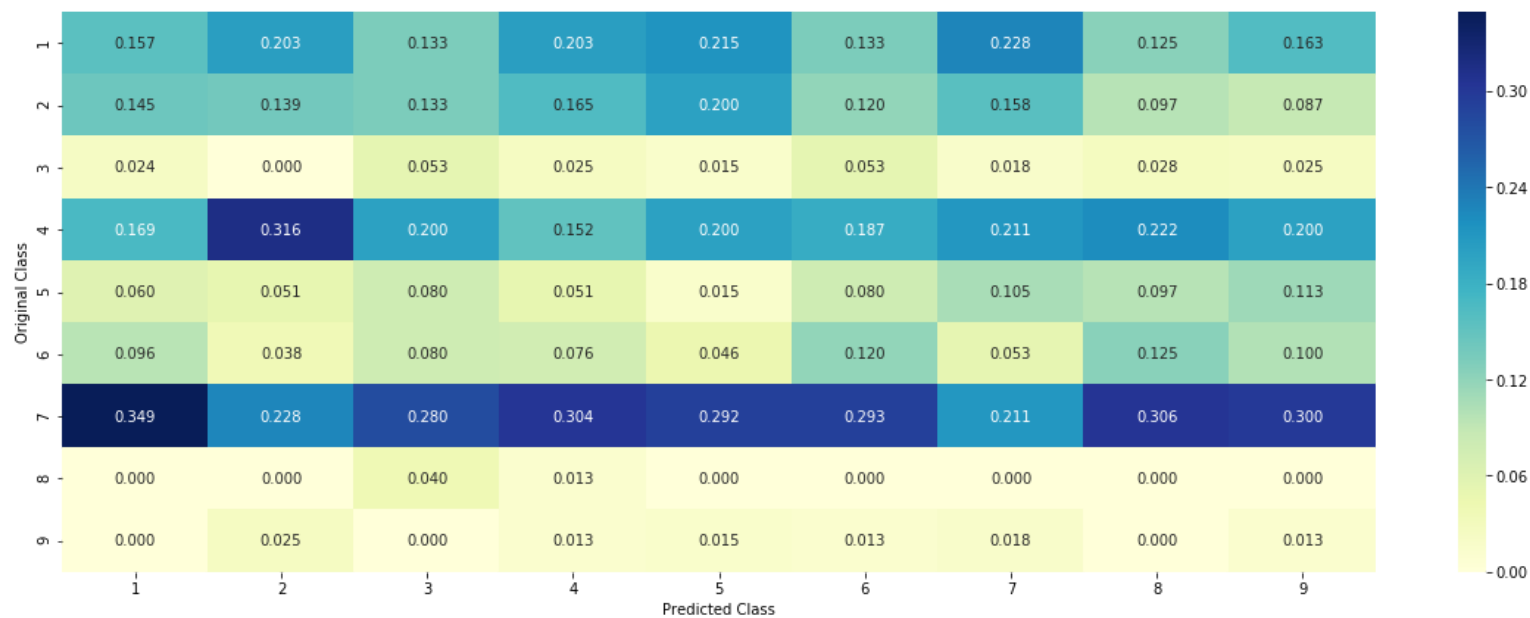
Log loss on Cross Validation Data using Random Model 2.503089399222508

Log loss on Test Data using Random Model 2.4987970579719736

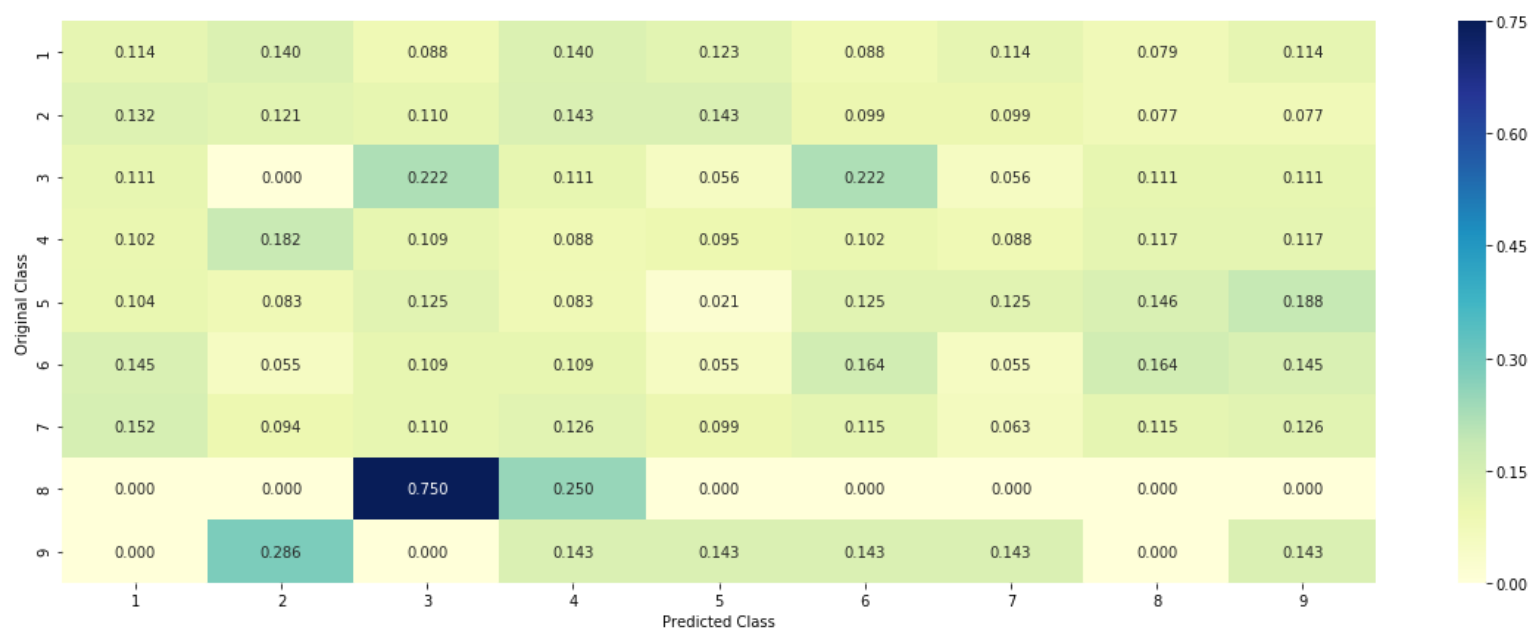
----- Confusion matrix -----



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



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



3.3 Univariate Analysis


```

1  # code for response coding with Laplace smoothing.
2  # alpha : used for laplace smoothing
3  # feature: ['gene', 'variation']
4  # df: ['train_df', 'test_df', 'cv_df']
5  # algorithm
6  # -----
7  # Consider all unique values and the number of occurrences of given feature in train data dataframe
8  # build a vector (1*9) , the first element = (number of times it occurred in class1 + 10*alpha / number of
9  # gv_dict is like a look up table, for every gene it store a (1*9) representation of it
10 # for a value of feature in df:
11 # if it is in train data:
12 # we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
13 # if it is not there is train:
14 # we add [1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
15 # return 'gv_fea'
16 # -----
17
18 # get_gv_fea_dict: Get Gene variation Feature Dict
19 def get_gv_fea_dict(alpha, feature, df):
20     # value_count: it contains a dict like
21     # print(train_df['Gene'].value_counts())
22     # output:
23     #      {BRCA1      174
24     #      TP53      106
25     #      EGFR       86
26     #      BRCA2       75
27     #      PTEN       69
28     #      KIT        61
29     #      BRAF        60
30     #      ERBB2       47
31     #      PDGFRA      46
32     #      ...}
33     # print(train_df['Variation'].value_counts())
34     # output:

```

```

35 # {
36 # Truncating_Mutations          63
37 # Deletion                      43
38 # Amplification                 43
39 # Fusions                       22
40 # Overexpression                3
41 # E17K                          3
42 # Q61L                          3
43 # S222D                         2
44 # P130S                         2
45 # ...
46 # }
47 value_count = train_df[feature].value_counts()
48
49 # gv_dict : Gene Variation Dict, which contains the probability array for each gene/variation
50 gv_dict = dict()
51
52 # denominator will contain the number of time that particular feature occurred in whole data
53 for i, denominator in value_count.items():
54     # vec will contain (p(yi==1/Gi) probability of gene/variation belongs to particular class
55     # vec is 9 dimensional vector
56     vec = []
57     for k in range(1,10):
58         # print(train_df.loc[(train_df['Class']==1) & (train_df['Gene']=='BRCA1')])
59         #          ID  Gene          Variation  Class
60         # 2470  2470  BRCA1          S1715C      1
61         # 2486  2486  BRCA1          S1841R      1
62         # 2614  2614  BRCA1           M1R        1
63         # 2432  2432  BRCA1          L1657P      1
64         # 2567  2567  BRCA1          T1685A      1
65         # 2583  2583  BRCA1          E1660G      1
66         # 2634  2634  BRCA1          W1718L      1
67         # cls_cnt.shape[0] will return the number of rows
68
69         cls_cnt = train_df.loc[(train_df['Class']==k) & (train_df[feature]==i)]

```

```

70
71     # cls_cnt.shape[0](numerator) will contain the number of time that particular feature occurs
72     vec.append((cls_cnt.shape[0] + alpha*10)/ (denominator + 90*alpha))
73
74     # we are adding the gene/variation to the dict as key and vec as value
75     gv_dict[i]=vec
76     return gv_dict
77
78 # Get Gene variation feature
79 def get_gv_feature(alpha, feature, df):
80     # print(gv_dict)
81     #     {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.0681818181818177, 0.13636363636363636],
82     #     'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.2704081632653061],
83     #     'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.0681818181818177, 0.0681818181818177],
84     #     'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.07878787878787878],
85     #     'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.465408805031446],
86     #     'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295],
87     #     'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333333333333334, 0.07333333333333333],
88     #     ...
89     # }
90     gv_dict = get_gv_fea_dict(alpha, feature, df)
91     # value_count is similar in get_gv_fea_dict
92     value_count = train_df[feature].value_counts()
93
94     # gv_fea: Gene_variation feature, it will contain the feature for each feature value in the data
95     gv_fea = []
96     # for every feature values in the given data frame we will check if it is there in the train data table
97     # 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
98     for index, row in df.iterrows():
99         if row[feature] in dict(value_count).keys():
100             gv_fea.append(gv_dict[row[feature]])
101         else:
102             gv_fea.append([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
103     #     gv_fea.append([-1, -1, -1, -1, -1, -1, -1, -1, -1])
104     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?

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

Number of Unique Genes : 227

BRCA1 172

TP53 107

EGFR 92

BRCA2 87

PTEN 84

KIT 68

BRAF 58

ERBB2 46

PDGFRA 40

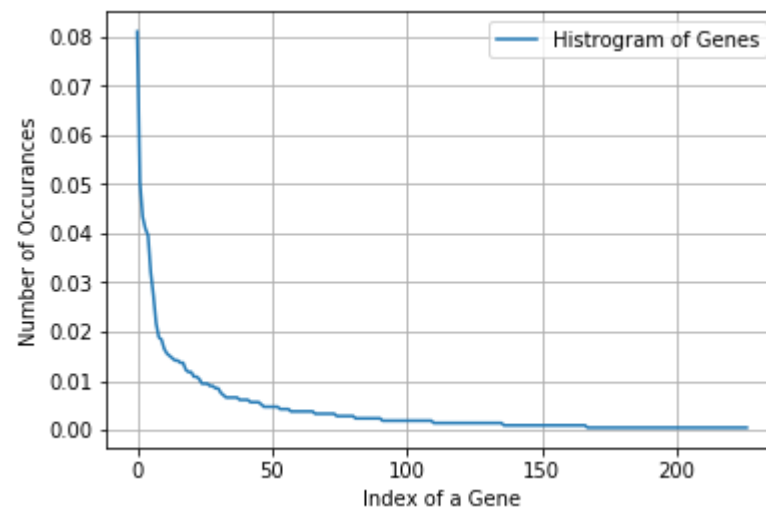
ALK 39

Name: Gene, dtype: int64

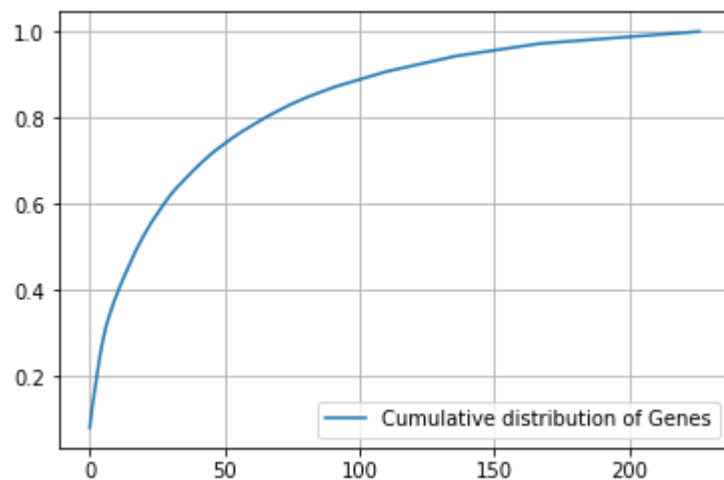
```
1 print("Ans: There are", unique_genes.shape[0], "different categories of genes in the train data, and the
```

Ans: There are 227 different categories of genes in the train data, and they are distributed as follows

```
1 s = sum(unique_genes.values);  
2 h = unique_genes.values/s;  
3 plt.plot(h, label="Histogram of Genes")  
4 plt.xlabel('Index of a Gene')  
5 plt.ylabel('Number of Occurances')  
6 plt.legend()  
7 plt.grid()  
8 plt.show()  
9
```



```
1 c = np.cumsum(h)
2 plt.plot(c,label='Cumulative distribution of Genes')
3 plt.grid()
4 plt.legend()
5 plt.show()
```



Q3. How to featurize this Gene feature ?

Ans.there are two ways we can featurize this variable check out this video:

<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

1. One hot Encoding
2. Response coding

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

```
1 #response-coding of the Gene feature
2 # alpha is used for laplace smoothing
3 alpha = 1
4 # train gene feature
5 train_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", train_df))
6 # test gene feature
7 test_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", test_df))
8 # cross validation gene feature
9 cv_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", cv_df))
```

```
1 print("train_gene_feature_responseCoding is converted feature using response coding method. The shape of
```

```
train_gene_feature_responseCoding is converted feature using response coding method. The shape of gene feature: (2124, 9)
```

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

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

```
2065    SOX9
2327    JAK2
953    PDGFRB
1028    TSC2
1794    AR
Name: Gene, dtype: object
```

```
1 gene_vectorizer.get_feature_names()

['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ar',
 'araf',
 'arid1a',
 'arid2',
 'arid5b',
 'asx11',
 'asx12',
 'atm',
 'atrx',
 'aurka',
 'aurkb',
 'axl',
 'b2m',
```

```
1 print("train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of

train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 226)
```

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 .

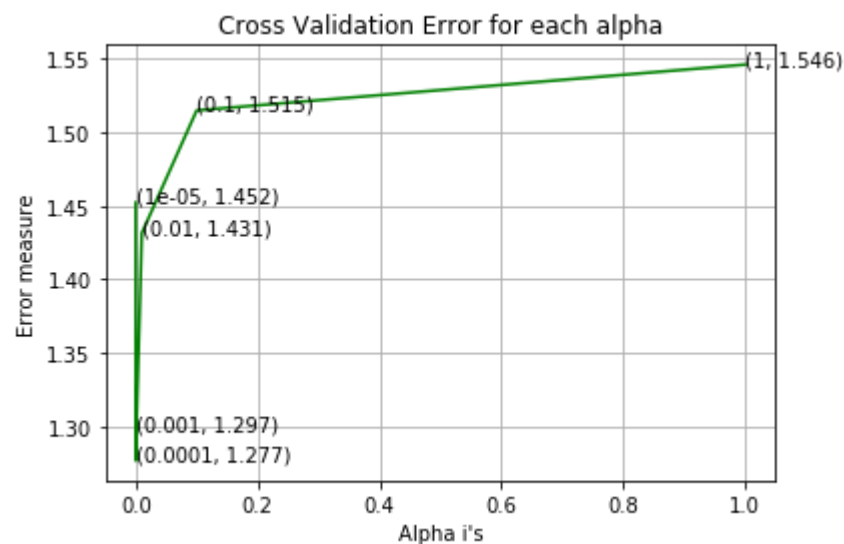

```

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

```

```
35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(cv_log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
42 clf.fit(train_gene_feature_onehotCoding, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(train_gene_feature_onehotCoding, y_train)
45
46 predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
48 predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
49 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_train, predict_y))
50 predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding)
51 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
52
```

```
For values of alpha = 1e-05 The log loss is: 1.4522395212582269
For values of alpha = 0.0001 The log loss is: 1.2765937139984063
For values of alpha = 0.001 The log loss is: 1.2973227647270822
For values of alpha = 0.01 The log loss is: 1.4314445419857231
For values of alpha = 0.1 The log loss is: 1.5149362420219377
For values of alpha = 1 The log loss is: 1.5461036406063464
```



For values of best alpha = 0.0001 The train log loss is: 1.0441095122666129

For values of best alpha = 0.0001 The cross validation log loss is: 1.2765937139984063

For values of best alpha = 0.0001 The test log loss is: 1.223038923188658

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.

```

1 print("Q6. How many data points in Test and CV datasets are covered by the ", unique_genes.shape[0], " {
2
3 test_coverage=test_df[test_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]
4 cv_coverage=cv_df[cv_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]
5
6 print('Ans\n1. In test data',test_coverage, 'out of',test_df.shape[0], ":",(test_coverage/test_df.shape
7 print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0],":", (cv_coverage/cv_df.shape

```

Q6. How many data points in Test and CV datasets are covered by the 227 genes in train dataset?

Ans

1. In test data 638 out of 665 : 95.93984962406014

2. In cross validation data 513 out of 532 : 96.42857142857143

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?

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

Number of Unique Variations : 1932

Deletion 52

Truncating_Mutations 51

Amplification 45

Fusions 22

Overexpression 3

Q61L 3

C618R 2

TPRSS2-ETV1_Fusion 2

G67R 2

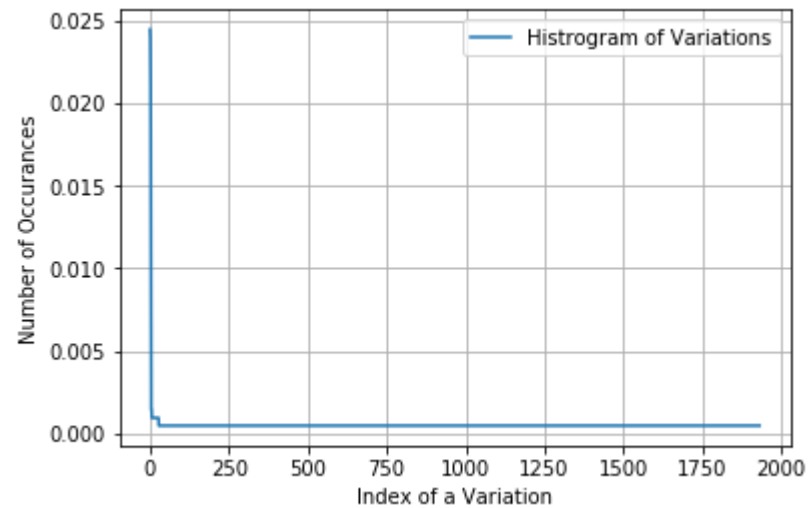
G12C 2

Name: Variation, dtype: int64

```
1 print("Ans: There are", unique_variations.shape[0], "different categories of variations in the train data")
```

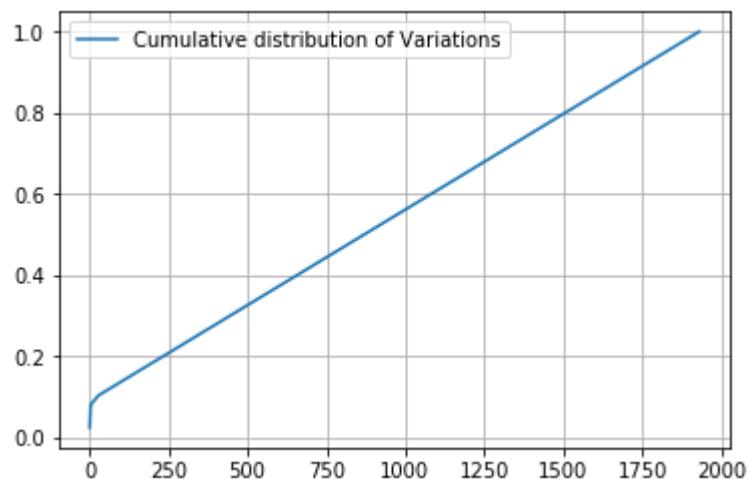
Ans: There are 1932 different categories of variations in the train data, and they are distributed as follows

```
1 s = sum(unique_variations.values);  
2 h = unique_variations.values/s;  
3 plt.plot(h, label="Histogram of Variations")  
4 plt.xlabel('Index of a Variation')  
5 plt.ylabel('Number of Occurances')  
6 plt.legend()  
7 plt.grid()  
8 plt.show()
```



```
1 c = np.cumsum(h)
2 print(c)
3 plt.plot(c,label='Cumulative distribution of Variations')
4 plt.grid()
5 plt.legend()
6 plt.show()
```

```
[0.02448211 0.04849341 0.06967985 ... 0.99905838 0.99952919 1. ]
```



Q9. How to featurize this Variation feature ?

Ans. There are two ways we can featurize this variable check out this video:

<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

1. One hot Encoding
2. Response coding

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

```
1 # alpha is used for Laplace smoothing
2 alpha = 1
3 # train gene feature
4 train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", train_df))
5 # test gene feature
6 test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", test_df))
7 # cross validation gene feature
8 cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", cv_df))
```

```
1 print("train_variation_feature_responseCoding is a converted feature using the response coding method. "
```

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

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

```
1 print("train_variation_feature_onehotEncoded is converted feature using the one-hot encoding method. Tl
```

train_variation_feature_onehotEncoded is converted feature using the one-hot encoding method. The shape of Variation feature: (2124, 1959)

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

Let's build a model just like the earlier!

```

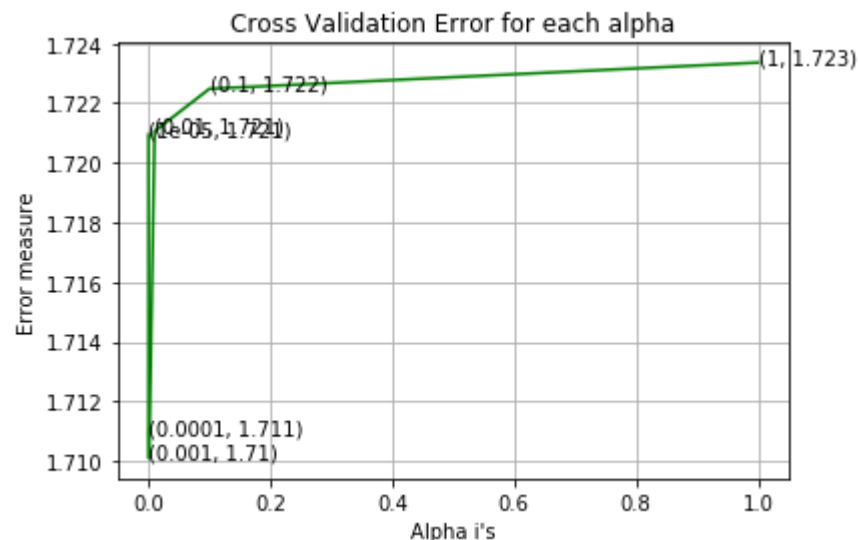
1  alpha = [10 ** x for x in range(-5, 1)]
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,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

```



```
35 plt.grid()
36 plt.title("Cross Validation Error for each alpha")
37 plt.xlabel("Alpha i's")
38 plt.ylabel("Error measure")
39 plt.show()
40
41
42 best_alpha = np.argmin(cv_log_error_array)
43 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
44 clf.fit(train_variation_feature_onehotCoding, y_train)
45 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
46 sig_clf.fit(train_variation_feature_onehotCoding, y_train)
47
48 predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
49 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
50 predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
51 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_train, predict_y))
52 predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding)
53 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
54
```

```
For values of alpha = 1e-05 The log loss is: 1.720943581527042
For values of alpha = 0.0001 The log loss is: 1.7109135510927833
For values of alpha = 0.001 The log loss is: 1.7100790201478915
For values of alpha = 0.01 The log loss is: 1.7210571125209235
For values of alpha = 0.1 The log loss is: 1.7224794460396793
For values of alpha = 1 The log loss is: 1.7233578731600068
```



For values of best alpha = 0.001 The train log loss is: 1.042227051115231

For values of best alpha = 0.001 The cross validation log loss is: 1.7100790201478915

For values of best alpha = 0.001 The test log loss is: 1.7012039305187143

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.

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

Q12. How many data points are covered by total 1932 genes in test and cross validation data sets?

Ans

1. In test data 74 out of 665 : 11.12781954887218

2. In cross validation data 54 out of 532 : 10.150375939849624

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?

```
1  # cls_text is a data frame
2  # for every row in data fram consider the 'TEXT'
3  # split the words by space
4  # make a dict with those words
5  # increment its count whenever we see that word
6
7  def extract_dictionary_paddle(cls_text):
8      dictionary = defaultdict(int)
9      for index, row in cls_text.iterrows():
10         for word in row['TEXT'].split():
11             dictionary[word] +=1
12     return dictionary
```

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

```
1 # building a CountVectorizer with all the words that occurred minimum 3 times in train data
2 text_vectorizer = CountVectorizer()
3 train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])
4 # getting all the feature names (words)
5 train_text_features= text_vectorizer.get_feature_names()
6
7 # train_text_feature_onehotCoding.sum(axis=0).A1 will sum every row and returns (1*number of features)
8 train_text_fea_counts = train_text_feature_onehotCoding.sum(axis=0).A1
9
10 # zip(list(text_features),text_fea_counts) will zip a word with its number of times it occurred
11 text_fea_dict = dict(zip(list(train_text_features),train_text_fea_counts))
12
13
14 print("Total number of unique words in train data :", len(train_text_features))
```

Total number of unique words in train data : 134512

```
1 # Collecting all the genes and variations in a single list
2 corpus = []
3 for word in data['Gene'].values:
4     corpus.append(word)
5 for word in data['Variation'].values:
6     corpus.append(word)
```

```
1  # Training TfidfVectorizer on the 'corpus' list
2  text1 = TfidfVectorizer()
3  text2 = text1.fit_transform(corpus)
4  text1_features = text1.get_feature_names()
5
6  # Transforming the train_df['TEXT']
7  train_text = text1.transform(train_df['TEXT'])
8
9  # Transforming the test_df['TEXT']
10 test_text = text1.transform(test_df['TEXT'])
11
12 # Transforming the cv_df['TEXT']
13 cv_text = text1.transform(cv_df['TEXT'])
14
15 # Normalizing the train_text
16 train_text = normalize(train_text,axis=0)
17
18 # Normalizing the test_text
19 test_text = normalize(test_text,axis=0)
20
21 # Normalizing the cv_text
22 cv_text = normalize(cv_text,axis=0)
```

```
1 dict_list = []
2 # dict_list =[] contains 9 dictionaries each corresponds to a class
3 for i in range(1,10):
4     cls_text = train_df[train_df['Class']==i]
5     # build a word dict based on the words in that class
6     dict_list.append(extract_dictionary_paddle(cls_text))
7     # append it to dict_list
8
9 # dict_list[i] is build on i'th class text data
10 # total_dict is build on whole training text data
11 total_dict = extract_dictionary_paddle(train_df)
12
13
14 confuse_array = []
15 for i in train_text_features:
16     ratios = []
17     max_val = -1
18     for j in range(0,9):
19         ratios.append((dict_list[j][i]+10)/(total_dict[i]+90))
20     confuse_array.append(ratios)
21 confuse_array = np.array(confuse_array)
```

```
1 #response coding of text features
2 train_text_feature_responseCoding = get_text_responsecoding(train_df)
3 test_text_feature_responseCoding = get_text_responsecoding(test_df)
4 cv_text_feature_responseCoding = get_text_responsecoding(cv_df)
```

```
1 # https://stackoverflow.com/a/16202486
2 # we convert each row values such that they sum to 1
3 train_text_feature_responseCoding = (train_text_feature_responseCoding.T/train_text_feature_responseCoding.sum(axis=1)).T
4 test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.sum(axis=1)).T
5 cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=1)).T
```

```
1  # don't forget to normalize every feature
2  train_text_feature_onehotCoding = normalize(train_text_feature_onehotCoding, axis=0)
3
4  # we use the same vectorizer that was trained on train data
5  test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
6  # don't forget to normalize every feature
7  test_text_feature_onehotCoding = normalize(test_text_feature_onehotCoding, axis=0)
8
9  # we use the same vectorizer that was trained on train data
10 cv_text_feature_onehotCoding = text_vectorizer.transform(cv_df['TEXT'])
11 # don't forget to normalize every feature
12 cv_text_feature_onehotCoding = normalize(cv_text_feature_onehotCoding, axis=0)
```

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

```

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

```

```

Counter({1: 46690, 2: 19055, 3: 8600, 4: 7286, 5: 4351, 6: 4332, 8: 2740, 7: 2498, 9: 2351, 10: 2004, 12: 1537, 11: 1486,
15: 1095, 14: 1095, 16: 1040, 13: 928, 18: 892, 17: 738, 20: 706, 19: 623, 21: 557, 22: 540, 48: 536, 24: 526, 30: 483, 2
3: 474, 25: 454, 26: 435, 32: 419, 28: 402, 27: 367, 33: 319, 36: 295, 51: 292, 29: 288, 34: 287, 31: 285, 40: 265, 35: 24
3, 38: 231, 42: 226, 37: 223, 44: 208, 45: 201, 39: 201, 50: 195, 41: 194, 54: 189, 49: 184, 52: 182, 46: 176, 55: 167, 5
7: 161, 43: 161, 64: 151, 60: 147, 56: 143, 53: 141, 47: 139, 58: 138, 63: 125, 62: 120, 72: 115, 65: 114, 61: 114, 59: 11
3, 66: 112, 69: 110, 68: 109, 70: 104, 96: 103, 80: 96, 67: 95, 71: 91, 78: 89, 83: 86, 73: 84, 84: 82, 75: 81, 90: 80, 10
2: 77, 74: 77, 86: 75, 76: 74, 79: 72, 100: 70, 77: 70, 98: 69, 87: 68, 89: 67, 91: 66, 85: 66, 88: 65, 112: 64, 93: 64, 8
1: 64, 99: 61, 92: 60, 144: 59, 103: 58, 82: 58, 108: 56, 113: 53, 105: 53, 94: 53, 107: 50, 114: 49, 128: 48, 109: 48, 9
5: 48, 120: 47, 101: 47, 106: 45, 147: 44, 126: 44, 119: 44, 104: 44, 118: 43, 110: 43, 131: 42, 97: 42, 152: 41, 136: 41,
135: 40, 134: 40, 130: 40, 125: 40, 115: 40, 140: 39, 132: 39, 111: 39, 133: 38, 148: 37, 146: 37, 139: 37, 123: 37, 121:
37, 164: 36, 143: 36, 116: 36, 155: 35, 141: 35, 127: 35, 117: 35, 149: 34, 153: 33, 138: 33, 185: 32, 184: 32, 122: 32, 1
51: 30, 145: 30, 129: 30, 181: 29, 174: 29, 150: 29, 124: 29, 179: 28, 170: 28, 216: 27, 210: 27, 187: 27, 183: 27, 175: 2
7, 162: 27, 171: 26, 160: 26, 156: 26, 142: 26, 223: 25, 206: 25, 200: 25, 182: 25, 177: 25, 167: 25, 178: 24, 165: 24, 15
4: 24, 228: 23, 195: 23, 186: 23, 180: 23, 176: 23, 245: 22, 235: 22, 205: 22, 203: 22, 168: 22, 157: 22, 240: 21, 202: 2
1, 201: 21, 199: 21, 198: 21, 172: 21, 166: 21, 161: 21, 158: 21, 288: 20, 268: 20, 254: 20, 233: 20, 222: 20, 215: 20, 21
1: 20, 192: 20, 191: 20, 188: 20, 274: 19, 250: 19, 249: 19, 217: 19, 214: 19, 212: 19, 209: 19, 207: 19, 159: 19, 277: 1
8, 260: 18, 258: 18, 255: 18, 234: 18, 221: 18, 220: 18, 194: 18, 190: 18, 163: 18, 395: 17, 264: 17, 262: 17, 241: 17, 23
9: 17, 208: 17, 196: 17, 169: 17, 137: 17, 285: 16, 282: 16, 237: 16, 231: 16, 224: 16, 193: 16, 400: 15, 267: 15, 266: 1
5, 252: 15, 246: 15, 244: 15, 242: 15, 230: 15, 225: 15, 219: 15, 189: 15, 173: 15, 371: 14, 298: 14, 294: 14, 291: 14, 27
6: 14, 269: 14, 248: 14, 236: 14, 227: 14, 226: 14, 204: 14, 394: 13, 375: 13, 342: 13, 328: 13, 317: 13, 314: 13, 305: 1
3, 275: 13, 259: 13, 256: 13, 232: 13, 354: 12, 345: 12, 340: 12, 325: 12, 324: 12, 311: 12, 303: 12, 299: 12, 281: 12, 27

```



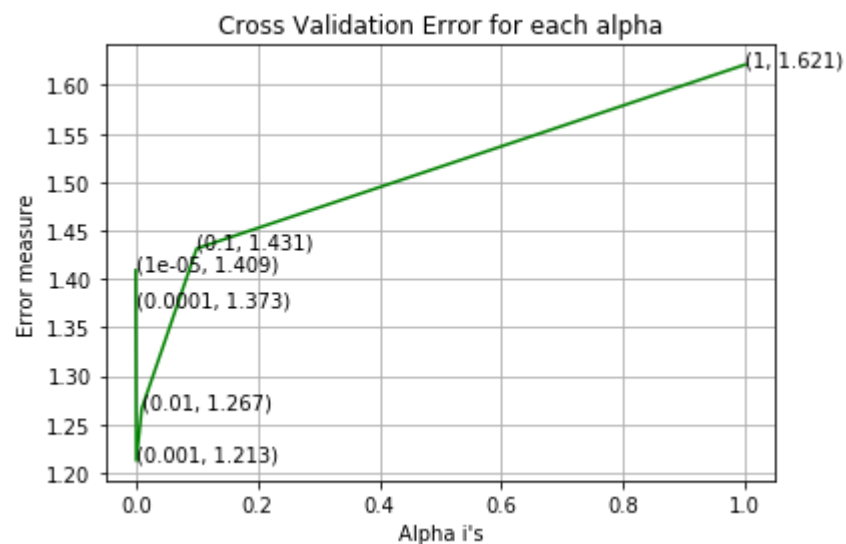
```

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

```

```
35 plt.grid()
36 plt.title("Cross Validation Error for each alpha")
37 plt.xlabel("Alpha i's")
38 plt.ylabel("Error measure")
39 plt.show()
40
41
42 best_alpha = np.argmin(cv_log_error_array)
43 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
44 clf.fit(train_text_feature_onehotCoding, y_train)
45 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
46 sig_clf.fit(train_text_feature_onehotCoding, y_train)
47
48 predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
49 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
50 predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
51 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_train, predict_y))
52 predict_y = sig_clf.predict_proba(test_text_feature_onehotCoding)
53 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
54
```

```
For values of alpha = 1e-05 The log loss is: 1.409032667418594
For values of alpha = 0.0001 The log loss is: 1.373092662164093
For values of alpha = 0.001 The log loss is: 1.212596575674792
For values of alpha = 0.01 The log loss is: 1.2665052982158738
For values of alpha = 0.1 The log loss is: 1.4314377504358973
For values of alpha = 1 The log loss is: 1.6211856455279456
```



For values of best alpha = 0.001 The train log loss is: 0.7909988562287417

For values of best alpha = 0.001 The cross validation log loss is: 1.212596575674792

For values of best alpha = 0.001 The test log loss is: 1.226921695988835

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

Ans. Yes, it seems like!

```

1  def get_intersec_text(df):
2      df_text_vec = CountVectorizer()
3      df_text_fea = df_text_vec.fit_transform(df['TEXT'])
4      df_text_features = df_text_vec.get_feature_names()
5
6      df_text_fea_counts = df_text_fea.sum(axis=0).A1
7      df_text_fea_dict = dict(zip(list(df_text_features), df_text_fea_counts))
8      len1 = len(set(df_text_features))
9      len2 = len(set(train_text_features) & set(df_text_features))
10     return len1, len2

```

```
1 len1,len2 = get_intersec_text(test_df)
2 print(np.round((len2/len1)*100, 3), "% of word of test data appeared in train data")
3 len1,len2 = get_intersec_text(cv_df)
4 print(np.round((len2/len1)*100, 3), "% of word of Cross Validation appeared in train data")
```

79.069 % of word of test data appeared in train data

79.153 % of word of Cross Validation appeared in train data

4. Machine Learning Models

```
1 #Data preparation for ML models.
2
3 #Misc. fonctionns for ML models
4
5
6 def predict_and_plot_confusion_matrix(train_x, train_y, test_x, test_y, clf):
7     clf.fit(train_x, train_y)
8     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
9     sig_clf.fit(train_x, train_y)
10    pred_y = sig_clf.predict(test_x)
11
12    # for calculating log_loss we willl provide the array of probabilities belongs to each class
13    print("Log loss :", log_loss(test_y, sig_clf.predict_proba(test_x)))
14    # calculating the number of data points that are misclassified
15    print("Number of mis-classified points :", np.count_nonzero((pred_y- test_y))/test_y.shape[0])
16    plot_confusion_matrix(test_y, pred_y)
```

```
1 def report_log_loss(train_x, train_y, test_x, test_y, clf):
2     clf.fit(train_x, train_y)
3     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
4     sig_clf.fit(train_x, train_y)
5     sig_clf_probs = sig_clf.predict_proba(test_x)
6     return log_loss(test_y, sig_clf_probs, eps=1e-15)
```

```
1  # this function will be used just for naive bayes
2  # for the given indices, we will print the name of the features
3  # and we will check whether the feature present in the test point text or not
4  def get_impfeature_names(indices, text, gene, var, no_features):
5      gene_count_vec = CountVectorizer()
6      var_count_vec = CountVectorizer()
7      text_count_vec = CountVectorizer()
8
9      gene_vec = gene_count_vec.fit(train_df['Gene'])
10     var_vec = var_count_vec.fit(train_df['Variation'])
11     text_vec = text_count_vec.fit(train_df['TEXT'])
12
13     fea1_len = len(gene_vec.get_feature_names())
14     fea2_len = len(var_count_vec.get_feature_names())
15
16     word_present = 0
17     for i,v in enumerate(indices):
18         if (v < fea1_len):
19             word = gene_vec.get_feature_names()[v]
20             yes_no = True if word == gene else False
21             if yes_no:
22                 word_present += 1
23                 print(i, "Gene feature [{}] present in test data point [{}]" .format(word,yes_no))
24         elif (v < fea1_len+fea2_len):
25             word = var_vec.get_feature_names()[v-(fea1_len)]
26             yes_no = True if word == var else False
27             if yes_no:
28                 word_present += 1
29                 print(i, "variation feature [{}] present in test data point [{}]" .format(word,yes_no))
30         else:
31             word = text_vec.get_feature_names()[v-(fea1_len+fea2_len)]
32             yes_no = True if word in text.split() else False
33             if yes_no:
34                 word_present += 1
```

```

35         print(i, "Text feature [{}] present in test data point [{}].format(word,yes_no))
36
37     print("Out of the top ",no_features," features ", word_present, "are present in query point")

```

Keeping top 15 created variables based on feature importance of Random Forest

```

1  #Keeping separate dataset for newly created feaures
2  train_df_other_feat=train_df[["stopwordst", "numerics", "Text_Word_Count", "Text_Avg_length", "upper".
3  test_df_other_feat=test_df[["stopwordst", "numerics", "Text_Word_Count", "Text_Avg_length", "upper".
4  cv_df_other_feat=cv_df[["stopwordst", "numerics", "Text_Word_Count", "Text_Avg_length", "upper",

```

```

1  train_df_other_feat.head(5)

```

	stopwordst	numerics	Text_Word_Count	Text_Avg_length	upper	Text_Character_Count	var_Avg_leng
2065	2609	84	7461	6.129875	812	45735	4.0
2327	6524	383	22515	6.917877	1968	155756	7.5
953	1468	73	4217	6.782073	191	28600	5.0
1028	1612	79	5158	6.701047	577	34564	6.0
1794	7504	510	25294	6.767573	1711	171179	13.0

```

1  from sklearn.preprocessing import MinMaxScaler
2  sc=MinMaxScaler()
3
4  train_df_other_feat=sc.fit_transform(train_df_other_feat)
5  test_df_other_feat=sc.transform(test_df_other_feat)
6  cv_df_other_feat=sc.transform(cv_df_other_feat)

```

Stacking the three types of features

```
1  # merging gene, variance and text features
2
3  # building train, test and cross validation data sets
4  # a = [[1, 2],
5  #      [3, 4]]
6  # b = [[4, 5],
7  #      [6, 7]]
8  # hstack(a, b) = [[1, 2, 4, 5],
9  #                 [ 3, 4, 6, 7]]
10
11 train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding, train_variation_feature_onehotCoding))
12 test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding, test_variation_feature_onehotCoding))
13 cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding, cv_variation_feature_onehotCoding))
14
15 train_x_onehotCoding = hstack((train_gene_var_onehotCoding, train_text_feature_onehotCoding)).tocsr()
16 train_y = np.array(list(train_df['Class']))
17
18 test_x_onehotCoding = hstack((test_gene_var_onehotCoding, test_text_feature_onehotCoding)).tocsr()
19 test_y = np.array(list(test_df['Class']))
20
21 cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
22 cv_y = np.array(list(cv_df['Class']))
```

```
1 train_x_onehotCoding=hstack((train_x_onehotCoding,train_df_other_feat)).tocsr()
2 test_x_onehotCoding = hstack((test_x_onehotCoding, test_df_other_feat)).tocsr()
3 cv_x_onehotCoding = hstack((cv_x_onehotCoding, cv_df_other_feat)).tocsr()
```

```
1 print("One hot encoding features :")
2 print("(number of data points * number of features) in train data = ", train_x_onehotCoding.shape)
3 print("(number of data points * number of features) in test data = ", test_x_onehotCoding.shape)
4 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, 136712)

(number of data points * number of features) in test data = (665, 136712)

(number of data points * number of features) in cross validation data = (532, 136712)

4.3. Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper paramter tuning


```

1
2 # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_mo
3 # -----
4 # default parameters
5 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=N
6 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
7 # class_weight=None, warm_start=False, average=False, n_iter=None)
8
9 # some of methods
10 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
11 # predict(X) Predict class labels for samples in X.
12
13 #-----
14 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuiti
15 #-----
16
17
18 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklea
19 # -----
20 # default paramters
21 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
22 #
23 # some of the methods of CalibratedClassifierCV()
24 # fit(X, y[, sample_weight]) Fit the calibrated model
25 # get_params([deep]) Get parameters for this estimator.
26 # predict(X) Predict the target of new samples.
27 # predict_proba(X) Posterior probabilities of classification
28 #-----
29 # video link:
30 #-----
31
32 alpha = [10 ** x for x in range(-8, 3)]
33 cv_log_error_array = []
34 for i in alpha:

```

```

35     print("for alpha =", i)
36     clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
37     clf.fit(train_x_onehotCoding, train_y)
38     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
39     sig_clf.fit(train_x_onehotCoding, train_y)
40     sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
41     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
42     # to avoid rounding error while multiplying probabilities we use log-probability estimates
43     print("Log Loss :", log_loss(cv_y, sig_clf_probs))
44
45     fig, ax = plt.subplots()
46     ax.plot(alpha, cv_log_error_array, c='g')
47     for i, txt in enumerate(np.round(cv_log_error_array, 3)):
48         ax.annotate((alpha[i], str(txt)), (alpha[i], cv_log_error_array[i]))
49     plt.grid()
50     plt.title("Cross Validation Error for each alpha")
51     plt.xlabel("Alpha i's")
52     plt.ylabel("Error measure")
53     plt.show()
54
55
56     best_alpha = np.argmin(cv_log_error_array)
57     clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_
58     clf.fit(train_x_onehotCoding, train_y)
59     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
60     sig_clf.fit(train_x_onehotCoding, train_y)
61
62     predict_y = sig_clf.predict_proba(train_x_onehotCoding)
63     print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predi
64     predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
65     print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_
66     predict_y = sig_clf.predict_proba(test_x_onehotCoding)
67     print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predic

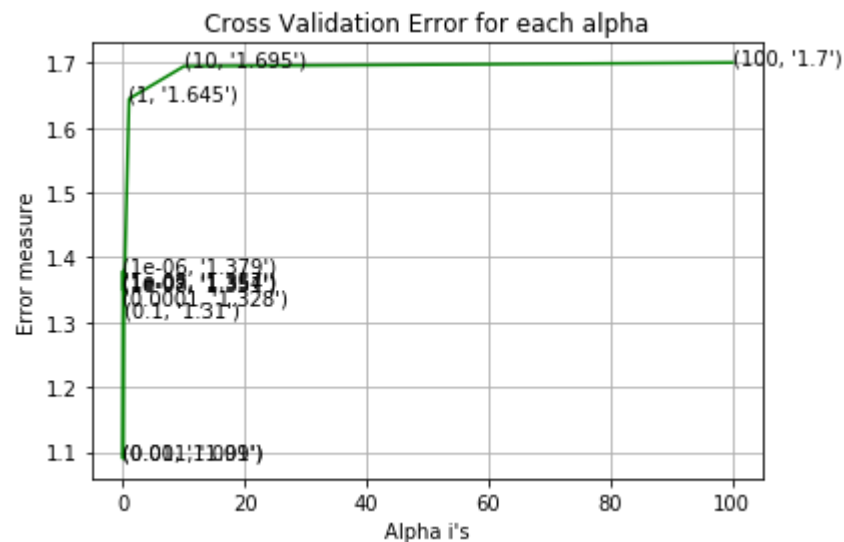
```

for alpha = 1e-08

```

Log Loss : 1.3508842922170041
for alpha = 1e-07
Log Loss : 1.35447047561646
for alpha = 1e-06
Log Loss : 1.3790610180468985
for alpha = 1e-05
Log Loss : 1.3565433625500933
for alpha = 0.0001
Log Loss : 1.3281267828415615
for alpha = 0.001
Log Loss : 1.0897809375901875
for alpha = 0.01
Log Loss : 1.090672560148139
for alpha = 0.1
Log Loss : 1.309603200733344
for alpha = 1
Log Loss : 1.6446133724366971
for alpha = 10
Log Loss : 1.6947507654550469
for alpha = 100
Log Loss : 1.7004145691906267

```



```

For values of best alpha = 0.001 The train log loss is: 0.6604857780958183
For values of best alpha = 0.001 The cross validation log loss is: 1.0897809375901875
For values of best alpha = 0.001 The test log loss is: 1.1094426053084374

```

4.3.1.2. Testing the model with best hyper paramters

```

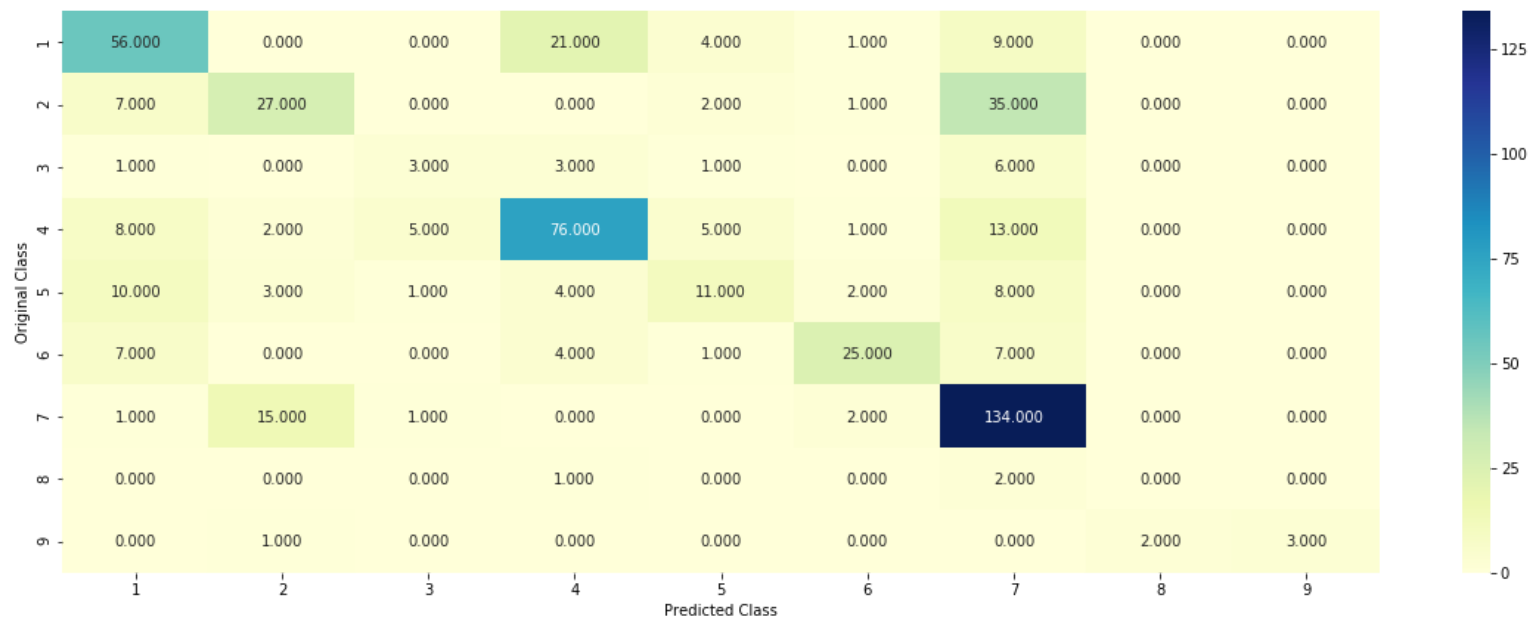
1  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model
2  # -----
3  # default parameters
4  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=1000,
5  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
6  # class_weight=None, warm_start=False, average=False, n_iter=None)
7
8  # some of methods
9  # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
10 # predict(X)    Predict class labels for samples in X.
11
12 #-----
13 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition
14 #-----
15 clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=None)
16 predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

Log loss : 1.0897809375901875

Number of mis-classified points : 0.37030075187969924

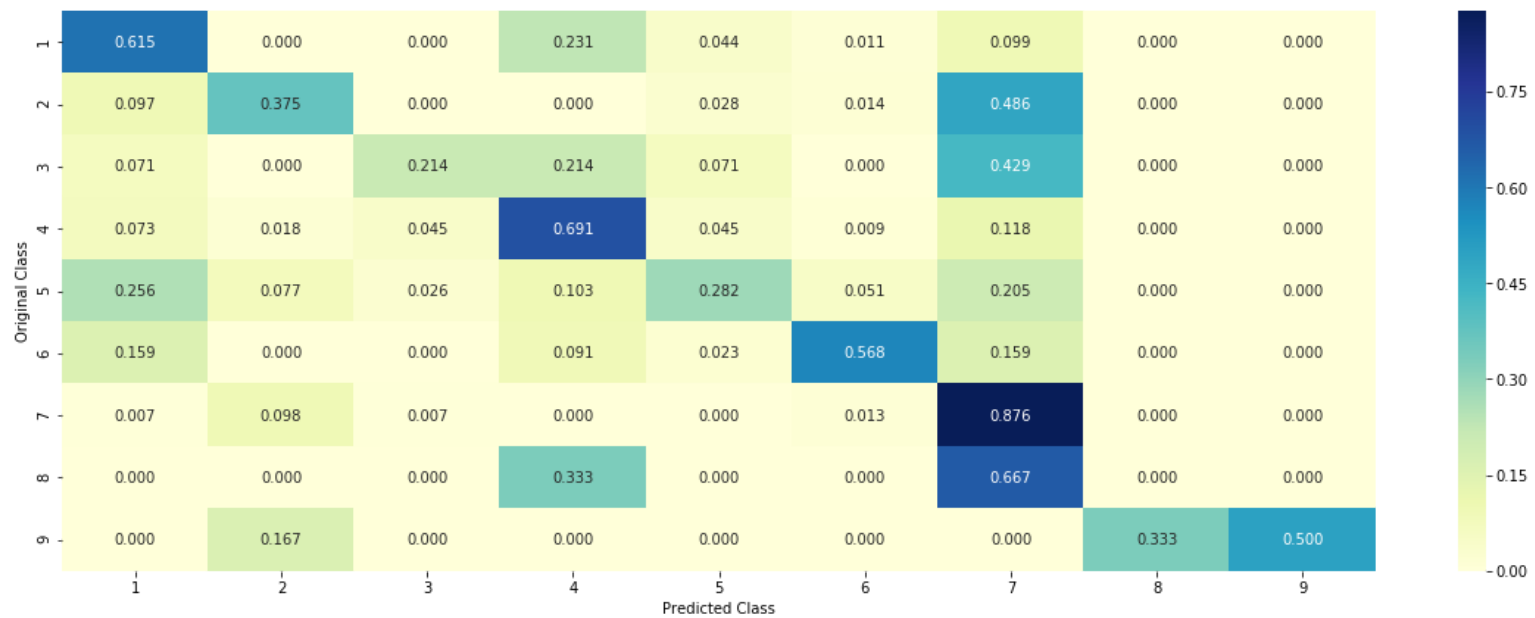
----- Confusion matrix -----



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



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



4.3.1.3. Feature Importance

```
1 def get_imp_feature_names(text, indices, removed_ind = []):
2     word_present = 0
3     tabulte_list = []
4     incresingorder_ind = 0
5     for i in indices:
6         if i < train_gene_feature_onehotCoding.shape[1]:
7             tabulte_list.append([incresingorder_ind, "Gene", "Yes"])
8         elif i < 18:
9             tabulte_list.append([incresingorder_ind, "Variation", "Yes"])
10        if ((i > 17) & (i not in removed_ind)) :
11            word = train_text_features[i]
12            yes_no = True if word in text.split() else False
13            if yes_no:
14                word_present += 1
15                tabulte_list.append([incresingorder_ind, train_text_features[i], yes_no])
16            incresingorder_ind += 1
17        print(word_present, "most important features are present in our query point")
18        print("-"*50)
19        print("The features that are most important of the ", predicted_cls[0], " class:")
20        print (tabulate(tabulte_list, headers=["Index", 'Feature name', 'Present or Not']))
```

4.3.1.3.1. Correctly Classified point


```

1  # from tabulate import tabulate
2  clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_
3  clf.fit(train_x_onehotCoding,train_y)
4  test_point_index = 1
5  no_feature = 500
6  predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
7  print("Predicted Class :", predicted_cls[0])
8  print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_i
9  print("Actual Class :", test_y[test_point_index])
10 indices = np.argsort(-clf.coef_)[predicted_cls-1][:,no_feature]
11 print("-"*50)
12 get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_

```

Predicted Class : 7

Predicted Class Probabilities: [[0.0045 0.1905 0.0012 0.0012 0.0047 0.0014 0.7872 0.0076 0.0017]]

Actual Class : 7

```

-----
23 Text feature [constitutively] present in test data point [True]
39 Text feature [flt1] present in test data point [True]
79 Text feature [oncogene] present in test data point [True]
80 Text feature [oncogenes] present in test data point [True]
84 Text feature [cysteine] present in test data point [True]
89 Text feature [inhibited] present in test data point [True]
137 Text feature [technology] present in test data point [True]
160 Text feature [dramatic] present in test data point [True]
162 Text feature [gaiix] present in test data point [True]
166 Text feature [ligand] present in test data point [True]
177 Text feature [downstream] present in test data point [True]
181 Text feature [concentrations] present in test data point [True]
182 Text feature [thyroid] present in test data point [True]
187 Text feature [expressing] present in test data point [True]
217 Text feature [activating] present in test data point [True]
241 Text feature [cdnas] present in test data point [True]
250 Text feature [manageable] present in test data point [True]
265 Text feature [axilla] present in test data point [True]
302 Text feature [inhibitor] present in test data point [True]
311 Text feature [cot] present in test data point [True]
313 Text feature [viability] present in test data point [True]
334 Text feature [activation] present in test data point [True]

```

```
352 Text feature [forced] present in test data point [True]
368 Text feature [subcutaneous] present in test data point [True]
371 Text feature [melanocyte] present in test data point [True]
376 Text feature [erk1] present in test data point [True]
388 Text feature [hours] present in test data point [True]
446 Text feature [procure] present in test data point [True]
448 Text feature [doses] present in test data point [True]
480 Text feature [mapk] present in test data point [True]
Out of the top 500 features 30 are present in query point
```

4.3.1.3.2. Incorrectly Classified point

```

1 test_point_index = 100
2 no_feature = 500
3 predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
4 print("Predicted Class :", predicted_cls[0])
5 print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]), 4))
6 print("Actual Class :", test_y[test_point_index])
7 indices = np.argsort(-clf.coef_)[predicted_cls-1][:, :no_feature]
8 print("-"*50)
9 get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_index])

```

Predicted Class : 7

Predicted Class Probabilities: [[0.0482 0.2032 0.0108 0.0446 0.071 0.0164 0.5932 0.0078 0.0046]]

Actual Class : 7

```

-----
23 Text feature [constitutively] present in test data point [True]
29 Text feature [constitutive] present in test data point [True]
47 Text feature [activated] present in test data point [True]
79 Text feature [oncogene] present in test data point [True]
89 Text feature [inhibited] present in test data point [True]
93 Text feature [transforming] present in test data point [True]
108 Text feature [transform] present in test data point [True]
148 Text feature [receptors] present in test data point [True]
177 Text feature [downstream] present in test data point [True]
210 Text feature [isozyme] present in test data point [True]
217 Text feature [activating] present in test data point [True]
232 Text feature [exchange] present in test data point [True]
326 Text feature [murine] present in test data point [True]
333 Text feature [agar] present in test data point [True]
334 Text feature [activation] present in test data point [True]
Out of the top 500 features 15 are present in query point

```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

```

1  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_mo
2  # -----
3  # default parameters
4  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=N
5  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
6  # class_weight=None, warm_start=False, average=False, n_iter=None)
7
8  # some of methods
9  # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
10 # predict(X)    Predict class labels for samples in X.
11
12 #-----
13 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuiti
14 #-----
15
16
17
18 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklea
19 # -----
20 # default paramters
21 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
22 #
23 # some of the methods of CalibratedClassifierCV()
24 # fit(X, y[, sample_weight])    Fit the calibrated model
25 # get_params([deep])    Get parameters for this estimator.
26 # predict(X)    Predict the target of new samples.
27 # predict_proba(X)    Posterior probabilities of classification
28 #-----
29 # video link:
30 #-----
31
32 alpha = [10 ** x for x in range(-8, 3)]
33 cv_log_error_array = []
34 for i in alpha:

```

```

35     print("for alpha =", i)
36     clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
37     clf.fit(train_x_onehotCoding, train_y)
38     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
39     sig_clf.fit(train_x_onehotCoding, train_y)
40     sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
41     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
42     print("Log Loss :", log_loss(cv_y, sig_clf_probs))
43
44     fig, ax = plt.subplots()
45     ax.plot(alpha, cv_log_error_array, c='g')
46     for i, txt in enumerate(np.round(cv_log_error_array, 3)):
47         ax.annotate((alpha[i], str(txt)), (alpha[i], cv_log_error_array[i]))
48     plt.grid()
49     plt.title("Cross Validation Error for each alpha")
50     plt.xlabel("Alpha i's")
51     plt.ylabel("Error measure")
52     plt.show()
53
54
55     best_alpha = np.argmin(cv_log_error_array)
56     clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
57     clf.fit(train_x_onehotCoding, train_y)
58     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
59     sig_clf.fit(train_x_onehotCoding, train_y)
60
61     predict_y = sig_clf.predict_proba(train_x_onehotCoding)
62     print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
63     predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
64     print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y))
65     predict_y = sig_clf.predict_proba(test_x_onehotCoding)
66     print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))

```

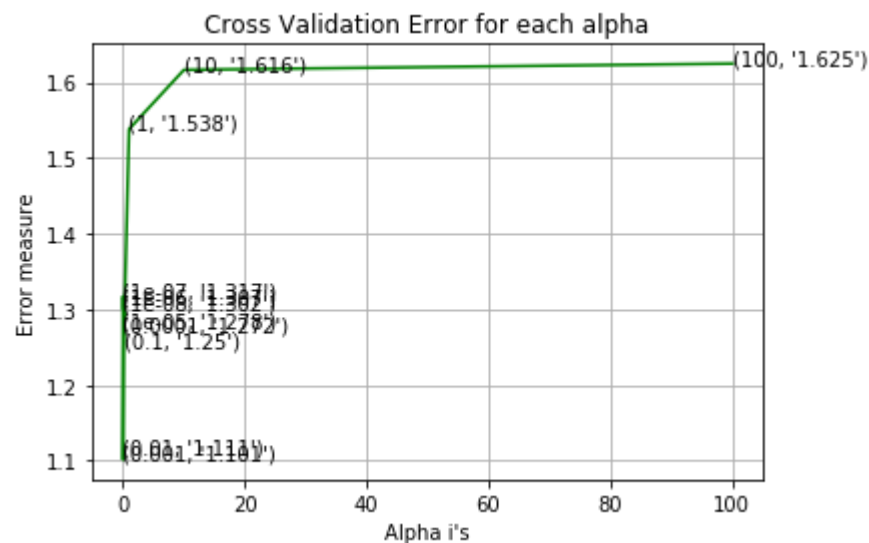
for alpha = 1e-08

Log Loss : 1.301898440354458

```

for alpha = 1e-07
Log Loss : 1.3171176156919957
for alpha = 1e-06
Log Loss : 1.3074648254683952
for alpha = 1e-05
Log Loss : 1.2778335725673895
for alpha = 0.0001
Log Loss : 1.2718700994639265
for alpha = 0.001
Log Loss : 1.1014973044058862
for alpha = 0.01
Log Loss : 1.1110127984170965
for alpha = 0.1
Log Loss : 1.2501172885189447
for alpha = 1
Log Loss : 1.5380488780721915
for alpha = 10
Log Loss : 1.6162602554168248
for alpha = 100
Log Loss : 1.6245371194726435

```



```

For values of best alpha = 0.001 The train log loss is: 0.648003444102412
For values of best alpha = 0.001 The cross validation log loss is: 1.1014973044058862
For values of best alpha = 0.001 The test log loss is: 1.1079192403984552

```

4.3.2.2. Testing model with best hyper parameters

```

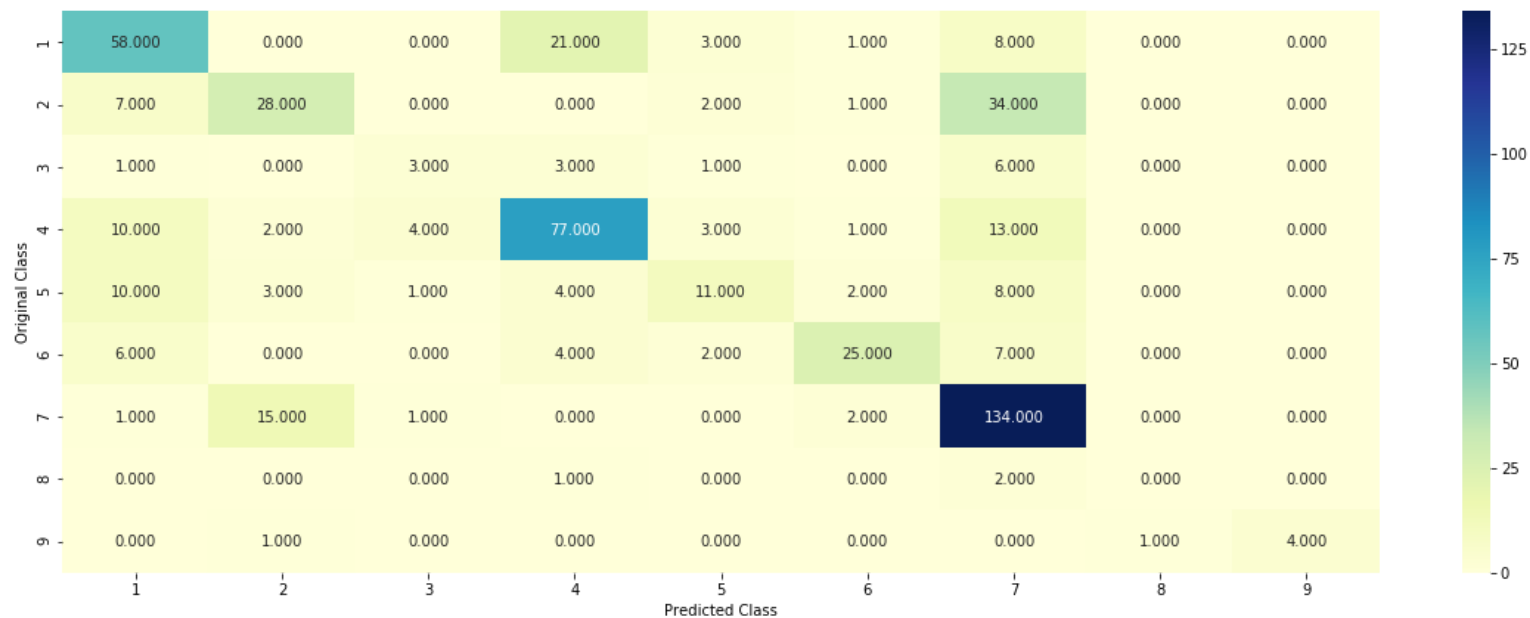
1  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model
2  # -----
3  # default parameters
4  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=N
5  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0,
6  # class_weight=None, warm_start=False, average=False, n_iter=None)
7
8  # some of methods
9  # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
10 # predict(X)    Predict class labels for samples in X.
11
12 #-----
13 # video link:
14 #-----
15
16 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
17 predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

Log loss : 1.1014973044058862

Number of mis-classified points : 0.3609022556390977

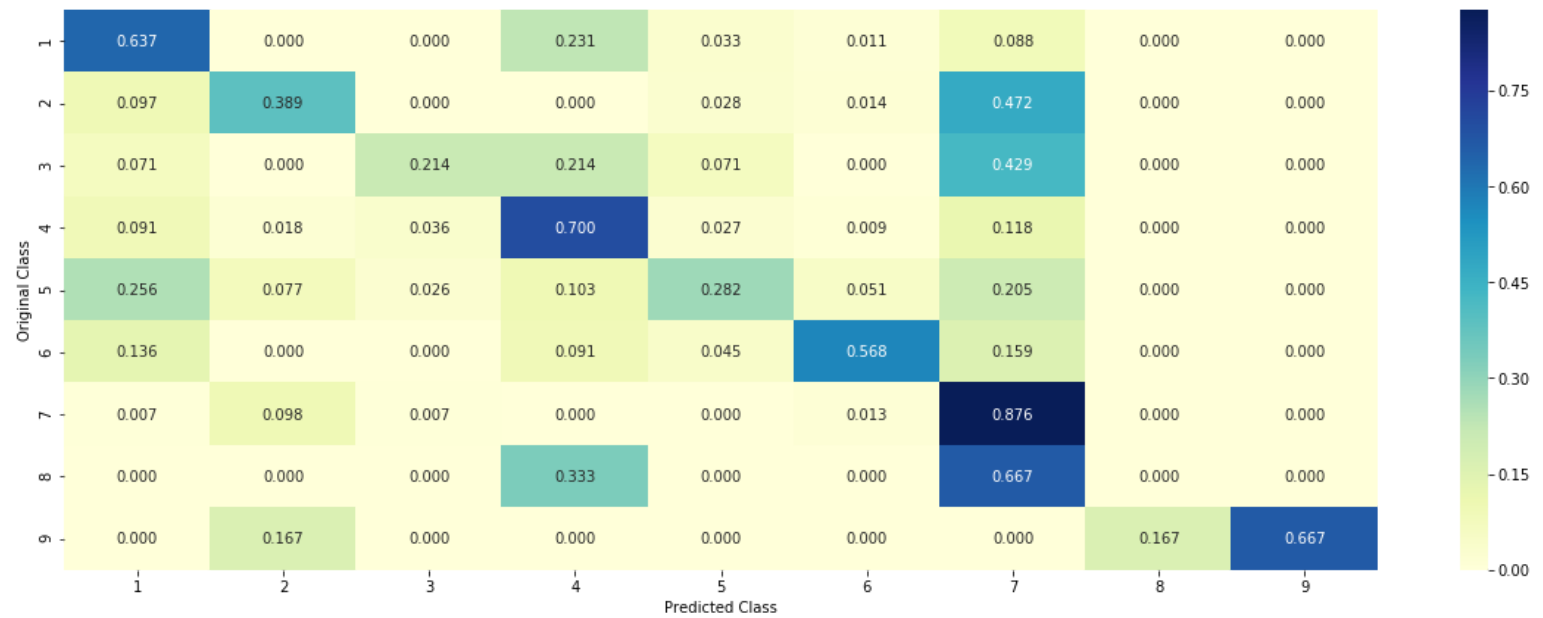
----- Confusion matrix -----



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



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



4.3.2.3. Feature Importance, Correctly Classified point

```

1  clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
2  clf.fit(train_x_onehotCoding,train_y)
3  test_point_index = 1
4  no_feature = 500
5  predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
6  print("Predicted Class :", predicted_cls[0])
7  print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_i
8  print("Actual Class :", test_y[test_point_index])
9  indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
10 print("-"*50)
11 get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_

```

Predicted Class : 7

Predicted Class Probabilities: [[5.100e-03 1.255e-01 2.000e-04 1.300e-03 2.300e-03 1.400e-03 8.556e-01
8.500e-03 1.000e-04]]

Actual Class : 7

```

-----
60 Text feature [constitutively] present in test data point [True]
107 Text feature [flt1] present in test data point [True]
124 Text feature [cysteine] present in test data point [True]
157 Text feature [oncogenes] present in test data point [True]
158 Text feature [inhibited] present in test data point [True]
195 Text feature [activating] present in test data point [True]
200 Text feature [ligand] present in test data point [True]
203 Text feature [oncogene] present in test data point [True]
204 Text feature [technology] present in test data point [True]
257 Text feature [gaiix] present in test data point [True]
260 Text feature [concentrations] present in test data point [True]
265 Text feature [downstream] present in test data point [True]
314 Text feature [hki] present in test data point [True]
316 Text feature [dramatic] present in test data point [True]
323 Text feature [expressing] present in test data point [True]
371 Text feature [cdnas] present in test data point [True]
380 Text feature [viability] present in test data point [True]
412 Text feature [thyroid] present in test data point [True]
459 Text feature [activation] present in test data point [True]
461 Text feature [manageable] present in test data point [True]
462 Text feature [ser473] present in test data point [True]
468 Text feature [axilla] present in test data point [True]

```

495 Text feature [extracellular] present in test data point [True]
 Out of the top 500 features 23 are present in query point

4.3.2.4. Feature Importance, Inorrectly Classified point

```

1 test_point_index = 100
2 no_feature = 500
3 predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
4 print("Predicted Class :", predicted_cls[0])
5 print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_i
6 print("Actual Class :", test_y[test_point_index])
7 indices = np.argsort(-clf.coef_)[predicted_cls-1][:, :no_feature]
8 print("-"*50)
9 get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_

```

Predicted Class : 7

Predicted Class Probabilities: [[0.0485 0.1851 0.0052 0.0442 0.0617 0.0143 0.6317 0.0072 0.0022]]

Actual Class : 7

```

-----
60 Text feature [constitutively] present in test data point [True]
89 Text feature [constitutive] present in test data point [True]
116 Text feature [activated] present in test data point [True]
158 Text feature [inhibited] present in test data point [True]
159 Text feature [transforming] present in test data point [True]
193 Text feature [receptors] present in test data point [True]
195 Text feature [activating] present in test data point [True]
203 Text feature [oncogene] present in test data point [True]
226 Text feature [transform] present in test data point [True]
241 Text feature [isozyme] present in test data point [True]
265 Text feature [downstream] present in test data point [True]
377 Text feature [agar] present in test data point [True]
442 Text feature [interatomic] present in test data point [True]
459 Text feature [activation] present in test data point [True]
Out of the top 500 features 14 are present in query point

```

Summary

```

1 from prettytable import PrettyTable
2 x=PrettyTable()
3 x.field_names = ["Model", "Train loss", "CV loss", "Test Loss", "Mis-classified pts"]
4 x.add_row(["LR(class balancing) with one hot encoding", 0.6604, 1.0897, 1.1094, 0.3703])
5 x.add_row(["LR(without class balancing) with one hot encoding", 0.6480, 1.1014, 1.1079, 0.3609])
6
7 print(x)

```

Model	Train loss	CV loss	Test Loss	Mis-classified pts
LR(class balancing) with one hot encoding	0.6604	1.0897	1.1094	0.3703
LR(without class balancing) with one hot encoding	0.648	1.1014	1.1079	0.3609

In case of BOW LR(without class balancing) performed slightly better.