

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

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

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

2.2.1. Type of Machine Learning Problem

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

2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

2.3. Train, CV and Test Datasets

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

3. Exploratory Data Analysis

```
In [4]: 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
20 from collections import Counter
21 from scipy.sparse import hstack
22 from sklearn.multiclass import OneVsRestClassifier
23 from sklearn.svm import SVC
24 from sklearn.model_selection import StratifiedKFold
25 from collections import Counter, defaultdict
26 from sklearn.calibration import CalibratedClassifierCV
27 from sklearn.naive_bayes import MultinomialNB
28 from sklearn.naive_bayes import GaussianNB
29 from sklearn.model_selection import train_test_split
30 from sklearn.model_selection import GridSearchCV
31 import math
32 from sklearn.metrics import normalized_mutual_info_score
33 from sklearn.ensemble import RandomForestClassifier
34 warnings.filterwarnings("ignore")
35
36 from mlxtend.classifier import StackingClassifier
37
38 from sklearn import model_selection
39 from sklearn.linear_model import LogisticRegression
40
```

```
In [128]: 1 # import dill
          2 # dill.load_session('notebook_env.db')
```

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

```
In [5]: 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']
```

Out[5]:

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

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

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

3.1.2. Reading Text Data

```
In [6]: 1 # note the separator in this file
2 data_text = pd.read_csv("training/training_text", sep="\|", engine="python", names=["ID", "TEXT"], skiprows=1)
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']

Out[6]:

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

3.1.3. Preprocessing of text

```
In [7]: 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
```

```
In [8]: 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 : 227.813003 seconds
```



```
In [10]: 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()
```

Out[10]:

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

```
In [81]: 1 result["TEXT"]
```

```
Out[81]: 0      cyclin dependent kinases cdks regulate variety...
1      abstract background non small cell lung cancer...
2      abstract background non small cell lung cancer...
3      recent evidence demonstrated acquired uniparen...
4      oncogenic mutations monomeric casitas b lineag...
5      oncogenic mutations monomeric casitas b lineag...
6      oncogenic mutations monomeric casitas b lineag...
7      cbl negative regulator activated receptor tyro...
8      abstract juvenile myelomonocytic leukemia jmml...
9      abstract juvenile myelomonocytic leukemia jmml...
10     oncogenic mutations monomeric casitas b lineag...
11     noonan syndrome autosomal dominant congenital ...
12     noonan syndrome autosomal dominant congenital ...
13     noonan syndrome autosomal dominant congenital ...
14     oncogenic mutations monomeric casitas b lineag...
15     noonan syndrome autosomal dominant congenital ...
16     determine residual cylindrical refractive erro...
17     acquired uniparental disomy aupd common featur...
18     oncogenic mutations monomeric casitas b lineag...
19     acquired uniparental disomy aupd common featur...
20     abstract background non small cell lung cancer...
21     oncogenic mutations monomeric casitas b lineag...
22     oncogenic mutations monomeric casitas b lineag...
23     recent evidence demonstrated acquired uniparen...
24     recent evidence demonstrated acquired uniparen...
25     recent evidence demonstrated acquired uniparen...
26     abstract n myristoylation common form co trans...
27     heterozygous mutations telomerase components t...
28     sequencing studies identified many recurrent c...
29     heterozygous mutations telomerase components t...
...
3291    investigatedthe transformingactivityofthe ret ...
3292    investigatedthe transformingactivityofthe ret ...
3293    ret transmembrane tyrosine kinase participatin...
3294    introduction inherited germ line activating mu...
3295    many missense mutations ret proto oncogene fou...
3296    ret proto oncogene encodes receptor tyrosine k...
3297    aml1 gene known frequent target chromosomal tr...
3298    introduction myelodysplastic syndromes mds het...
```

```

3299 bcr abl fusion protein generated 9 22 q34 q11 ...
3300 frequent mutations associated leukemia recurre...
3301 frequent mutations associated leukemia recurre...
3302 familial platelet disorder predisposition acut...
3303 introduction myelodysplastic syndromes mds het...
3304 familial platelet disorder predisposition acut...
3305 introduction myelodysplastic syndromes mds het...
3306 report two new runx1 mutations one patient con...
3307 runx genes come prominence recently roles esse...
3308 familial platelet disorder propensity acute my...
3309 bcr abl fusion protein generated 9 22 q34 q11 ...
3310 runx proteins belong family metazoan transcrip...
3311 aml1 evi 1 chimeric gene generated 3 21 q26 q2...
3312 balanced chromosomal translocations frequently...
3313 bcr abl fusion protein generated 9 22 q34 q11 ...
3314 introduction myelodysplastic syndromes mds het...
3315 runx gene family includes three evolutionarily...
3316 introduction myelodysplastic syndromes mds het...
3317 introduction myelodysplastic syndromes mds het...
3318 runt related transcription factor 1 gene runx1...
3319 runx1 aml1 gene frequent target chromosomal tr...
3320 frequent mutations associated leukemia recurre...
Name: TEXT, Length: 3321, dtype: object

```

```
In [11]: 1 result[result.isnull().any(axis=1)]
```

Out[11]:

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

```
In [12]: 1 result.loc[result['TEXT'].isnull(), 'TEXT'] = result['Gene'] + ' '+result['Variation']
```

```
In [13]: 1 result[result['ID']==1109]
```

```
Out[13]:
```

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

3.1.4. Test, Train and Cross Validation Split

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

```
In [15]: 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' [stratify
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

```
In [16]: 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

```

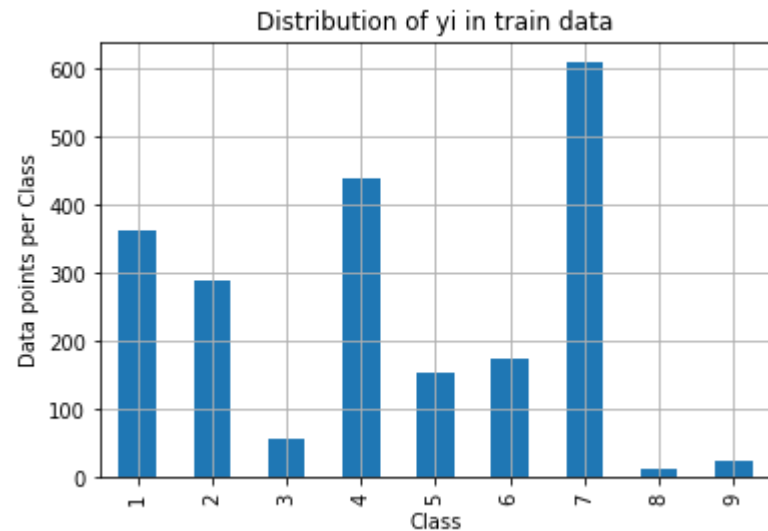
In [17]: 1 # it returns a dict, keys as class labels and values as the number of data points in that class
2 train_class_distribution = train_df['Class'].value_counts().sort_index()
3 test_class_distribution = test_df['Class'].value_counts().sort_index()
4 cv_class_distribution = cv_df['Class'].value_counts().sort_index()
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((tra
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((test
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_cla
50

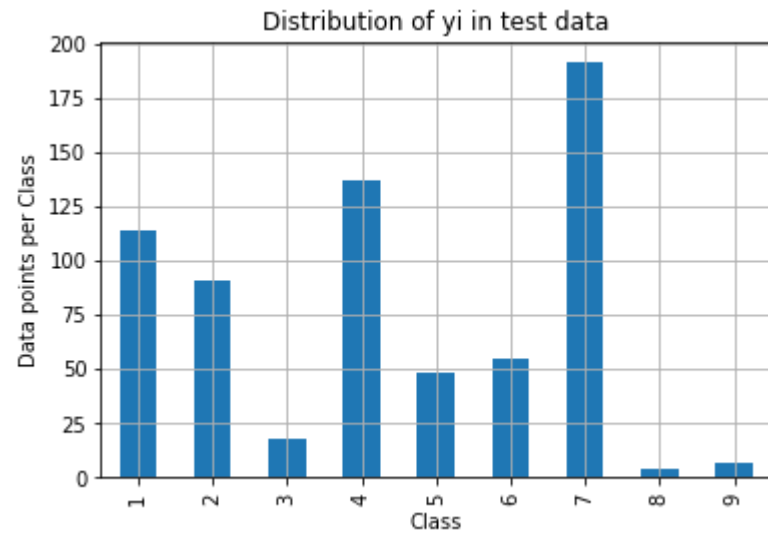
```



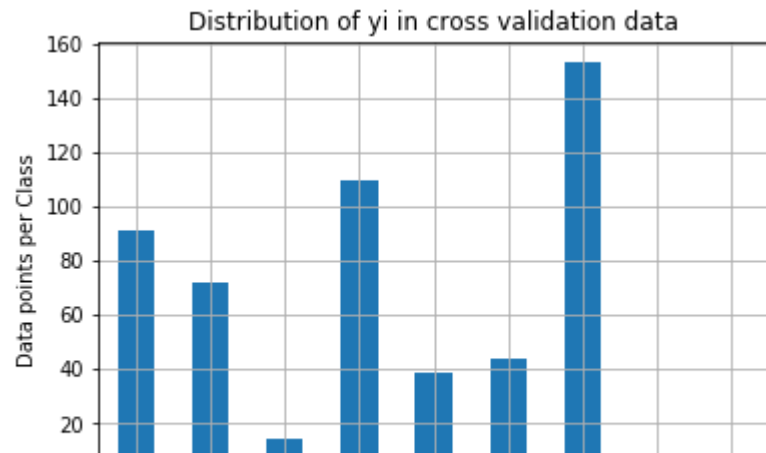
```

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.


```

In [18]: 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 corresponds to columns and axis=1 corresponds to rows in two dimensional array
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 corresponds to columns and axis=1 corresponds to rows in two dimensional array
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()
```

```

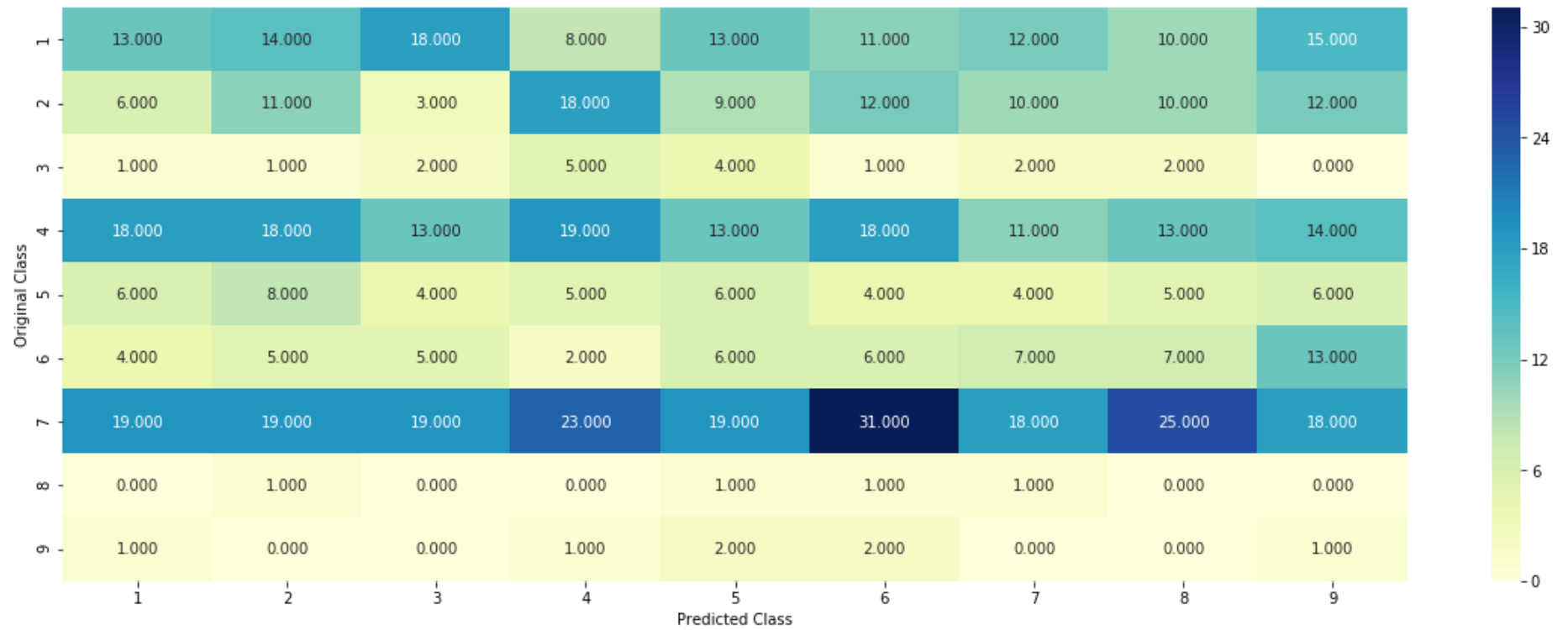
In [19]: 1 # we need to generate 9 numbers and the sum of numbers should be 1
2 # one solution is to generate 9 numbers and divide each of the numbers by their sum
3 # ref: https://stackoverflow.com/a/18662466/4084039
4 test_data_len = test_df.shape[0]
5 cv_data_len = cv_df.shape[0]
6
7 # we create a 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)

```

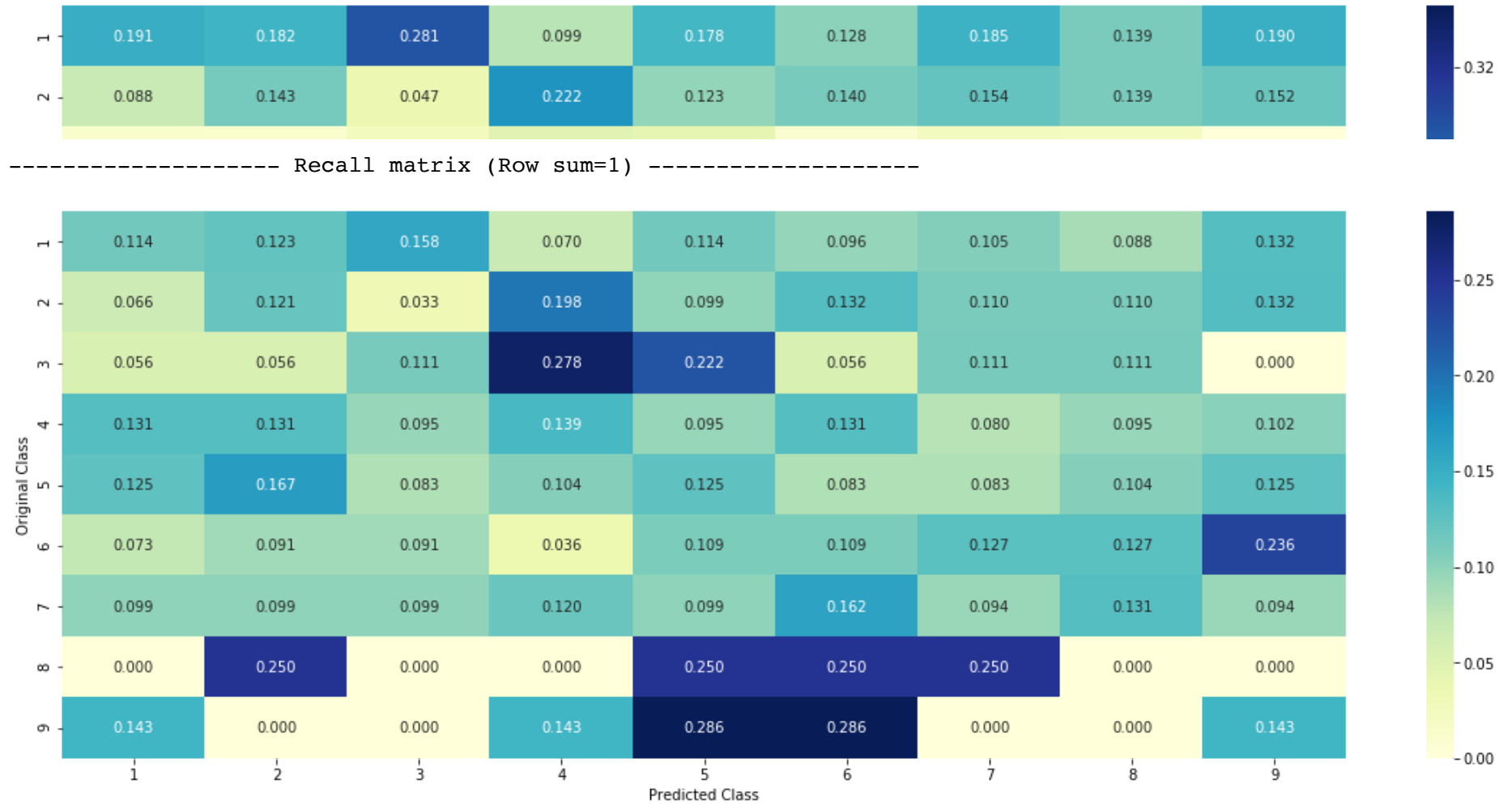
Log loss on Cross Validation Data using Random Model 2.4695406766150847

Log loss on Test Data using Random Model 2.475479077789838

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



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



3.3 Univariate Analysis

```

In [20]: 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 occurances 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 t
          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             #
60             # ID      Gene      Variation      Class
61             # 2470    2470    BRCA1          S1715C      1
62             # 2486    2486    BRCA1          S1841R      1
63             # 2614    2614    BRCA1          M1R        1
64             # 2432    2432    BRCA1          L1657P      1
65             # 2567    2567    BRCA1          T1685A      1
66             # 2583    2583    BRCA1          E1660G      1
67             # 2634    2634    BRCA1          W1718L      1
68             # cls_cnt.shape[0] will return the number of rows
69
70             cls_cnt = train_df.loc[(train_df['Class']==k) & (train_df[feature]==i)]
71
72             # cls_cnt.shape[0](numerator) will contain the number of time that particular feature occurred
73             vec.append((cls_cnt.shape[0] + alpha*10)/ (denominator + 90*alpha))
74
75             # we are adding the gene/variation to the dict as key and vec as value
76             gv_dict[i]=vec
77         return gv_dict
78
79 # Get Gene variation feature
80 def get_gv_feature(alpha, feature, df):
81     # print(gv_dict)
82     # {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.06818181818181817, 0.13636363636363635, 0.06818181818181818, 0.06818181818181818, 0.06818181818181818, 0.06818181818181818, 0.06818181818181818],
83     # 'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.27040816326530615, 0.061224489795918366, 0.061224489795918366, 0.061224489795918366, 0.061224489795918366, 0.061224489795918366],
84     # 'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.06818181818181817, 0.06818181818181818, 0.06818181818181818, 0.06818181818181818, 0.06818181818181818, 0.06818181818181818]}

```

```

84     # 'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.078787878787878782,
85     # 'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.46540880503144655,
86     # 'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295,
87     # 'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333333333333334, 0.073333333333333334,
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 then
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 caculate the probability of a feature belongs to any particular class, we apply laplace smoothing

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

3.2.1 Univariate Analysis on Gene Feature

Q1. Gene, What type of feature it is ?

Ans. Gene is a categorical variable

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


```
In [21]: 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 : 236

BRCA1 177

TP53 106

EGFR 94

BRCA2 81

PTEN 78

BRAF 67

KIT 52

ALK 48

ERBB2 38

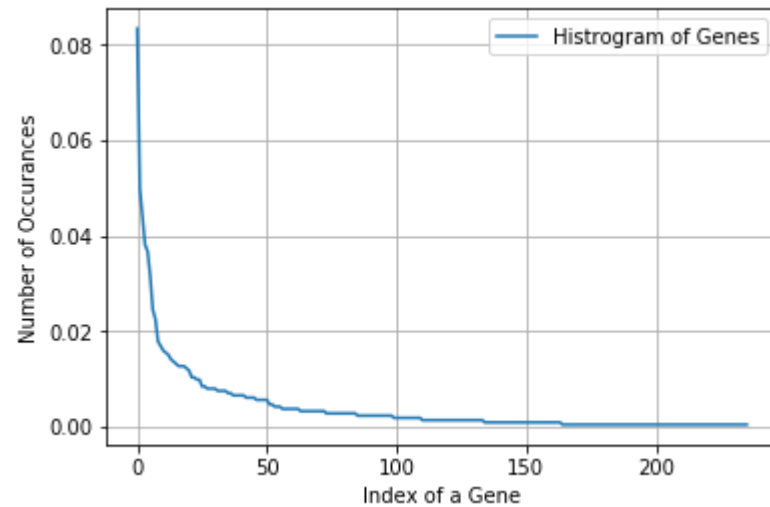
PIK3CA 36

Name: Gene, dtype: int64

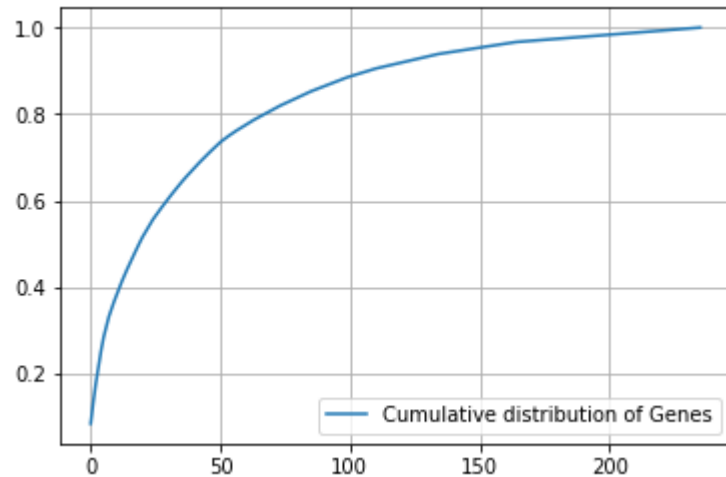
```
In [22]: 1 print("Ans: There are", unique_genes.shape[0] , "different categories of genes in the train data, and they a
```

Ans: There are 236 different categories of genes in the train data, and they are distributed as follows

```
In [23]: 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
```



```
In [24]: 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.

```
In [25]: 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))
```

```
In [26]: 1 print("train_gene_feature_responseCoding is converted feature using response coding method. The shape of gene feature is: (2124, 9)")

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

```
In [27]: 1 # TFIDF
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 gene_vectorizer = TfidfVectorizer()
4 train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])
5 test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])
6 cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

```
In [28]: 1 train_df['Gene'].head()
```

```
Out[28]: 116      KMT2D
2513     BRCA1
2693     BRAF
578     SMAD4
3009     KIT
Name: Gene, dtype: object
```

```
In [29]: 1 gene_vectorizer.get_feature_names()
```

```
Out[29]: ['abl1',  
          'acvr1',  
          'ago2',  
          'akt1',  
          'akt2',  
          'akt3',  
          'alk',  
          'apc',  
          'ar',  
          'araf',  
          'arid1a',  
          'arid1b',  
          'arid2',  
          'arid5b',  
          'asx11',  
          'atm',  
          'atr',  
          'atrx',  
          'aurka',  
          ...]
```

```
In [30]: 1 print("train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene
```

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature:  
e: (2124, 235)
```

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

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

```

In [31]: 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.
4 # -----
5 # default parameters
6 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None,
7 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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, n_jobs=3)
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_, ep
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, n_jobs=3)

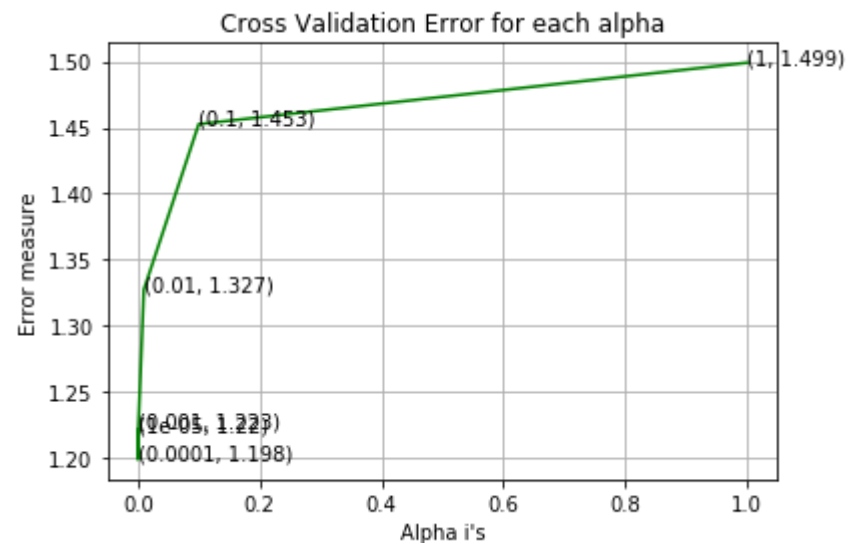
```

```

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_cv, 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: 1.2202789742572298
 For values of alpha = 0.0001 The log loss is: 1.1982950255130638
 For values of alpha = 0.001 The log loss is: 1.222517090954132
 For values of alpha = 0.01 The log loss is: 1.3269440830250716
 For values of alpha = 0.1 The log loss is: 1.4525111676082707
 For values of alpha = 1 The log loss is: 1.499096866060869



For values of best alpha = 0.0001 The train log loss is: 1.007164617292122
 For values of best alpha = 0.0001 The cross validation log loss is: 1.1982950255130638
 For values of best alpha = 0.0001 The test log loss is: 1.1654562649337685

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

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

```
In [33]: 1 print("Q6. How many data points in Test and CV datasets are covered by the ", unique_genes.shape[0], " gene
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[0])
7 print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0],":", (cv_coverage/cv_df.shape[0])
```

Q6. How many data points in Test and CV datasets are covered by the 236 genes in train dataset?

Ans

1. In test data 651 out of 665 : 97.89473684210527
2. In cross validation data 514 out of 532 : 96.61654135338345

3.2.2 Univariate Analysis on Variation Feature

Q7. Variation, What type of feature is it ?

Ans. Variation is a categorical variable

Q8. How many categories are there?

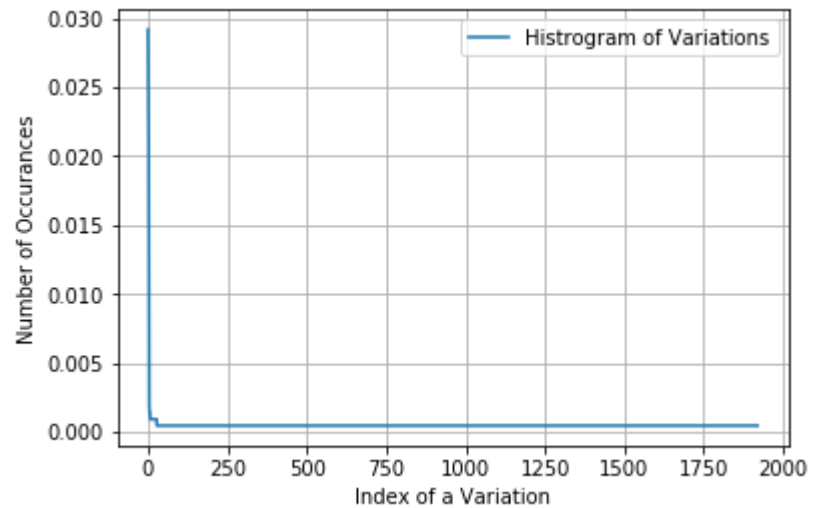

```
In [34]: 1 unique_variations = train_df['Variation'].value_counts()
2 print('Number of Unique Variations :', unique_variations.shape[0])
3 # the top 10 variations that occurred most
4 print(unique_variations.head(10))
```

```
Number of Unique Variations : 1920
Truncating_Mutations      62
Deletion                  52
Amplification             45
Fusions                   21
Overexpression            4
G12V                      3
E17K                      3
F28L                      2
T73I                      2
A146T                     2
Name: Variation, dtype: int64
```

```
In [35]: 1 print("Ans: There are", unique_variations.shape[0] , "different categories of variations in the train data,
```

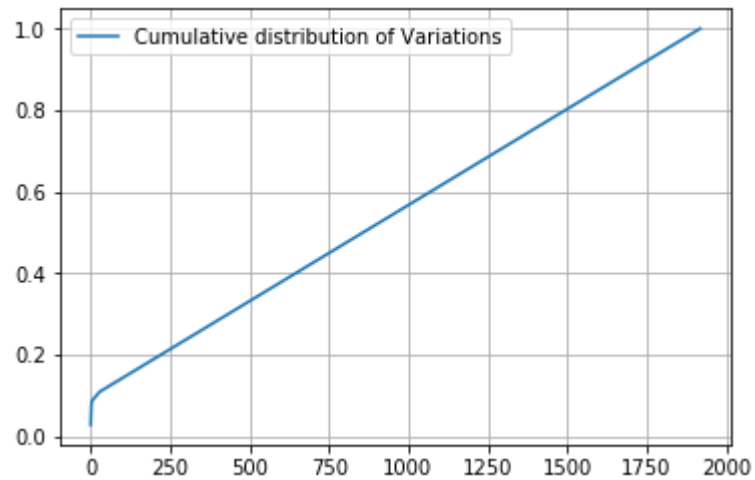
```
Ans: There are 1920 different categories of variations in the train data, and they are distributed as follows
```

```
In [36]: 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()
```



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

```
[0.02919021 0.05367232 0.07485876 ... 0.99905838 0.99952919 1.          ]
```



Q9. How to featurize this Variation feature ?

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

1. One hot Encoding
2. Response coding

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

```
In [0]: 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))
```

```
In [0]: 1 print("train_variation_feature_responseCoding is a converted feature using the response coding method. The
train_variation_feature_responseCoding is a converted feature using the response coding method. The shape of
Variation feature: (2124, 9)
```

```
In [38]: 1 # TFIDF
2 variation_vectorizer = TfidfVectorizer()
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'])
```

```
In [39]: 1 print("train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The s
train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape of V
ariation feature: (2124, 1948)
```

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

Let's build a model just like the earlier!

```

In [40]: 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.
4 # -----
5 # default parameters
6 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None,
7 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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, n_jobs=3)
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_, ep
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

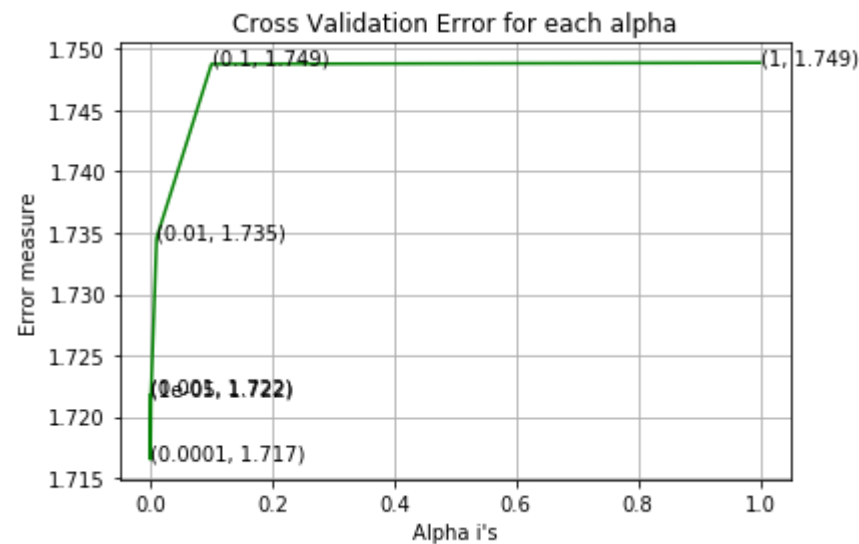
```

```

42 best_alpha = np.argmin(cv_log_error_array)
43 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42, n_jobs=3)
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_cv, 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: 1.7217921451034903
 For values of alpha = 0.0001 The log loss is: 1.716520021064393
 For values of alpha = 0.001 The log loss is: 1.7219560644207113
 For values of alpha = 0.01 The log loss is: 1.7345769217678837
 For values of alpha = 0.1 The log loss is: 1.748764356802458
 For values of alpha = 1 The log loss is: 1.748872154328207



For values of best alpha = 0.0001 The train log loss is: 0.7578118654422595
 For values of best alpha = 0.0001 The cross validation log loss is: 1.716520021064393
 For values of best alpha = 0.0001 The test log loss is: 1.7176841315254827

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

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

```
In [41]: 1 print("Q12. How many data points are covered by total ", unique_variations.shape[0], " genes in test and cr
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[0])
5 print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0],":", (cv_coverage/cv_df.shape[0])
```

Q12. How many data points are covered by total 1920 genes in test and cross validation data sets?

Ans

1. In test data 70 out of 665 : 10.526315789473683
2. In cross validation data 49 out of 532 : 9.210526315789473

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 prediciting y_i?
5. Is the text feature stable across train, test and CV datasets?

```
In [0]: 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
```

```
In [0]: 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
```



```
In [82]: 1 # building a CountVectorizer with all the words that occurred minimum 3 times in train data
2 text_vectorizer = TfidfVectorizer(min_df=3, max_features=1000)
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) vect
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 : 1000

```
In [43]: 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)
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-43-6d69a1102340> in <module>
      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
```

NameError: name 'extract_dictionary_paddle' is not defined

```
In [0]: 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)
```

```
In [0]: 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.
4 test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.sum
5 cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=
```

```
In [83]: 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)
```

```
In [84]: 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()))
```

```
In [85]: 1 # Number of words for a given frequency.
        2 print(Counter(sorted_text_occur))
```

```
Counter({253.45612854557848: 1, 181.8561472374477: 1, 141.81578713152044: 1, 132.4119790600596: 1, 132.248561
3245489: 1, 118.82156729160252: 1, 118.47323464437058: 1, 114.21120064713215: 1, 112.33018004794411: 1, 107.6
71276275327: 1, 107.06482767967276: 1, 89.1263711122647: 1, 88.35227467864152: 1, 86.7254319818862: 1, 81.641
82663564429: 1, 81.46677217209721: 1, 79.43685103101126: 1, 79.33372393634076: 1, 78.82645319750993: 1, 77.99
997212017877: 1, 77.73834725495779: 1, 75.16572289005073: 1, 69.6619948558994: 1, 68.69049717413029: 1, 67.88
022151261575: 1, 67.75392879950992: 1, 67.07620365755038: 1, 65.58130505874054: 1, 65.20226276415035: 1, 63.7
65880603490835: 1, 63.457339760323656: 1, 62.80280020966844: 1, 62.54808981645129: 1, 60.470597635312835: 1,
59.54634341135907: 1, 58.56814084030761: 1, 56.868959718213254: 1, 56.835933427076434: 1, 56.00033536171236:
1, 54.974205192575575: 1, 52.51137471639824: 1, 50.71620160884344: 1, 49.231602410746525: 1, 48.5162219160665
8: 1, 47.34244737784118: 1, 46.46346391448383: 1, 45.89412565164875: 1, 45.61689695648438: 1, 44.560208931307
01: 1, 43.90986592280627: 1, 43.74834766086299: 1, 43.720834052367266: 1, 43.6635037667005: 1, 43.62042213852
648: 1, 42.85457592372548: 1, 42.69494270138875: 1, 42.36474321367222: 1, 42.127512183105054: 1, 41.781929560
19826: 1, 41.58399107872969: 1, 41.124372893441084: 1, 41.05810960675291: 1, 41.00073002392298: 1, 40.5253785
7903626: 1, 39.61881318803275: 1, 39.41481629416075: 1, 39.402334019248535: 1, 39.33721055039442: 1, 38.22811
5566851095: 1, 38.07748568594918: 1, 37.882403866262: 1, 37.870855607768505: 1, 37.161319920779924: 1, 36.717
0867600636: 1, 36.693707887314176: 1, 36.3359922393846: 1, 35.989411003186724: 1, 35.81540382208122: 1, 35.71
451539590343: 1, 35.62467322679815: 1, 35.26610066155905: 1, 34.996741397988245: 1, 34.94013138955001: 1, 34.
63432109591924: 1, 34.50489211328151: 1, 34.3964684020224: 1, 34.283202692168494: 1, 34.08132556614205: 1, 3
4.04406894264587: 1, 33.498236231193: 1, 33.09774659467206: 1, 33.09276623545926: 1, 32.85685527883516: 1, 3
2.7500000000000001: 1, 32.600000000000001: 1, 32.500000000000001: 1, 32.350000000000001: 1, 32.250000000000001: 1, 32.150000000000001: 1, 32.050000000000001: 1, 31.950000000000001: 1, 31.850000000000001: 1, 31.750000000000001: 1, 31.650000000000001: 1, 31.550000000000001: 1, 31.450000000000001: 1, 31.350000000000001: 1, 31.250000000000001: 1, 31.150000000000001: 1, 31.050000000000001: 1, 30.950000000000001: 1, 30.850000000000001: 1, 30.750000000000001: 1, 30.650000000000001: 1, 30.550000000000001: 1, 30.450000000000001: 1, 30.350000000000001: 1, 30.250000000000001: 1, 30.150000000000001: 1, 30.050000000000001: 1, 30.000000000000001: 1, 29.950000000000001: 1, 29.900000000000001: 1, 29.850000000000001: 1, 29.800000000000001: 1, 29.750000000000001: 1, 29.700000000000001: 1, 29.650000000000001: 1, 29.600000000000001: 1, 29.550000000000001: 1, 29.500000000000001: 1, 29.450000000000001: 1, 29.400000000000001: 1, 29.350000000000001: 1, 29.300000000000001: 1, 29.250000000000001: 1, 29.200000000000001: 1, 29.150000000000001: 1, 29.100000000000001: 1, 29.050000000000001: 1, 29.000000000000001: 1, 28.950000000000001: 1, 28.900000000000001: 1, 28.850000000000001: 1, 28.800000000000001: 1, 28.750000000000001: 1, 28.700000000000001: 1, 28.650000000000001: 1, 28.600000000000001: 1, 28.550000000000001: 1, 28.500000000000001: 1, 28.450000000000001: 1, 28.400000000000001: 1, 28.350000000000001: 1, 28.300000000000001: 1, 28.250000000000001: 1, 28.200000000000001: 1, 28.150000000000001: 1, 28.100000000000001: 1, 28.050000000000001: 1, 28.000000000000001: 1, 27.950000000000001: 1, 27.900000000000001: 1, 27.850000000000001: 1, 27.800000000000001: 1, 27.750000000000001: 1, 27.700000000000001: 1, 27.650000000000001: 1, 27.600000000000001: 1, 27.550000000000001: 1, 27.500000000000001: 1, 27.450000000000001: 1, 27.400000000000001: 1, 27.350000000000001: 1, 27.300000000000001: 1, 27.250000000000001: 1, 27.200000000000001: 1, 27.150000000000001: 1, 27.100000000000001: 1, 27.050000000000001: 1, 27.000000000000001: 1, 26.950000000000001: 1, 26.900000000000001: 1, 26.850000000000001: 1, 26.800000000000001: 1, 26.750000000000001: 1, 26.700000000000001: 1, 26.650000000000001: 1, 26.600000000000001: 1, 26.550000000000001: 1, 26.500000000000001: 1, 26.450000000000001: 1, 26.400000000000001: 1, 26.350000000000001: 1, 26.300000000000001: 1, 26.250000000000001: 1, 26.200000000000001: 1, 26.150000000000001: 1, 26.100000000000001: 1, 26.050000000000001: 1, 26.000000000000001: 1, 25.950000000000001: 1, 25.900000000000001: 1, 25.850000000000001: 1, 25.800000000000001: 1, 25.750000000000001: 1, 25.700000000000001: 1, 25.650000000000001: 1, 25.600000000000001: 1, 25.550000000000001: 1, 25.500000000000001: 1, 25.450000000000001: 1, 25.400000000000001: 1, 25.350000000000001: 1, 25.300000000000001: 1, 25.250000000000001: 1, 25.200000000000001: 1, 25.150000000000001: 1, 25.100000000000001: 1, 25.050000000000001: 1, 25.000000000000001: 1, 24.950000000000001: 1, 24.900000000000001: 1, 24.850000000000001: 1, 24.800000000000001: 1, 24.750000000000001: 1, 24.700000000000001: 1, 24.650000000000001: 1, 24.600000000000001: 1, 24.550000000000001: 1, 24.500000000000001: 1, 24.450000000000001: 1, 24.400000000000001: 1, 24.350000000000001: 1, 24.300000000000001: 1, 24.250000000000001: 1, 24.200000000000001: 1, 24.150000000000001: 1, 24.100000000000001: 1, 24.050000000000001: 1, 24.000000000000001: 1, 23.950000000000001: 1, 23.900000000000001: 1, 23.850000000000001: 1, 23.800000000000001: 1, 23.750000000000001: 1, 23.700000000000001: 1, 23.650000000000001: 1, 23.600000000000001: 1, 23.550000000000001: 1, 23.500000000000001: 1, 23.450000000000001: 1, 23.400000000000001: 1, 23.350000000000001: 1, 23.300000000000001: 1, 23.250000000000001: 1, 23.200000000000001: 1, 23.150000000000001: 1, 23.100000000000001: 1, 23.050000000000001: 1, 23.000000000000001: 1, 22.950000000000001: 1, 22.900000000000001: 1, 22.850000000000001: 1, 22.800000000000001: 1, 22.750000000000001: 1, 22.700000000000001: 1, 22.650000000000001: 1, 22.600000000000001: 1, 22.550000000000001: 1, 22.500000000000001: 1, 22.450000000000001: 1, 22.400000000000001: 1, 22.350000000000001: 1, 22.300000000000001: 1, 22.250000000000001: 1, 22.200000000000001: 1, 22.150000000000001: 1, 22.100000000000001: 1, 22.050000000000001: 1, 22.000000000000001: 1, 21.950000000000001: 1, 21.900000000000001: 1, 21.850000000000001: 1, 21.800000000000001: 1, 21.750000000000001: 1, 21.700000000000001: 1, 21.650000000000001: 1, 21.600000000000001: 1, 21.550000000000001: 1, 21.500000000000001: 1, 21.450000000000001: 1, 21.400000000000001: 1, 21.350000000000001: 1, 21.300000000000001: 1, 21.250000000000001: 1, 21.200000000000001: 1, 21.150000000000001: 1, 21.100000000000001: 1, 21.050000000000001: 1, 21.000000000000001: 1, 20.950000000000001: 1, 20.900000000000001: 1, 20.850000000000001: 1, 20.800000000000001: 1, 20.750000000000001: 1, 20.700000000000001: 1, 20.650000000000001: 1, 20.600000000000001: 1, 20.550000000000001: 1, 20.500000000000001: 1, 20.450000000000001: 1, 20.400000000000001: 1, 20.350000000000001: 1, 20.300000000000001: 1, 20.250000000000001: 1, 20.200000000000001: 1, 20.150000000000001: 1, 20.100000000000001: 1, 20.050000000000001: 1, 20.000000000000001: 1, 19.950000000000001: 1, 19.900000000000001: 1, 19.850000000000001: 1, 19.800000000000001: 1, 19.750000000000001: 1, 19.700000000000001: 1, 19.650000000000001: 1, 19.600000000000001: 1, 19.550000000000001: 1, 19.500000000000001: 1, 19.450000000000001: 1, 19.400000000000001: 1, 19.350000000000001: 1, 19.300000000000001: 1, 19.250000000000001: 1, 19.200000000000001: 1, 19.150000000000001: 1, 19.100000000000001: 1, 19.050000000000001: 1, 19.000000000000001: 1, 18.950000000000001: 1, 18.900000000000001: 1, 18.850000000000001: 1, 18.800000000000001: 1, 18.750000000000001: 1, 18.700000000000001: 1, 18.650000000000001: 1, 18.600000000000001: 1, 18.550000000000001: 1, 18.500000000000001: 1, 18.450000000000001: 1, 18.400000000000001: 1, 18.350000000000001: 1, 18.300000000000001: 1, 18.250000000000001: 1, 18.200000000000001: 1, 18.150000000000001: 1, 18.100000000000001: 1, 18.050000000000001: 1, 18.000000000000001: 1, 17.950000000000001: 1, 17.900000000000001: 1, 17.850000000000001: 1, 17.800000000000001: 1, 17.750000000000001: 1, 17.700000000000001: 1, 17.650000000000001: 1, 17.600000000000001: 1, 17.550000000000001: 1, 17.500000000000001: 1, 17.450000000000001: 1, 17.400000000000001: 1, 17.350000000000001: 1, 17.300000000000001: 1, 17.250000000000001: 1, 17.200000000000001: 1, 17.150000000000001: 1, 17.100000000000001: 1, 17.050000000000001: 1, 17.000000000000001: 1, 16.950000000000001: 1, 16.900000000000001: 1, 16.850000000000001: 1, 16.800000000000001: 1, 16.750000000000001: 1, 16.700000000000001: 1, 16.650000000000001: 1, 16.600000000000001: 1, 16.550000000000001: 1, 16.500000000000001: 1, 16.450000000000001: 1, 16.400000000000001: 1, 16.350000000000001: 1, 16.300000000000001: 1, 16.250000000000001: 1, 16.200000000000001: 1, 16.150000000000001: 1, 16.100000000000001: 1, 16.050000000000001: 1, 16.000000000000001: 1, 15.950000000000001: 1, 15.900000000000001: 1, 15.850000000000001: 1, 15.800000000000001: 1, 15.750000000000001: 1, 15.700000000000001: 1, 15.650000000000001: 1, 15.600000000000001: 1, 15.550000000000001: 1, 15.500000000000001: 1, 15.450000000000001: 1, 15.400000000000001: 1, 15.350000000000001: 1, 15.300000000000001: 1, 15.250000000000001: 1, 15.200000000000001: 1, 15.150000000000001: 1, 15.100000000000001: 1, 15.050000000000001: 1, 15.000000000000001: 1, 14.950000000000001: 1, 14.900000000000001: 1, 14.850000000000001: 1, 14.800000000000001: 1, 14.750000000000001: 1, 14.700000000000001: 1, 14.650000000000001: 1, 14.600000000000001: 1, 14.550000000000001: 1, 14.500000000000001: 1, 14.450000000000001: 1, 14.400000000000001: 1, 14.350000000000001: 1, 14.300000000000001: 1, 14.250000000000001: 1, 14.200000000000001: 1, 14.150000000000001: 1, 14.100000000000001: 1, 14.050000000000001: 1, 14.000000000000001: 1, 13.950000000000001: 1, 13.900000000000001: 1, 13.850000000000001: 1, 13.800000000000001: 1, 13.750000000000001: 1, 13.700000000000001: 1, 13.650000000000001: 1, 13.600000000000001: 1, 13.550000000000001: 1, 13.500000000000001: 1, 13.450000000000001: 1, 13.400000000000001: 1, 13.350000000000001: 1, 13.300000000000001: 1, 13.250000000000001: 1, 13.200000000000001: 1, 13.150000000000001: 1, 13.100000000000001: 1, 13.050000000000001: 1, 13.000000000000001: 1, 12.950000000000001: 1, 12.900000000000001: 1, 12.850000000000001: 1, 12.800000000000001: 1, 12.750000000000001: 1, 12.700000000000001: 1, 12.650000000000001: 1, 12.600000000000001: 1, 12.550000000000001: 1, 12.500000000000001: 1, 12.450000000000001: 1, 12.400000000000001: 1, 12.350000000000001: 1, 12.300000000000001: 1, 12.250000000000001: 1, 12.200000000000001: 1, 12.150000000000001: 1, 12.100000000000001: 1, 12.050000000000001: 1, 12.000000000000001: 1, 11.950000000000001: 1, 11.900000000000001: 1, 11.850000000000001: 1, 11.800000000000001: 1, 11.750000000000001: 1, 11.700000000000001: 1, 11.650000000000001: 1, 11.600000000000001: 1, 11.550000000000001: 1, 11.500000000000001: 1, 11.450000000000001: 1, 11.400000000000001: 1, 11.350000000000001: 1, 11.300000000000001: 1, 11.250000000000001: 1, 11.200000000000001: 1, 11.150000000000001: 1, 11.100000000000001: 1, 11.050000000000001: 1, 11.000000000000001: 1, 10.950000000000001: 1, 10.900000000000001: 1, 10.850000000000001: 1, 10.800000000000001: 1, 10.750000000000001: 1, 10.700000000000001: 1, 10.650000000000001: 1, 10.600000000000001: 1, 10.550000000000001: 1, 10.500000000000001: 1, 10.450000000000001: 1, 10.400000000000001: 1, 10.350000000000001: 1, 10.300000000000001: 1, 10.250000000000001: 1, 10.200000000000001: 1, 10.150000000000001: 1, 10.100000000000001: 1, 10.050000000000001: 1, 10.000000000000001: 1, 9.950000000000001: 1, 9.900000000000001: 1, 9.850000000000001: 1, 9.800000000000001: 1, 9.750000000000001: 1, 9.700000000000001: 1, 9.650000000000001: 1, 9.600000000000001: 1, 9.550000000000001: 1, 9.500000000000001: 1, 9.450000000000001: 1, 9.400000000000001: 1, 9.350000000000001: 1, 9.30000000000
```

```

In [86]: 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.
5 # -----
6 # default parameters
7 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None,
8 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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, n_jobs=3)
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_, ep
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

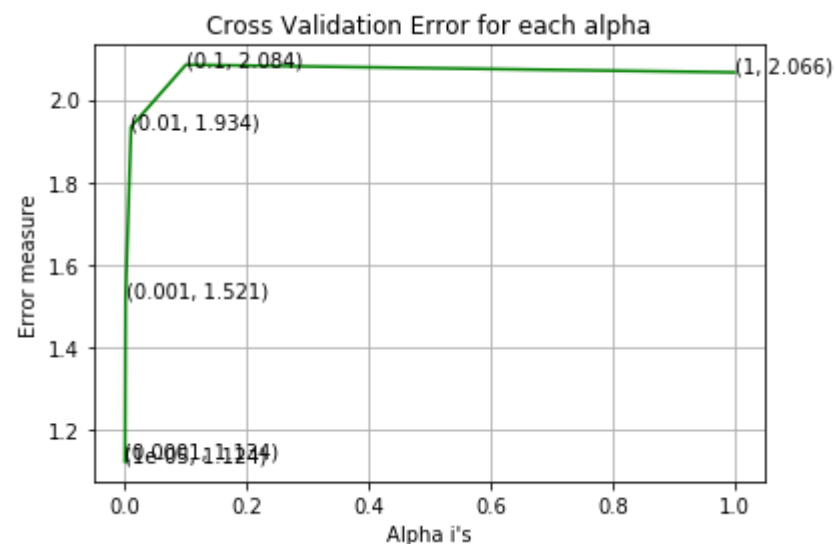
```

```

42 best_alpha = np.argmin(cv_log_error_array)
43 clf = SGDCClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42, n_jobs=3)
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_cv, 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: 1.1236270404130342
 For values of alpha = 0.0001 The log loss is: 1.1339643177809415
 For values of alpha = 0.001 The log loss is: 1.5214866424038809
 For values of alpha = 0.01 The log loss is: 1.9340658920508387
 For values of alpha = 0.1 The log loss is: 2.0844881800102137
 For values of alpha = 1 The log loss is: 2.0659298902039547



For values of best alpha = 1e-05 The train log loss is: 0.7136163056681504
 For values of best alpha = 1e-05 The cross validation log loss is: 1.1236270404130342
 For values of best alpha = 1e-05 The test log loss is: 1.1432071674613324

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

Ans. Yes, it seems like!

```
In [87]: 1 def get_intersec_text(df):
2         df_text_vec = TfidfVectorizer(min_df=3)
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
```

```
In [88]: 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")
```

3.414 % of word of test data appeared in train data

3.792 % of word of Cross Validation appeared in train data

4. Machine Learning Models

```
In [51]: 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)
```

```
In [52]: 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)
```



```

In [53]: 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(min_df=3)
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 [{}]"
24                      .format(word,yes_no))
25        elif (v < fea1_len+fea2_len):
26            word = var_vec.get_feature_names()[v-(fea1_len)]
27            yes_no = True if word == var else False
28            if yes_no:
29                word_present += 1
30                print(i, "variation feature [{}] present in test data point [{}]"
31                      .format(word,yes_no))
32        else:
33            word = text_vec.get_feature_names()[v-(fea1_len+fea2_len)]
34            yes_no = True if word in text.split() else False
35            if yes_no:
36                word_present += 1
37                print(i, "Text feature [{}] present in test data point [{}]"
38                      .format(word,yes_no))
39    print("Out of the top ",no_features," features ", word_present, "are present in query point")

```

Stacking the three types of features

```
In [89]: 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']))
23
24
25 # train_gene_var_responseCoding = np.hstack((train_gene_feature_responseCoding, train_variation_feature_resp
26 # test_gene_var_responseCoding = np.hstack((test_gene_feature_responseCoding, test_variation_feature_respons
27 # cv_gene_var_responseCoding = np.hstack((cv_gene_feature_responseCoding, cv_variation_feature_responseCodin
28
29 # train_x_responseCoding = np.hstack((train_gene_var_responseCoding, train_text_feature_responseCoding))
30 # test_x_responseCoding = np.hstack((test_gene_var_responseCoding, test_text_feature_responseCoding))
31 # cv_x_responseCoding = np.hstack((cv_gene_var_responseCoding, cv_text_feature_responseCoding))
32
```

```
In [90]: 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, 3183)
(number of data points * number of features) in test data = (665, 3183)
(number of data points * number of features) in cross validation data = (532, 3183)
```

```
In [56]: 1 print(" Response encoding features :")
2 print("(number of data points * number of features) in train data = ", train_x_responseCoding.shape)
3 print("(number of data points * number of features) in test data = ", test_x_responseCoding.shape)
4 print("(number of data points * number of features) in cross validation data =", cv_x_responseCoding.shape)
```

Response encoding features :

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-56-7d7a26188f88> in <module>
      1 print(" Response encoding features :")
----> 2 print("(number of data points * number of features) in train data = ", train_x_responseCoding.shape)
      3 print("(number of data points * number of features) in test data = ", test_x_responseCoding.shape)
      4 print("(number of data points * number of features) in cross validation data =", cv_x_responseCoding.
shape)
```

NameError: name 'train_x_responseCoding' is not defined

4.1. Base Line Model

4.1.1. Naive Bayes

4.1.1.1. Hyper parameter tuning

```

In [91]: 1 # find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/skl
2 # -----
3 # default paramters
4 # sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)
5
6 # some of methods of MultinomialNB()
7 # fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
8 # predict(X) Perform classification on an array of test vectors X.
9 # predict_log_proba(X) Return log-probability estimates for the test vector X.
10 # -----
11 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm
12 # -----
13
14
15 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
16 # -----
17 # default paramters
18 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
19 #
20 # some of the methods of CalibratedClassifierCV()
21 # fit(X, y[, sample_weight]) Fit the calibrated model
22 # get_params([deep]) Get parameters for this estimator.
23 # predict(X) Predict the target of new samples.
24 # predict_proba(X) Posterior probabilities of classification
25 # -----
26 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm
27 # -----
28
29
30 alpha = [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100,1000]
31 cv_log_error_array = []
32 for i in alpha:
33     print("for alpha =", i)
34     clf = MultinomialNB(alpha=i)
35     clf.fit(train_x_onehotCoding, train_y)
36     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
37     sig_clf.fit(train_x_onehotCoding, train_y)
38     sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
39     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
40     # to avoid rounding error while multiplying probabilites we use log-probability estimates
41     print("Log Loss :",log_loss(cv_y, sig_clf_probs))

```

```

42
43 fig, ax = plt.subplots()
44 ax.plot(np.log10(alpha), cv_log_error_array,c='g')
45 for i, txt in enumerate(np.round(cv_log_error_array,3)):
46     ax.annotate((alpha[i],str(txt)), (np.log10(alpha[i]),cv_log_error_array[i]))
47 plt.grid()
48 plt.xticks(np.log10(alpha))
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 = MultinomialNB(alpha=alpha[best_alpha])
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
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, predict_
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_cv,
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, predict_y,
68

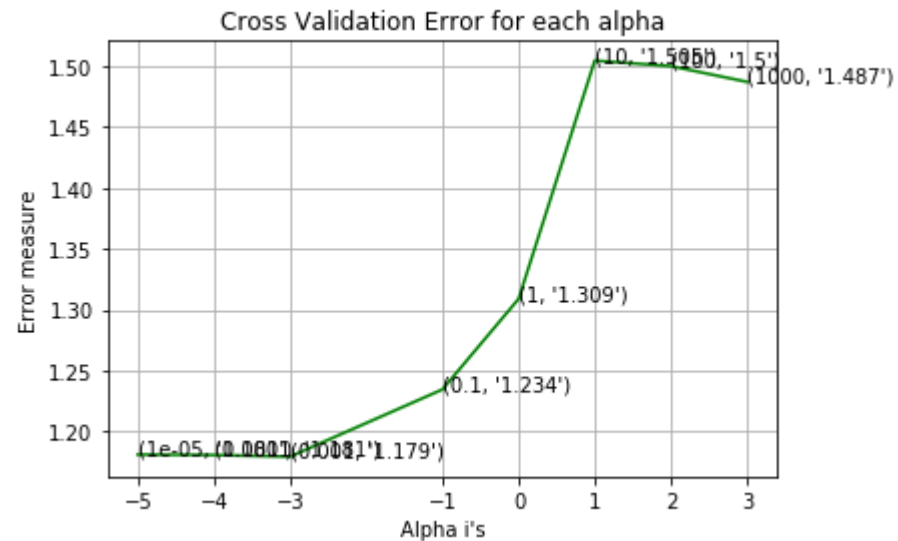
```

```

for alpha = 1e-05
Log Loss : 1.1808361651709762
for alpha = 0.0001
Log Loss : 1.1805077618041426
for alpha = 0.001
Log Loss : 1.1791699326060952
for alpha = 0.1
Log Loss : 1.2343585317850052
for alpha = 1
Log Loss : 1.3087969709305698
for alpha = 10
Log Loss : 1.5045900910017453
for alpha = 100
Log Loss : 1.5000429623188356

```

```
for alpha = 1000
Log Loss : 1.4871093699387652
```



```
For values of best alpha = 0.001 The train log loss is: 0.5142704513553946
For values of best alpha = 0.001 The cross validation log loss is: 1.1791699326060952
For values of best alpha = 0.001 The test log loss is: 1.2041049705163938
```

4.1.1.2. Testing the model with best hyper paramters

```

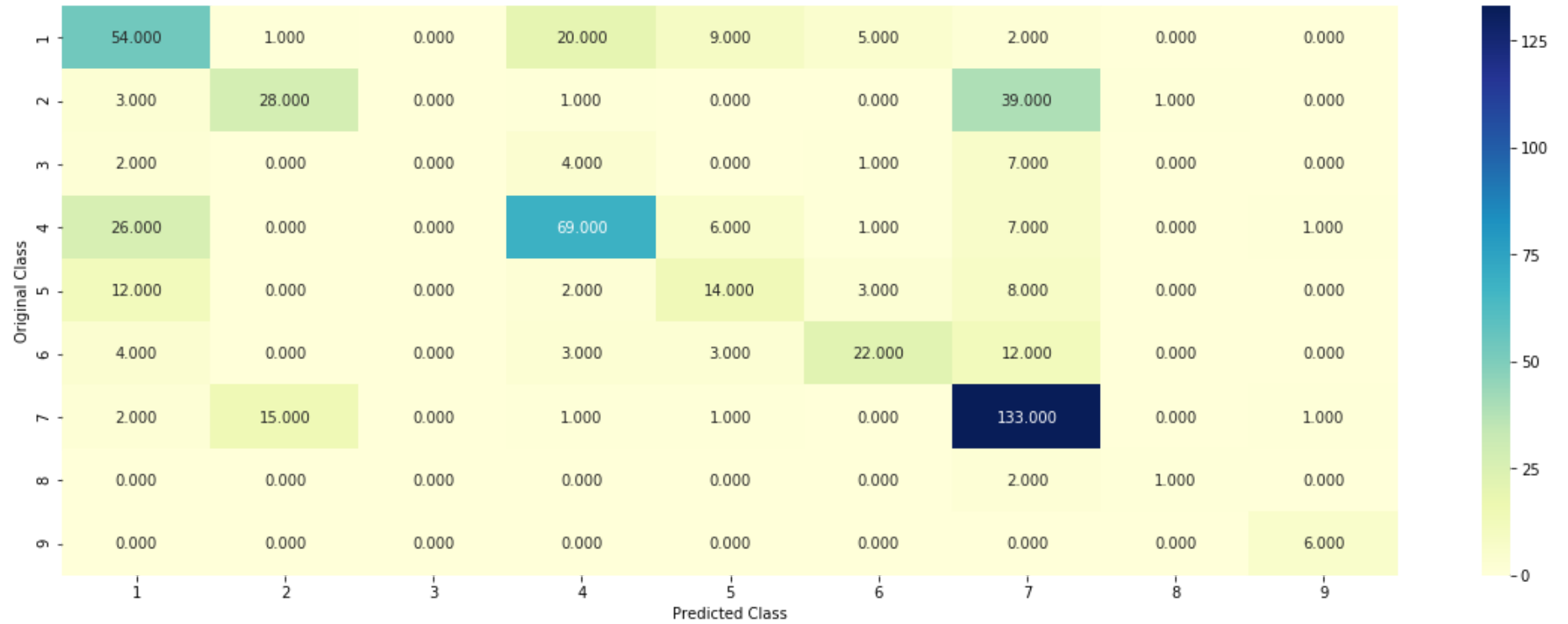
In [92]: 1 # find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/skl
2 # -----
3 # default paramters
4 # sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)
5
6 # some of methods of MultinomialNB()
7 # fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
8 # predict(X) Perform classification on an array of test vectors X.
9 # predict_log_proba(X) Return log-probability estimates for the test vector X.
10 # -----
11 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm
12 # -----
13
14
15 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
16 # -----
17 # default paramters
18 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
19 #
20 # some of the methods of CalibratedClassifierCV()
21 # fit(X, y[, sample_weight]) Fit the calibrated model
22 # get_params([deep]) Get parameters for this estimator.
23 # predict(X) Predict the target of new samples.
24 # predict_proba(X) Posterior probabilities of classification
25 # -----
26
27 clf = MultinomialNB(alpha=alpha[best_alpha])
28 clf.fit(train_x_onehotCoding, train_y)
29 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
30 sig_clf.fit(train_x_onehotCoding, train_y)
31 sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
32 # to avoid rounding error while multiplying probabillites we use log-probability estimates
33 print("Log Loss :", log_loss(cv_y, sig_clf_probs))
34 print("Number of missclassified point :", np.count_nonzero((sig_clf.predict(cv_x_onehotCoding) - cv_y))/cv_y)
35 plot_confusion_matrix(cv_y, sig_clf.predict(cv_x_onehotCoding.toarray()))

```

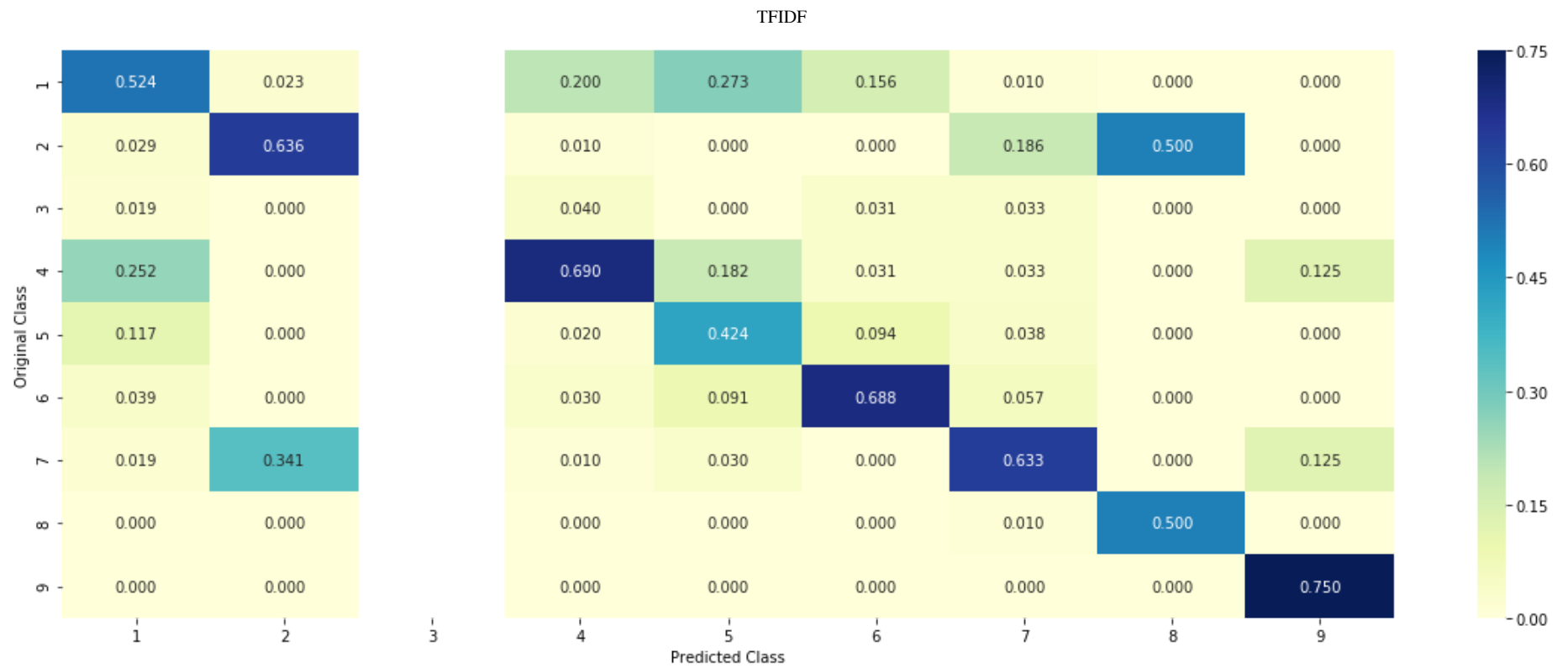
Log Loss : 1.1791699326060952

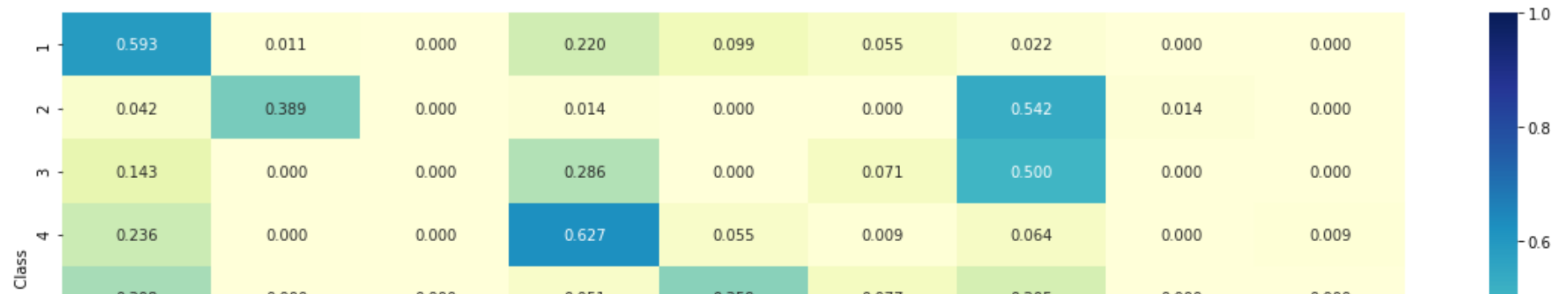
Number of missclassified point : 0.38533834586466165

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



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





4.1.1.3. Feature Importance, Correctly classified point

```
In [93]: 1 test_point_index = 1
2 no_feature = 100
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])
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 : 6

Predicted Class Probabilities: [[0.111 0.0628 0.0172 0.1298 0.1621 0.3874 0.1205 0.0048 0.0045]]

Actual Class : 4

Out of the top 100 features 0 are present in query point

4.1.1.4. Feature Importance, Incorrectly classified point

```
In [94]: 1 test_point_index = 100
2 no_feature = 100
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]), 2))
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 : 8

Predicted Class Probabilities: [[0.2506 0.0746 0.0195 0.1115 0.0526 0.0475 0.1427 0.2909 0.0102]]

Actual Class : 8

```
-----
65 Text feature [1000] present in test data point [True]
Out of the top 100 features 1 are present in query point
```

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

```

In [0]: 1 # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neig
2 # -----
3 # default parameter
4 # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
5 # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)
6
7 # methods of
8 # fit(X, y) : Fit the model using X as training data and y as target values
9 # predict(X):Predict the class labels for the provided data
10 # predict_proba(X):Return probability estimates for the test data X.
11 #-----
12 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-g
13 #-----
14
15
16 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
17 # -----
18 # default paramters
19 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
20 #
21 # some of the methods of CalibratedClassifierCV()
22 # fit(X, y[, sample_weight]) Fit the calibrated model
23 # get_params([deep]) Get parameters for this estimator.
24 # predict(X) Predict the target of new samples.
25 # predict_proba(X) Posterior probabilities of classification
26 #-----
27 # video link:
28 #-----
29
30
31 alpha = [5, 11, 15, 21, 31, 41, 51, 99]
32 cv_log_error_array = []
33 for i in alpha:
34     print("for alpha =", i)
35     clf = KNeighborsClassifier(n_neighbors=i)
36     clf.fit(train_x_responseCoding, train_y)
37     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
38     sig_clf.fit(train_x_responseCoding, train_y)
39     sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
40     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
41     # to avoid rounding error while multiplying probabilities we use log-probability estimates

```

```

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 = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
57     clf.fit(train_x_responseCoding, train_y)
58     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
59     sig_clf.fit(train_x_responseCoding, train_y)
60
61     predict_y = sig_clf.predict_proba(train_x_responseCoding)
62     print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_
63     predict_y = sig_clf.predict_proba(cv_x_responseCoding)
64     print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv,
65     predict_y = sig_clf.predict_proba(test_x_responseCoding)
66     print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y,
67

```

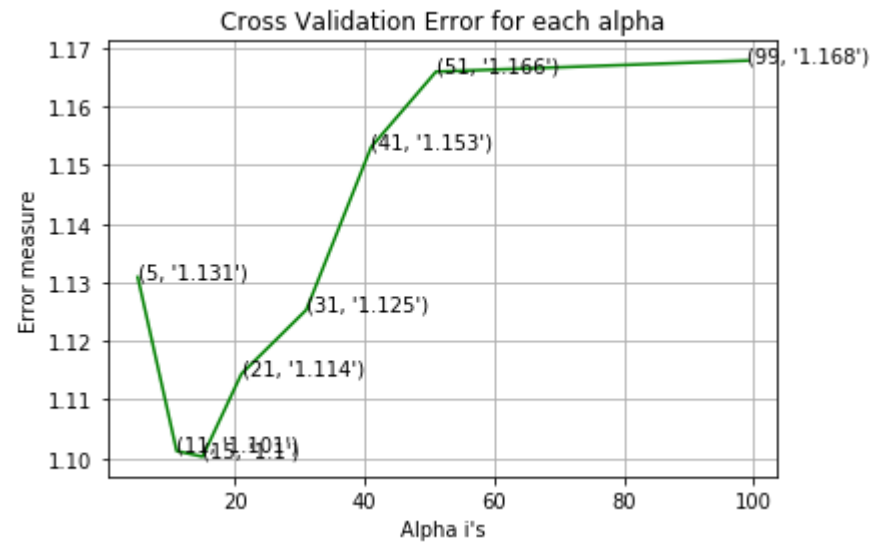
```

for alpha = 5
Log Loss : 1.1309170126692265
for alpha = 11
Log Loss : 1.1012397762291362
for alpha = 15
Log Loss : 1.1002748803755749
for alpha = 21
Log Loss : 1.1144110925957647
for alpha = 31
Log Loss : 1.1253206995500455
for alpha = 41
Log Loss : 1.1530939773909168
for alpha = 51
Log Loss : 1.1659585643007098

```

```
for alpha = 99
```

For alpha = 15
Log Loss : 1.1678505034822115



For values of best alpha = 15 The train log loss is: 0.7056892871225193
For values of best alpha = 15 The cross validation log loss is: 1.1002748803755749
For values of best alpha = 15 The test log loss is: 1.0911901980302394

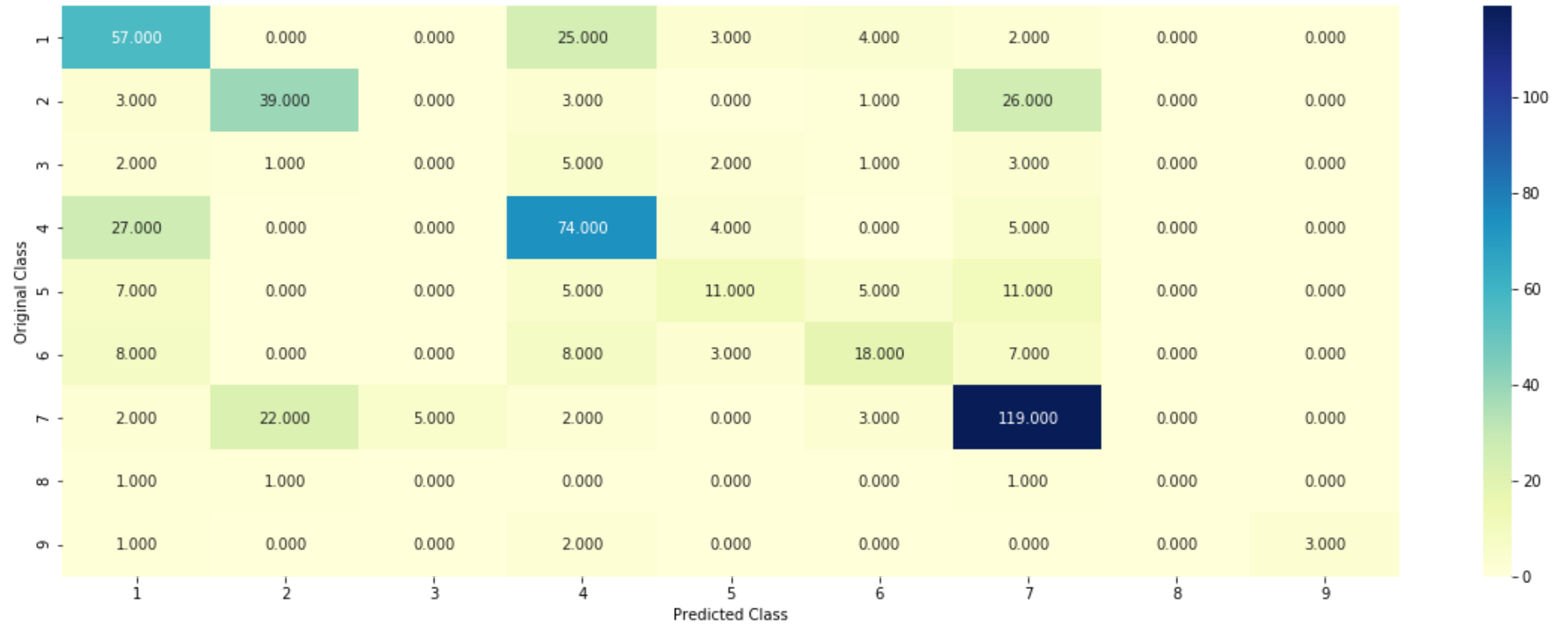
4.2.2. Testing the model with best hyper paramters

```
In [0]: 1 # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neig
2 # -----
3 # default parameter
4 # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
5 # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)
6
7 # methods of
8 # fit(X, y) : Fit the model using X as training data and y as target values
9 # predict(X):Predict the class labels for the provided data
10 # predict_proba(X):Return probability estimates for the test data X.
11 #-----
12 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-g
13 #-----
14 clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
15 predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseCoding, cv_y, clf)
```

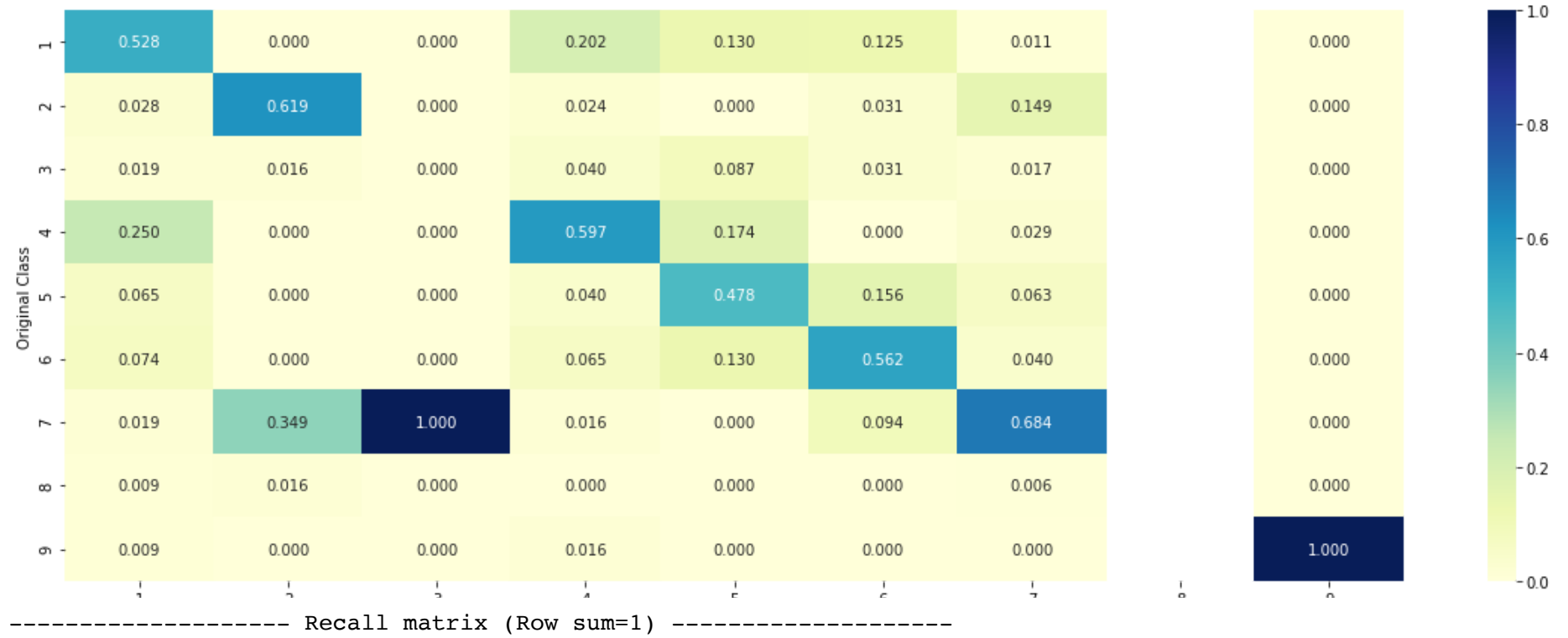
Log loss : 1.1002748803755749

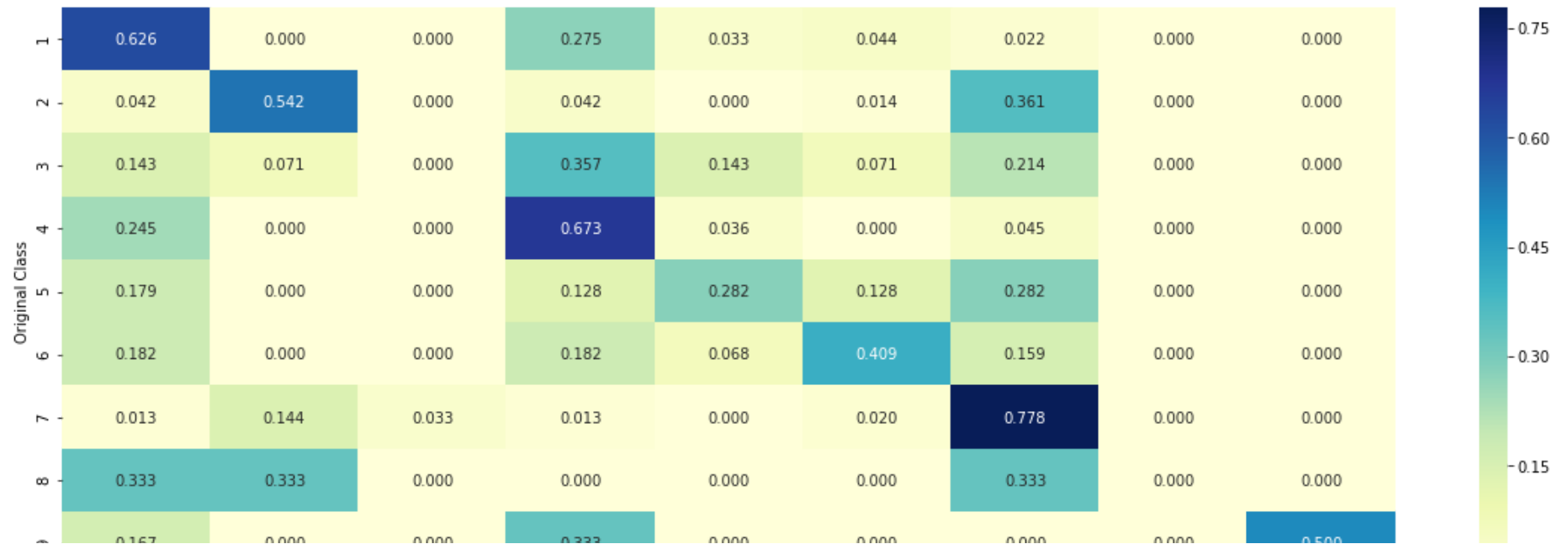
Number of mis-classified points : 0.3966165413533835

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



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





4.2.3. Sample Query point -1

```
In [0]: 1 clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
2 clf.fit(train_x_responseCoding, train_y)
3 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
4 sig_clf.fit(train_x_responseCoding, train_y)
5
6 test_point_index = 1
7 predicted_cls = sig_clf.predict(test_x_responseCoding[0].reshape(1,-1))
8 print("Predicted Class :", predicted_cls[0])
9 print("Actual Class :", test_y[test_point_index])
10 neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
11 print("The ", alpha[best_alpha], " nearest neighbours of the test points belongs to classes", train_y[neighbors[1][0]])
12 print("Fequency of nearest points :", Counter(train_y[neighbors[1][0]]))
```

Predicted Class : 6

Actual Class : 7

The 15 nearest neighbours of the test points belongs to classes [7 7 7 6 6 7 6 7 6 7 7 7 7 7 6]

Fequency of nearest points : Counter({7: 10, 6: 5})

4.2.4. Sample Query Point-2

```
In [0]: 1 clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
2 clf.fit(train_x_responseCoding, train_y)
3 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
4 sig_clf.fit(train_x_responseCoding, train_y)
5
6 test_point_index = 100
7
8 predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
9 print("Predicted Class :", predicted_cls[0])
10 print("Actual Class :", test_y[test_point_index])
11 neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
12 print("the k value for knn is",alpha[best_alpha],"and the nearest neighbours of the test points belongs to classes",test_y[neighbors[1][0]])
13 print("Frequency of nearest points :",Counter(train_y[neighbors[1][0]]))
```

Predicted Class : 7

Actual Class : 7

the k value for knn is 15 and the nearest neighbours of the test points belongs to classes [2 7 5 7 7 2 7 7 7
2 7 7 7 7 7]

Frequency of nearest points : Counter({7: 11, 2: 3, 5: 1})

4.3. Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper parameter tuning

```

In [61]: 1
2 # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.
3 # -----
4 # default parameters
5 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None,
6 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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-intuition-1
15 #-----
16
17
18 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
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(-6, 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, n_jobs=
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_state=0)
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, predict_y))
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_cv, predict_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, predict_y))

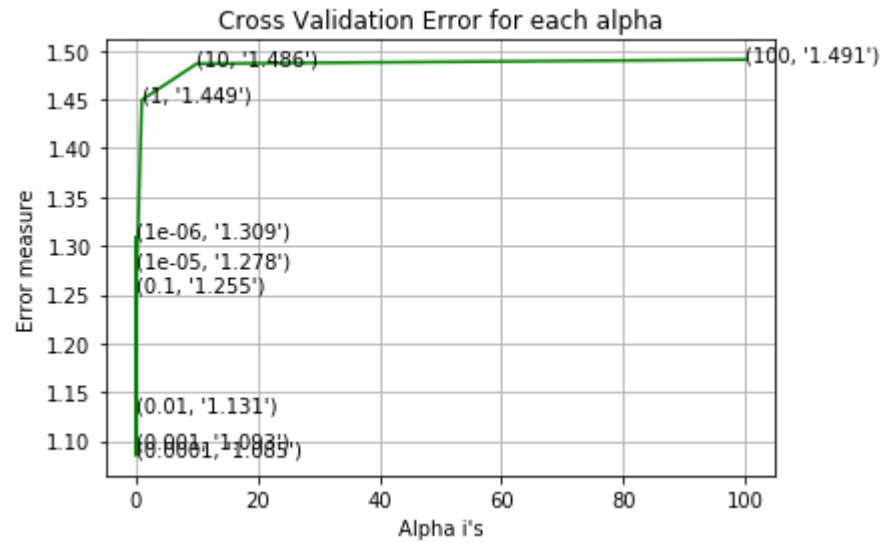
```

```

for alpha = 1e-06
Log Loss : 1.3088182435952065
for alpha = 1e-05
Log Loss : 1.2776692006389703
for alpha = 0.0001
Log Loss : 1.08501525947154
for alpha = 0.001
Log Loss : 1.093454308037274
for alpha = 0.01
Log Loss : 1.1313606853109164
for alpha = 0.1
Log Loss : 1.2545938717297804
for alpha = 1
Log Loss : 1.449333904723763
for alpha = 10

```

Log Loss : 1.4864632116895264
 for alpha = 100
 Log Loss : 1.4906691508630068



For values of best alpha = 0.0001 The train log loss is: 0.5563652337873516
 For values of best alpha = 0.0001 The cross validation log loss is: 1.08501525947154
 For values of best alpha = 0.0001 The test log loss is: 1.1115249520602726

4.3.1.2. Testing the model with best hyper paramters

```

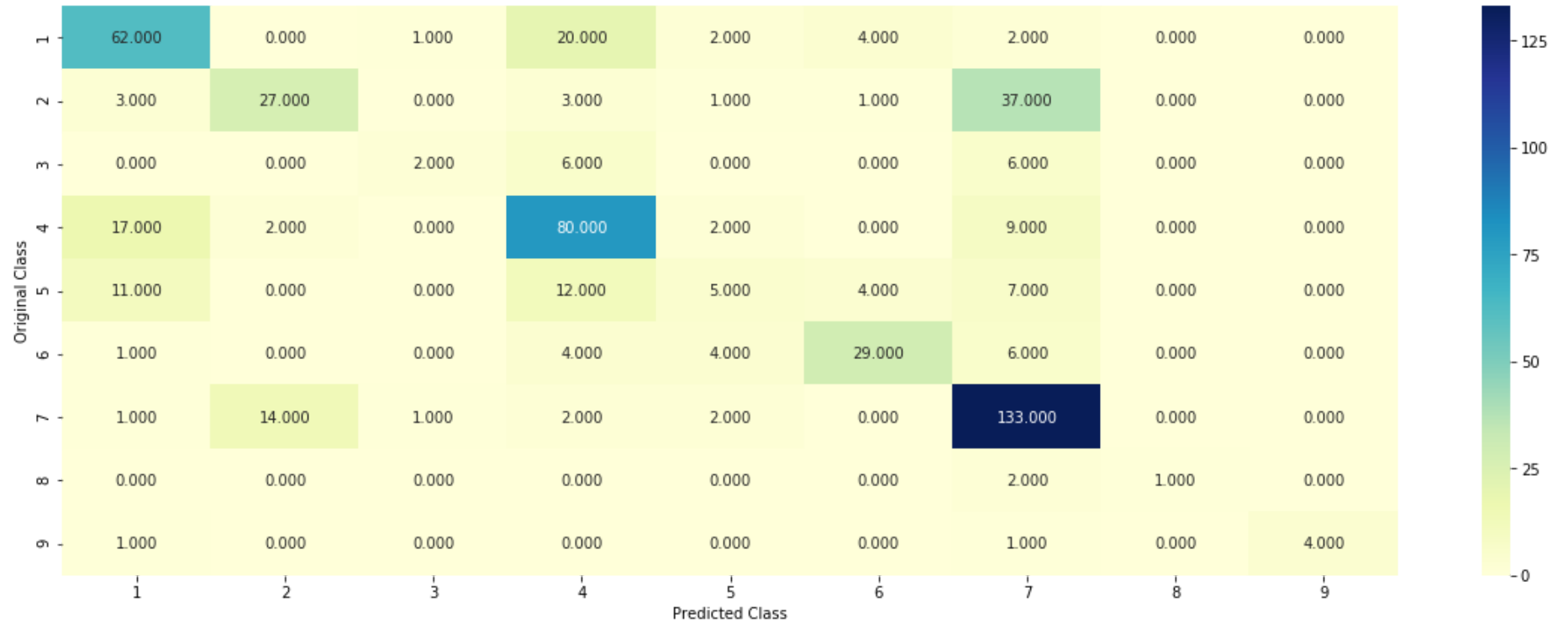
In [62]: 1 # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.
2 # -----
3 # default parameters
4 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None,
5 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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-1
14 #-----
15 clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_stat
16 predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

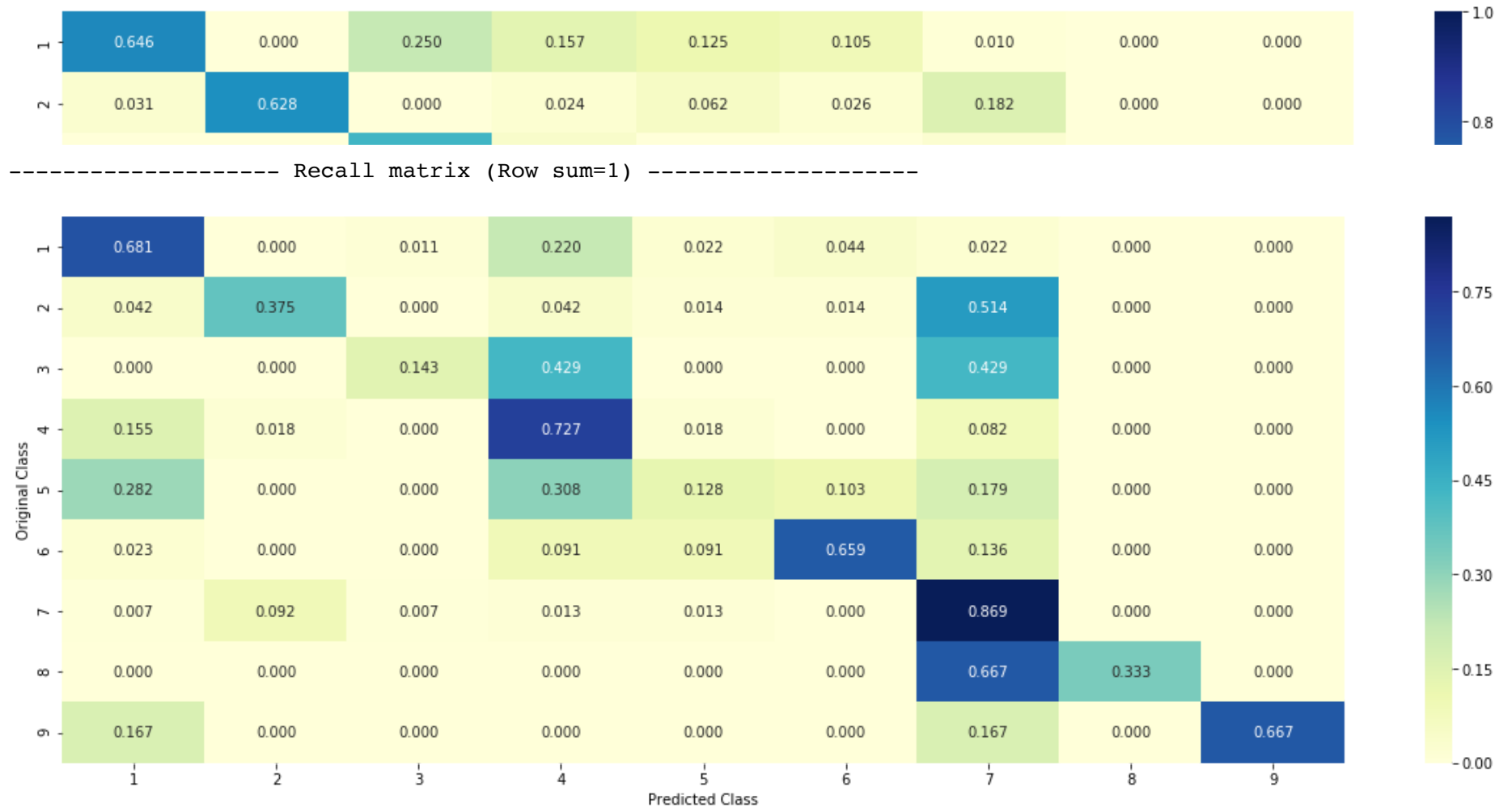
Log loss : 1.08501525947154

Number of mis-classified points : 0.35526315789473684

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



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



4.3.1.3. Feature Importance

```

In [63]: 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

```
In [64]: 1 # from tabulate import tabulate
2 clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=1)
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_index]), 4))
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_index])
```

Predicted Class : 4

Predicted Class Probabilities: [[0.161 0.0456 0.0365 0.5917 0.0637 0.0146 0.0478 0.0292 0.0099]]

Actual Class : 4

```
-----
150 Text feature [382] present in test data point [True]
236 Text feature [truncate] present in test data point [True]
244 Text feature [adapted] present in test data point [True]
284 Text feature [inactivation] present in test data point [True]
393 Text feature [instability] present in test data point [True]
436 Text feature [chenevix] present in test data point [True]
437 Text feature [hereditary] present in test data point [True]
439 Text feature [germline] present in test data point [True]
440 Text feature [deficiencies] present in test data point [True]
452 Text feature [inactivating] present in test data point [True]
473 Text feature [pms2] present in test data point [True]
485 Text feature [material] present in test data point [True]
Out of the top 500 features 12 are present in query point
```

4.3.1.3.2. Incorrectly Classified point

```
In [65]: 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]), 3))
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 : 8

Predicted Class Probabilities: [[0.107 0.1074 0.0566 0.0831 0.0171 0.002 0.2638 0.3446 0.0185]]

Actual Class : 8

```
-----
23 Text feature [g34v] present in test data point [True]
28 Text feature [schiffman] present in test data point [True]
30 Text feature [paediatric] present in test data point [True]
33 Text feature [k27] present in test data point [True]
34 Text feature [schwartzentruber] present in test data point [True]
35 Text feature [h3f3a] present in test data point [True]
38 Text feature [somers] present in test data point [True]
47 Text feature [permissiveness] present in test data point [True]
49 Text feature [g34r] present in test data point [True]
53 Text feature [thalamic] present in test data point [True]
59 Text feature [1004] present in test data point [True]
64 Text feature [wake] present in test data point [True]
66 Text feature [extensions] present in test data point [True]
71 Text feature [wolffe] present in test data point [True]
73 Text feature [foxg1] present in test data point [True]
.. - - - - -
```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

```

In [66]: 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=None,
5 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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-1
14 #-----
15
16
17
18 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
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(-6, 1)]
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, n_jobs=3)
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))

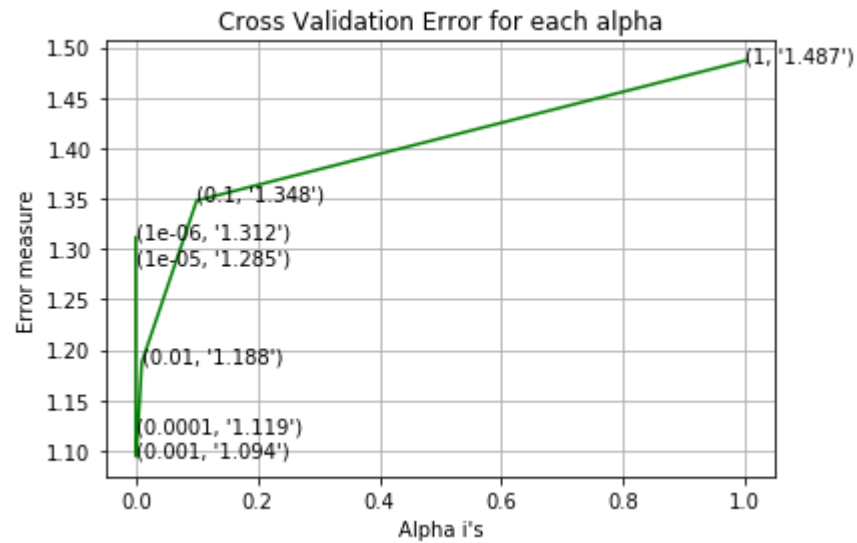
```

```

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, n_jobs=3)
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_
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,
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-06
Log Loss : 1.3115387654652857
for alpha = 1e-05
Log Loss : 1.28525478407198
for alpha = 0.0001
Log Loss : 1.1185234997053703
for alpha = 0.001
Log Loss : 1.0941479320488363
for alpha = 0.01
Log Loss : 1.1878002654430508
for alpha = 0.1
Log Loss : 1.3480381452471564
for alpha = 1
Log Loss : 1.4867946160872978

```



For values of best alpha = 0.001 The train log loss is: 0.5492333718423743

For values of best alpha = 0.001 The cross validation log loss is: 1.0941479320488363

For values of best alpha = 0.001 The test log loss is: 1.0838709350672702

4.3.2.2. Testing model with best hyper parameters

```

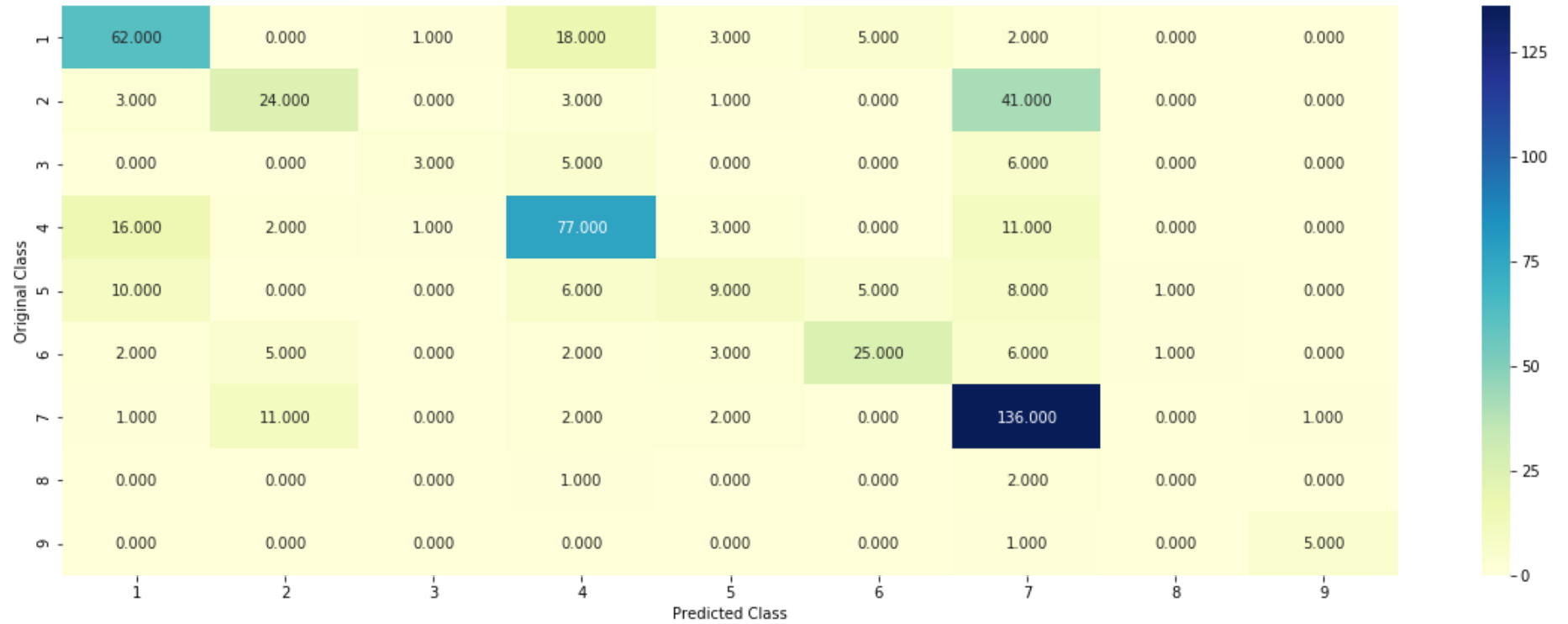
In [67]: 1 # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.
2 # -----
3 # default parameters
4 # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None,
5 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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, n_jobs=3)
17 predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

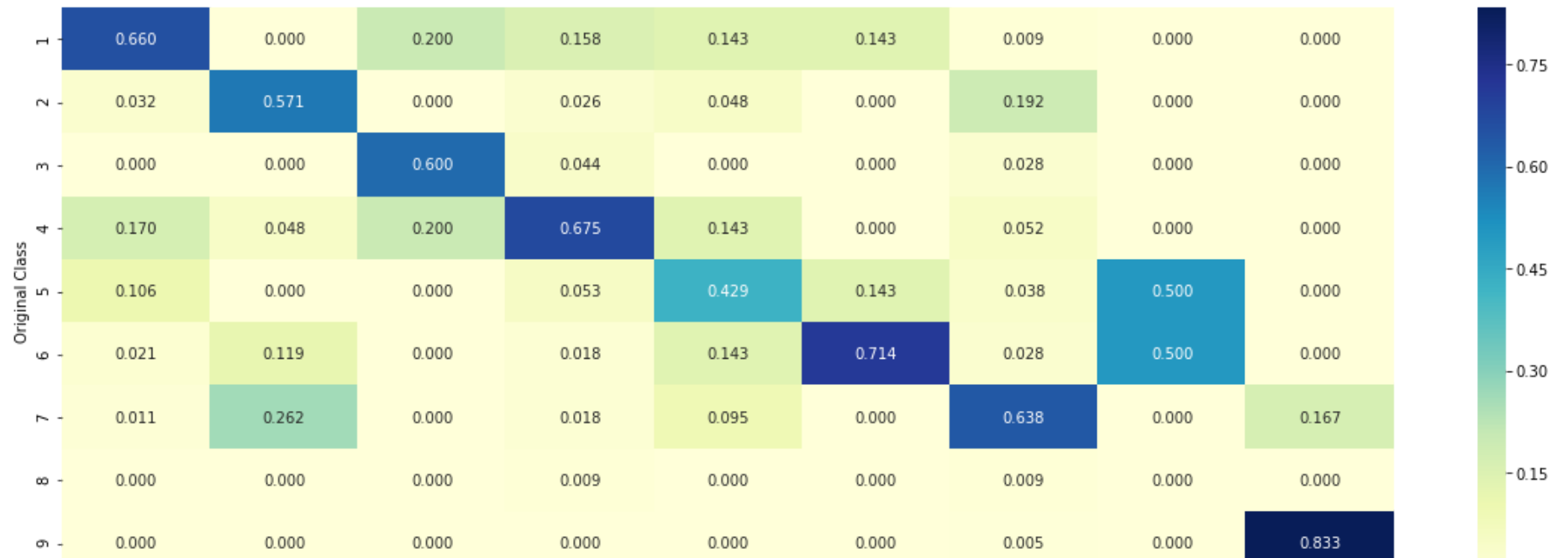
Log loss : 1.0941479320488363

Number of mis-classified points : 0.35902255639097747

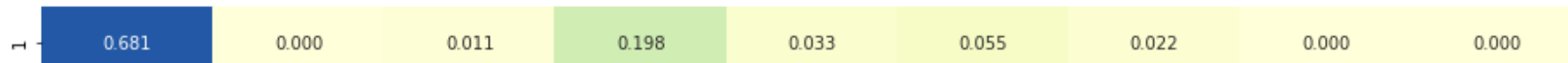
----- Confusion matrix -----



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



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



4.3.2.3. Feature Importance, Correctly Classified point

```
In [68]: 1 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42, n_jobs=3)
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_index]), 4))
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_index])
```

Predicted Class : 4

Predicted Class Probabilities: [[1.509e-01 1.410e-02 1.000e-02 6.896e-01 7.280e-02 1.370e-02 1.460e-02 3.420e-02 1.000e-04]]

Actual Class : 4

```
-----
99 Text feature [382] present in test data point [True]
144 Text feature [truncate] present in test data point [True]
177 Text feature [inactivation] present in test data point [True]
207 Text feature [adapted] present in test data point [True]
242 Text feature [deficiencies] present in test data point [True]
273 Text feature [instability] present in test data point [True]
333 Text feature [germline] present in test data point [True]
373 Text feature [pms2] present in test data point [True]
376 Text feature [hereditary] present in test data point [True]
407 Text feature [material] present in test data point [True]
412 Text feature [carriers] present in test data point [True]
489 Text feature [inactivating] present in test data point [True]
Out of the top 500 features 12 are present in query point
```

4.3.2.4. Feature Importance, Inorrectly Classified point

```
In [69]: 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: [[0.1685 0.099 0.003 0.1324 0.0099 0.0032 0.3391 0.2445 0.0004]]

Actual Class : 8

```
-----
45 Text feature [activated] present in test data point [True]
58 Text feature [constitutive] present in test data point [True]
99 Text feature [activation] present in test data point [True]
113 Text feature [constitutively] present in test data point [True]
114 Text feature [expressing] present in test data point [True]
120 Text feature [recurrence] present in test data point [True]
138 Text feature [transformation] present in test data point [True]
142 Text feature [downstream] present in test data point [True]
240 Text feature [agar] present in test data point [True]
243 Text feature [serum] present in test data point [True]
297 Text feature [conventional] present in test data point [True]
335 Text feature [hazards] present in test data point [True]
379 Text feature [oncogenesis] present in test data point [True]
381 Text feature [transform] present in test data point [True]
455 Text feature [phosphorylation] present in test data point [True]
```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

```

In [95]: 1 # read more about support vector machines with linear kernals here http://scikit-learn.org/stable/modules/g
2
3 # -----
4 # default parameters
5 # SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
6 # cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_stat
7
8 # Some of methods of SVM()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # -----
12 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivati
13 # -----
14
15
16
17 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
18 # -----
19 # default paramters
20 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
21 #
22 # some of the methods of CalibratedClassifierCV()
23 # fit(X, y[, sample_weight]) Fit the calibrated model
24 # get_params([deep]) Get parameters for this estimator.
25 # predict(X) Predict the target of new samples.
26 # predict_proba(X) Posterior probabilities of classification
27 #-----
28 # video link:
29 #-----
30
31 alpha = [10 ** x for x in range(-5, 3)]
32 cv_log_error_array = []
33 for i in alpha:
34     print("for C =", i)
35     # clf = SVC(C=i, kernel='linear', probability=True, class_weight='balanced')
36     clf = SGDClassifier( class_weight='balanced', alpha=i, penalty='l2', loss='hinge', random_state=42, n_j
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))

```

```

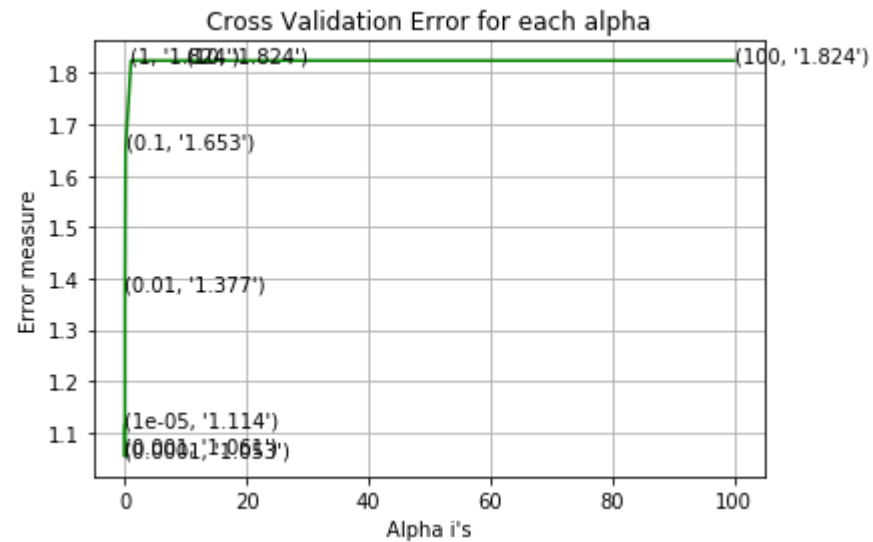
42     print("Log Loss :",log_loss(cv_y, sig_clf_probs))
43
44     fig, ax = plt.subplots()
45     ax.plot(alpha, cv_log_error_array,c='g')
46     for i, txt in enumerate(np.round(cv_log_error_array,3)):
47         ax.annotate((alpha[i],str(txt)), (alpha[i],cv_log_error_array[i]))
48     plt.grid()
49     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 = SVC(C=i,kernel='linear',probability=True, class_weight='balanced')
57     clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_st
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, predict_
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_cv,
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, predict_y,

for C = 1e-05
Log Loss : 1.1142159368824716
for C = 0.0001
Log Loss : 1.0530701415497956
for C = 0.001
Log Loss : 1.0607930820778635
for C = 0.01
Log Loss : 1.3770671358668958
for C = 0.1
Log Loss : 1.6533655406696115
for C = 1
Log Loss : 1.8243564312423943
for C = 10
Log Loss : 1.824353879737443

```

```
for C = 100
```

```
Log Loss : 1.8243536192153422
```



```
For values of best alpha = 0.0001 The train log loss is: 0.3836468326570109
```

```
For values of best alpha = 0.0001 The cross validation log loss is: 1.0530701415497956
```

```
For values of best alpha = 0.0001 The test log loss is: 1.0328584955264182
```

4.4.2. Testing model with best hyper parameters

```

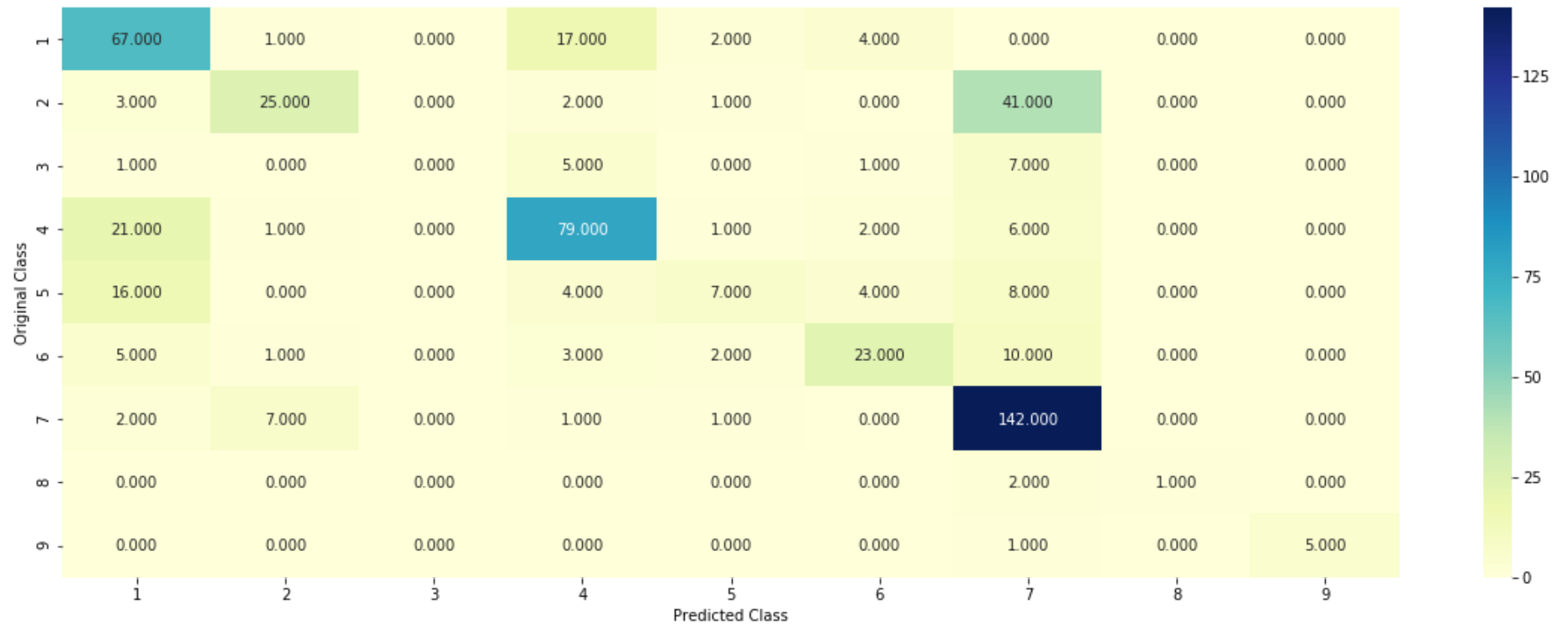
In [96]: 1 # read more about support vector machines with linear kernals here http://scikit-learn.org/stable/modules/g
2
3 # -----
4 # default parameters
5 # SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
6 # cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_stat
7
8 # Some of methods of SVM()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # -----
12 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivati
13 # -----
14
15
16 # clf = SVC(C=alpha[best_alpha],kernel='linear',probability=True, class_weight='balanced')
17 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42,class_weight='bala
18 predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf)

```

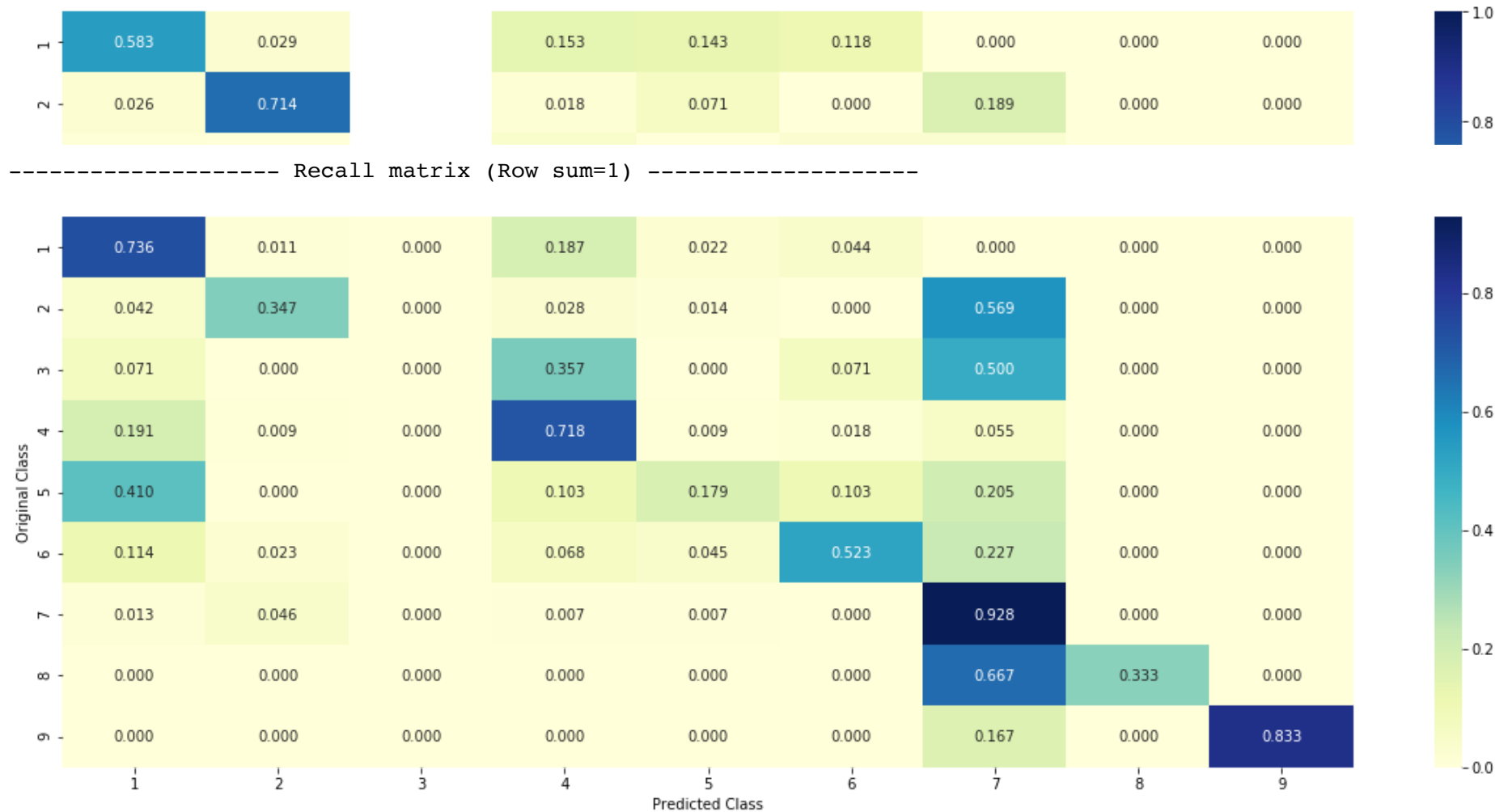
Log loss : 1.0530701415497956

Number of mis-classified points : 0.34398496240601506

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



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



4.3.3. Feature Importance

4.3.3.1. For Correctly classified point

```
In [97]: 1 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
2 clf.fit(train_x_onehotCoding,train_y)
3 test_point_index = 1
4 # test_point_index = 100
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_index]), 4))
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_index])
```

Predicted Class : 6

Predicted Class Probabilities: [[0.1798 0.0563 0.0237 0.1681 0.1823 0.3664 0.0152 0.0047 0.0035]]

Actual Class : 4

```
-----
112 Text feature [119] present in test data point [True]
284 Text feature [006] present in test data point [True]
287 Text feature [111] present in test data point [True]
379 Text feature [000249] present in test data point [True]
492 Text feature [05] present in test data point [True]
Out of the top 500 features 5 are present in query point
```

4.3.3.2. For Incorrectly classified point

```
In [98]: 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 : 8

Predicted Class Probabilities: [[0.1644 0.0823 0.0055 0.1492 0.0165 0.0088 0.0933 0.4707 0.0094]]

Actual Class : 8

```
-----
54 Text feature [100] present in test data point [True]
56 Text feature [1000] present in test data point [True]
246 Text feature [120] present in test data point [True]
312 Text feature [0016] present in test data point [True]
367 Text feature [10] present in test data point [True]
374 Text feature [113] present in test data point [True]
Out of the top 500 features 6 are present in query point
```

4.5 Random Forest Classifier

4.5.1. Hyper paramter tuning (With One hot Encoding)

```

In [0]: 1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_sp
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_st
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # predict_proba (X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-the
19 # -----
20
21
22 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
23 # -----
24 # default paramters
25 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
26 #
27 # some of the methods of CalibratedClassifierCV()
28 # fit(X, y[, sample_weight]) Fit the calibrated model
29 # get_params([deep]) Get parameters for this estimator.
30 # predict(X) Predict the target of new samples.
31 # predict_proba(X) Posterior probabilities of classification
32 #-----
33 # video link:
34 #-----
35
36 alpha = [100,200,500,1000,2000]
37 max_depth = [5, 10]
38 cv_log_error_array = []
39 for i in alpha:
40     for j in max_depth:
41         print("for n_estimators = ", i,"and max depth = ", j)

```

```

42     clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=
43     clf.fit(train_x_onehotCoding, train_y)
44     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
45     sig_clf.fit(train_x_onehotCoding, train_y)
46     sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
47     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
48     print("Log Loss :", log_loss(cv_y, sig_clf_probs))
49
50     '''fig, ax = plt.subplots()
51     features = np.dot(np.array(alpha)[: ,None], np.array(max_depth)[None]).ravel()
52     ax.plot(features, cv_log_error_array, c='g')
53     for i, txt in enumerate(np.round(cv_log_error_array,3)):
54         ax.annotate((alpha[int(i/2)], max_depth[int(i%2)], str(txt)), (features[i], cv_log_error_array[i]))
55     plt.grid()
56     plt.title("Cross Validation Error for each alpha")
57     plt.xlabel("Alpha i's")
58     plt.ylabel("Error measure")
59     plt.show()
60     '''
61
62     best_alpha = np.argmin(cv_log_error_array)
63     clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[i
64     clf.fit(train_x_onehotCoding, train_y)
65     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
66     sig_clf.fit(train_x_onehotCoding, train_y)
67
68     predict_y = sig_clf.predict_proba(train_x_onehotCoding)
69     print('For values of best estimator = ', alpha[int(best_alpha/2)], "The train log loss is:", log_loss(y_train,
70     predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
71     print('For values of best estimator = ', alpha[int(best_alpha/2)], "The cross validation log loss is:", log_
72     predict_y = sig_clf.predict_proba(test_x_onehotCoding)
73     print('For values of best estimator = ', alpha[int(best_alpha/2)], "The test log loss is:", log_loss(y_test,

```

```

for n_estimators = 100 and max depth = 5
Log Loss : 1.2572535683354957
for n_estimators = 100 and max depth = 10
Log Loss : 1.1868414223711878
for n_estimators = 200 and max depth = 5
Log Loss : 1.2378734502517341
for n_estimators = 200 and max depth = 10
Log Loss : 1.1811031780258958
for n_estimators = 500 and max depth = 5

```

```
Log Loss : 1.2368241894319212
for n_estimators = 500 and max depth = 10
Log Loss : 1.176754594516683
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2357829533963691
for n_estimators = 1000 and max depth = 10
Log Loss : 1.174993079576866
for n_estimators = 2000 and max depth = 5
Log Loss : 1.236042392554891
for n_estimators = 2000 and max depth = 10
Log Loss : 1.1759745074379755
For values of best estimator = 1000 The train log loss is: 0.7095396732082752
For values of best estimator = 1000 The cross validation log loss is: 1.174993079576866
For values of best estimator = 1000 The test log loss is: 1.1630923149103904
```

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

```

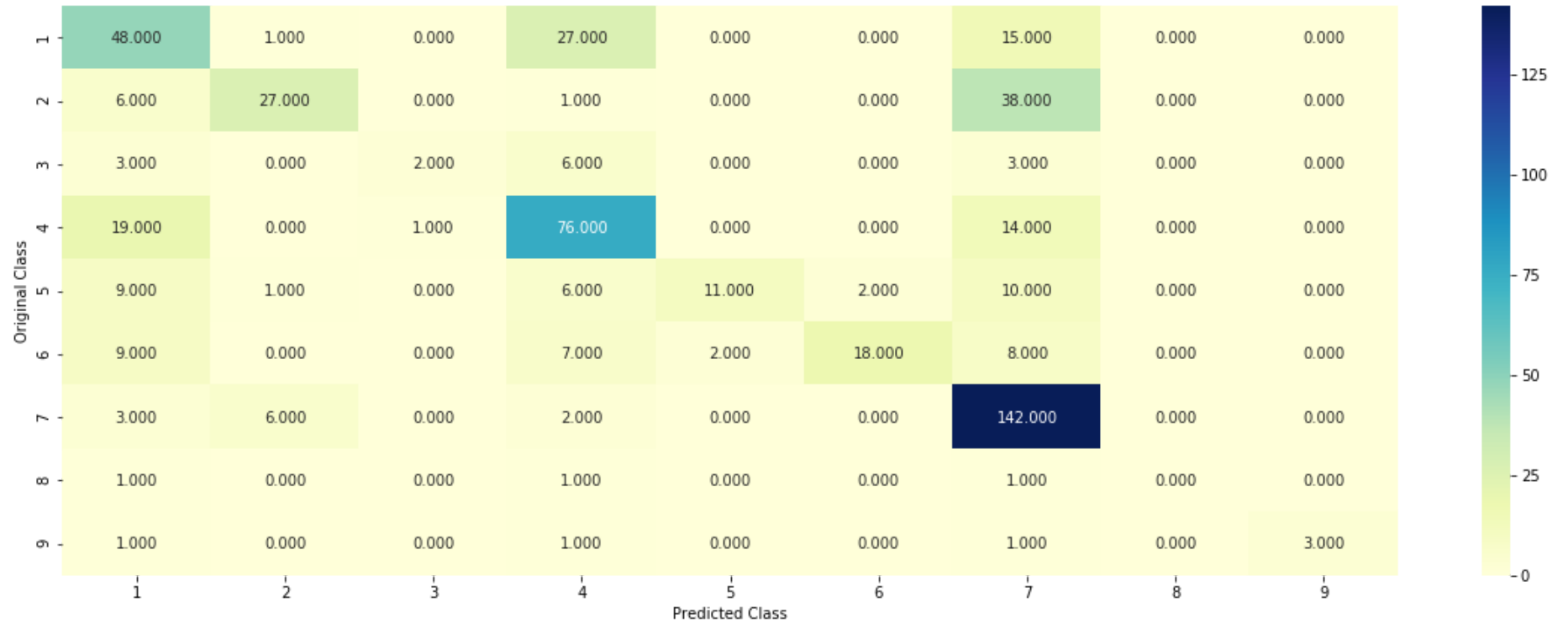
In [0]: 1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_sp
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_st
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight])    Fit the SVM model according to the given training data.
10 # predict(X)    Perform classification on samples in X.
11 # predict_proba (X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-the
19 # -----
20
21 clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[i
22 predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf)

```

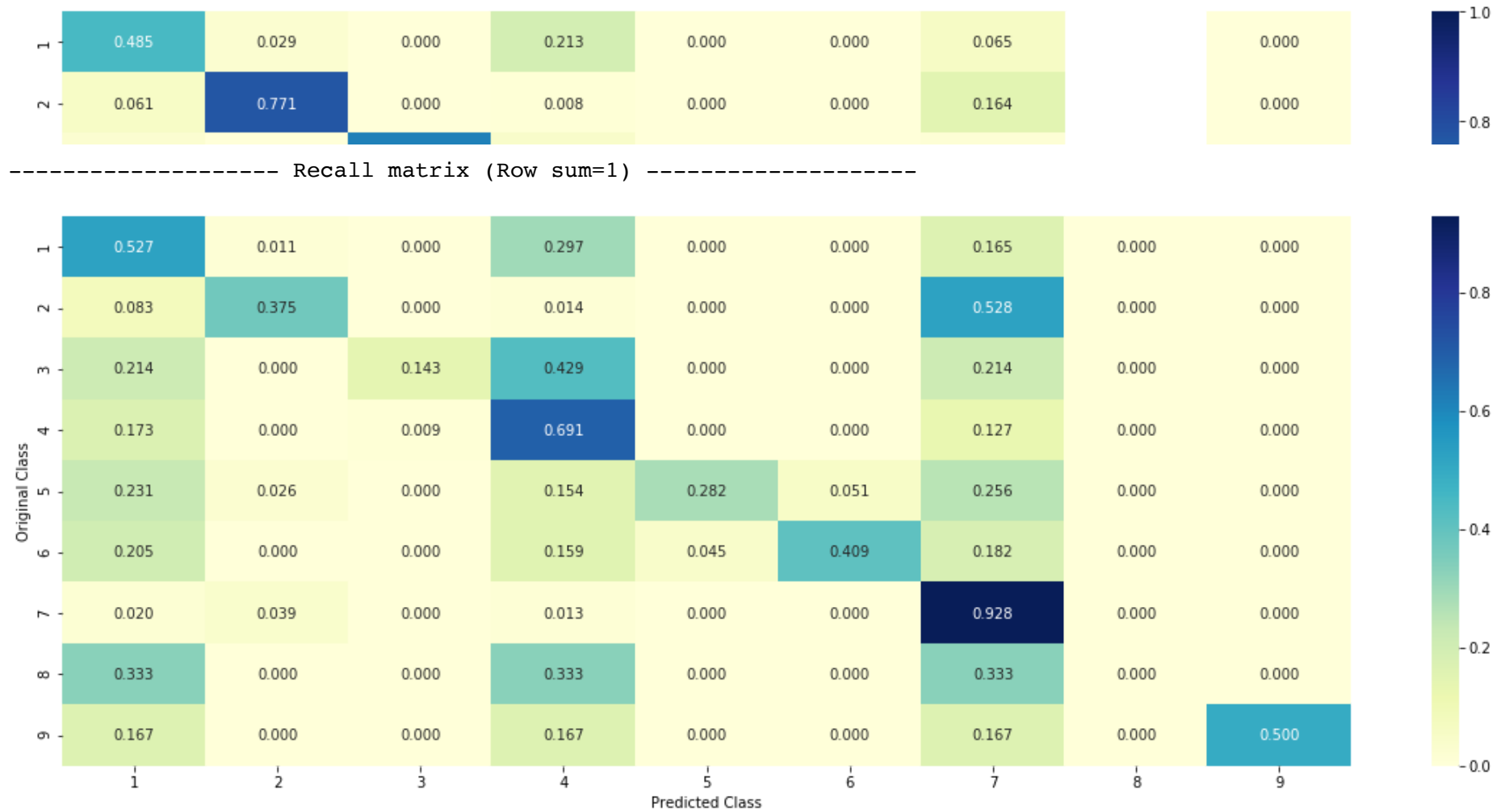
Log loss : 1.174993079576866

Number of mis-classified points : 0.38533834586466165

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



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



4.5.3. Feature Importance

4.5.3.1. Correctly Classified point

```
In [0]: 1 # test_point_index = 10
2 clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[i
3 clf.fit(train_x_onehotCoding, train_y)
4 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
5 sig_clf.fit(train_x_onehotCoding, train_y)
6
7 test_point_index = 1
8 no_feature = 100
9 predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
10 print("Predicted Class :", predicted_cls[0])
11 print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index
12 print("Actual Class :", test_y[test_point_index])
13 indices = np.argsort(-clf.feature_importances_)
14 print("-"*50)
15 get_impfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0454 0.1404 0.0133 0.029 0.036 0.0294 0.6977 0.005 0.004]]

Actual Class : 7

```
-----
0 Text feature [inhibitors] present in test data point [True]
1 Text feature [kinase] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [tyrosine] present in test data point [True]
4 Text feature [missense] present in test data point [True]
5 Text feature [inhibitor] present in test data point [True]
7 Text feature [treatment] present in test data point [True]
8 Text feature [oncogenic] present in test data point [True]
9 Text feature [suppressor] present in test data point [True]
10 Text feature [activation] present in test data point [True]
11 Text feature [phosphorylation] present in test data point [True]
12 Text feature [kinases] present in test data point [True]
13 Text feature [nonsense] present in test data point [True]
14 Text feature [akt] present in test data point [True]
15 Text feature [function] present in test data point [True]
17 Text feature [erk] present in test data point [True]
19 Text feature [growth] present in test data point [True]
20 Text feature [variants] present in test data point [True]
22 Text feature [frameshift] present in test data point [True]
24 Text feature [therapeutic] present in test data point [True]
25 Text feature [functional] present in test data point [True]
```

28 Text feature [signaling] present in test data point [True]
30 Text feature [patients] present in test data point [True]
31 Text feature [cells] present in test data point [True]
32 Text feature [constitutively] present in test data point [True]
34 Text feature [trials] present in test data point [True]
35 Text feature [therapy] present in test data point [True]
37 Text feature [erk1] present in test data point [True]
38 Text feature [activate] present in test data point [True]
39 Text feature [downstream] present in test data point [True]
41 Text feature [efficacy] present in test data point [True]
42 Text feature [protein] present in test data point [True]
43 Text feature [loss] present in test data point [True]
44 Text feature [inhibited] present in test data point [True]
45 Text feature [expressing] present in test data point [True]
46 Text feature [pten] present in test data point [True]
48 Text feature [lines] present in test data point [True]
49 Text feature [treated] present in test data point [True]
50 Text feature [proliferation] present in test data point [True]
51 Text feature [drug] present in test data point [True]
57 Text feature [mek] present in test data point [True]
59 Text feature [inhibition] present in test data point [True]
61 Text feature [repair] present in test data point [True]
62 Text feature [sensitivity] present in test data point [True]
64 Text feature [receptor] present in test data point [True]
66 Text feature [assays] present in test data point [True]
68 Text feature [survival] present in test data point [True]
69 Text feature [cell] present in test data point [True]
71 Text feature [ligand] present in test data point [True]
73 Text feature [expression] present in test data point [True]
74 Text feature [variant] present in test data point [True]
75 Text feature [oncogene] present in test data point [True]
78 Text feature [extracellular] present in test data point [True]
79 Text feature [doses] present in test data point [True]
80 Text feature [mapk] present in test data point [True]
81 Text feature [hours] present in test data point [True]
84 Text feature [information] present in test data point [True]
86 Text feature [harboring] present in test data point [True]
90 Text feature [dna] present in test data point [True]
91 Text feature [concentrations] present in test data point [True]
92 Text feature [likelihood] present in test data point [True]
93 Text feature [months] present in test data point [True]
94 Text feature [binding] present in test data point [True]

```
96 Text feature [imatinib] present in test data point [True]
98 Text feature [preclinical] present in test data point [True]
Out of the top 100 features 65 are present in query point
```

4.5.3.2. Inorrectly Classified point

```
In [0]: 1 test_point_index = 100
2 no_feature = 100
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]), 2))
6 print("Actual Class :", test_y[test_point_index])
7 indices = np.argsort(-clf.feature_importances_)
8 print("-"*50)
9 get_impfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_index])
```

Predicted Class : 7

Predicted Class Probabilities: [[0.1337 0.116 0.0224 0.1773 0.0674 0.0545 0.4156 0.0071 0.0059]]

Actual Class : 7

```
-----
0 Text feature [inhibitors] present in test data point [True]
1 Text feature [kinase] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [tyrosine] present in test data point [True]
6 Text feature [activated] present in test data point [True]
8 Text feature [oncogenic] present in test data point [True]
10 Text feature [activation] present in test data point [True]
11 Text feature [phosphorylation] present in test data point [True]
12 Text feature [kinases] present in test data point [True]
14 Text feature [akt] present in test data point [True]
15 Text feature [function] present in test data point [True]
19 Text feature [growth] present in test data point [True]
21 Text feature [constitutive] present in test data point [True]
25 Text feature [functional] present in test data point [True]
28 Text feature [signaling] present in test data point [True]
31 Text feature [cells] present in test data point [True]
32 Text feature [constitutively] present in test data point [True]
38 Text feature [activate] present in test data point [True]
39 Text feature [downstream] present in test data point [True]
42 Text feature [protein] present in test data point [True]
43 Text feature [loss] present in test data point [True]
44 Text feature [inhibited] present in test data point [True]
46 Text feature [pten] present in test data point [True]
47 Text feature [transforming] present in test data point [True]
48 Text feature [lines] present in test data point [True]
50 Text feature [proliferation] present in test data point [True]
53 Text feature [neutral] present in test data point [True]
```

```
55 Text feature [transform] present in test data point [True]
56 Text feature [stability] present in test data point [True]
58 Text feature [transformation] present in test data point [True]
59 Text feature [inhibition] present in test data point [True]
62 Text feature [sensitivity] present in test data point [True]
64 Text feature [receptor] present in test data point [True]
66 Text feature [assays] present in test data point [True]
69 Text feature [cell] present in test data point [True]
75 Text feature [oncogene] present in test data point [True]
84 Text feature [information] present in test data point [True]
90 Text feature [dna] present in test data point [True]
94 Text feature [binding] present in test data point [True]
Out of the top 100 features 39 are present in query point
```

4.5.3. Hyper paramter tuning (With Response Coding)


```

In [0]: 1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_sp
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_st
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # predict_proba (X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-the
19 # -----
20
21
22 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.c
23 # -----
24 # default paramters
25 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
26 #
27 # some of the methods of CalibratedClassifierCV()
28 # fit(X, y[, sample_weight]) Fit the calibrated model
29 # get_params([deep]) Get parameters for this estimator.
30 # predict(X) Predict the target of new samples.
31 # predict_proba(X) Posterior probabilities of classification
32 #-----
33 # video link:
34 #-----
35
36 alpha = [10,50,100,200,500,1000]
37 max_depth = [2,3,5,10]
38 cv_log_error_array = []
39 for i in alpha:
40     for j in max_depth:
41         print("for n_estimators = ", i,"and max depth = ", j)

```

```

42     clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=
43     clf.fit(train_x_responseCoding, train_y)
44     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
45     sig_clf.fit(train_x_responseCoding, train_y)
46     sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
47     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
48     print("Log Loss :", log_loss(cv_y, sig_clf_probs))
49     '''
50     fig, ax = plt.subplots()
51     features = np.dot(np.array(alpha)[: ,None], np.array(max_depth)[None]).ravel()
52     ax.plot(features, cv_log_error_array, c='g')
53     for i, txt in enumerate(np.round(cv_log_error_array,3)):
54         ax.annotate((alpha[int(i/4)], max_depth[int(i%4)], str(txt)), (features[i], cv_log_error_array[i]))
55     plt.grid()
56     plt.title("Cross Validation Error for each alpha")
57     plt.xlabel("Alpha i's")
58     plt.ylabel("Error measure")
59     plt.show()
60     '''
61
62     best_alpha = np.argmin(cv_log_error_array)
63     clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_depth=max_depth[i
64     clf.fit(train_x_responseCoding, train_y)
65     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
66     sig_clf.fit(train_x_responseCoding, train_y)
67
68     predict_y = sig_clf.predict_proba(train_x_responseCoding)
69     print('For values of best alpha = ', alpha[int(best_alpha/4)], "The train log loss is:", log_loss(y_train, p
70     predict_y = sig_clf.predict_proba(cv_x_responseCoding)
71     print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log loss is:", log_loss
72     predict_y = sig_clf.predict_proba(test_x_responseCoding)
73     print('For values of best alpha = ', alpha[int(best_alpha/4)], "The test log loss is:", log_loss(y_test, pre

```

```

for n_estimators = 10 and max depth = 2
Log Loss : 2.2657048897349608
for n_estimators = 10 and max depth = 3
Log Loss : 1.7459205010556096
for n_estimators = 10 and max depth = 5
Log Loss : 1.4368353925512503
for n_estimators = 10 and max depth = 10
Log Loss : 1.904597809032912
for n_estimators = 50 and max depth = 2

```

```
Log Loss : 1.7221951095007484
for n_estimators = 50 and max depth = 3
Log Loss : 1.4984825877845531
for n_estimators = 50 and max depth = 5
Log Loss : 1.4593628982873716
for n_estimators = 50 and max depth = 10
Log Loss : 1.8434939703555409
for n_estimators = 100 and max depth = 2
Log Loss : 1.6182209245331227
for n_estimators = 100 and max depth = 3
Log Loss : 1.5199297988828253
for n_estimators = 100 and max depth = 5
Log Loss : 1.4177501184246677
for n_estimators = 100 and max depth = 10
Log Loss : 1.8227504417195126
for n_estimators = 200 and max depth = 2
Log Loss : 1.6622571648074496
for n_estimators = 200 and max depth = 3
Log Loss : 1.4800771339141767
for n_estimators = 200 and max depth = 5
Log Loss : 1.4412060242341358
for n_estimators = 200 and max depth = 10
Log Loss : 1.7892406351442258
for n_estimators = 500 and max depth = 2
Log Loss : 1.715950314170445
for n_estimators = 500 and max depth = 3
Log Loss : 1.5658682738699774
for n_estimators = 500 and max depth = 5
Log Loss : 1.4445360301518217
for n_estimators = 500 and max depth = 10
Log Loss : 1.8421097596928397
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6834927870864949
for n_estimators = 1000 and max depth = 3
Log Loss : 1.5631973035931377
for n_estimators = 1000 and max depth = 5
Log Loss : 1.4449980792724129
for n_estimators = 1000 and max depth = 10
Log Loss : 1.85233132619749
For values of best alpha = 100 The train log loss is: 0.060702709444608406
For values of best alpha = 100 The cross validation log loss is: 1.417750118424668
For values of best alpha = 100 The test log loss is: 1.3806278998341923
```

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

```

In [0]: 1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_sp
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_st
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight])    Fit the SVM model according to the given training data.
10 # predict(X)    Perform classification on samples in X.
11 # predict_proba (X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-the
19 # -----
20
21 clf = RandomForestClassifier(max_depth=max_depth[int(best_alpha%4)], n_estimators=alpha[int(best_alpha/4)],
22 predict_and_plot_confusion_matrix(train_x_responseCoding, train_y,cv_x_responseCoding,cv_y, clf)

```

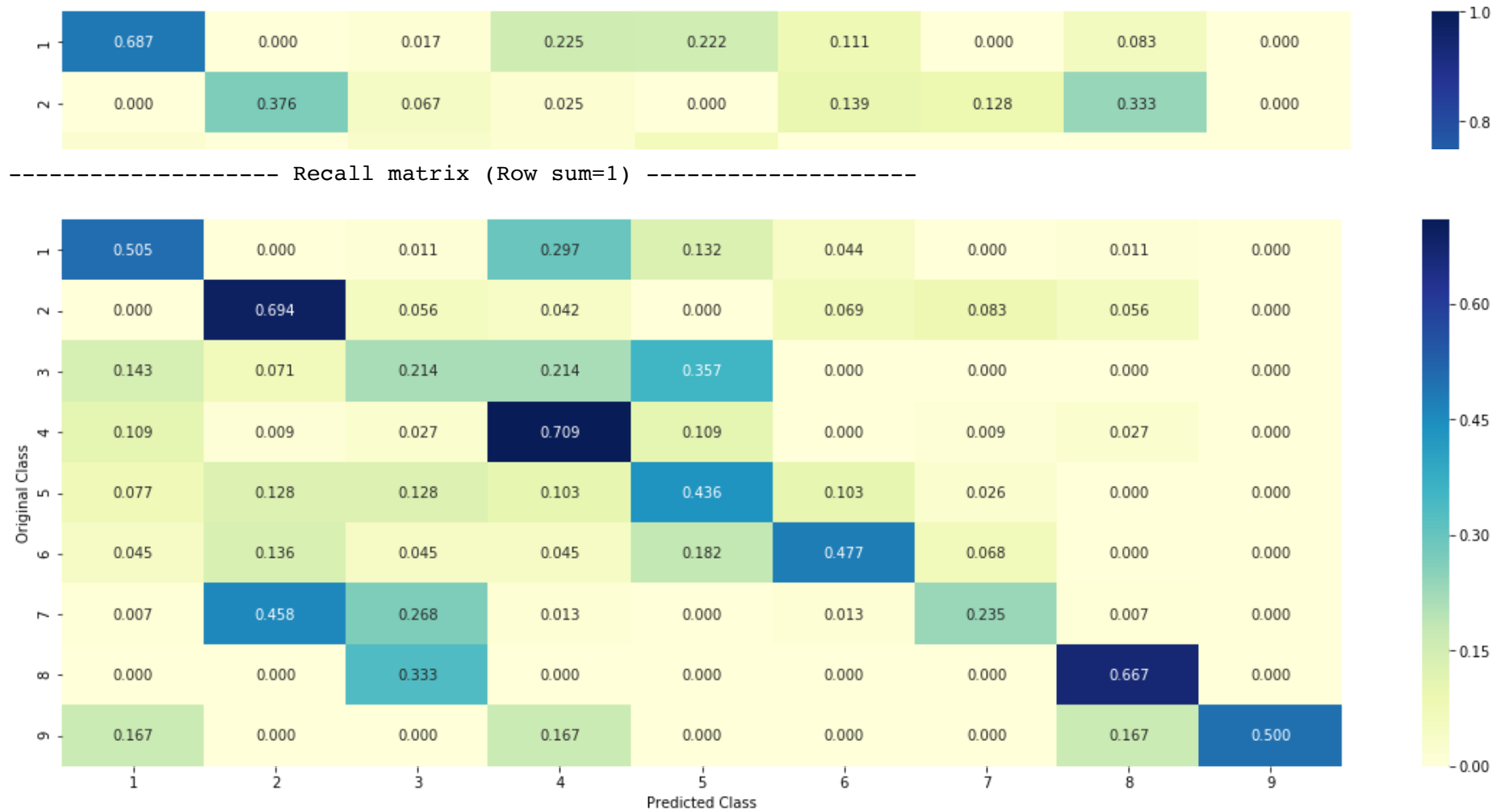
Log loss : 1.4177501184246677

Number of mis-classified points : 0.518796992481203

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



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



4.5.5. Feature Importance

4.5.5.1. Correctly Classified point


```

In [0]: 1 clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_depth=max_depth[i]
2 clf.fit(train_x_responseCoding, train_y)
3 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
4 sig_clf.fit(train_x_responseCoding, train_y)
5
6
7 test_point_index = 1
8 no_feature = 27
9 predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
10 print("Predicted Class :", predicted_cls[0])
11 print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_responseCoding[test_point_index]).r
12 print("Actual Class :", test_y[test_point_index])
13 indices = np.argsort(-clf.feature_importances_)
14 print("-"*50)
15 for i in indices:
16     if i<9:
17         print("Gene is important feature")
18     elif i<18:
19         print("Variation is important feature")
20     else:
21         print("Text is important feature")

```

Predicted Class : 2

Predicted Class Probabilities: [[0.0143 0.5044 0.1471 0.0191 0.0245 0.065 0.1724 0.039 0.0142]]

Actual Class : 7

```

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

```

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

4.5.5.2. Incorrectly Classified point

```
In [0]: 1 test_point_index = 100
2 predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
3 print("Predicted Class :", predicted_cls[0])
4 print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_responseCoding[test_point_index]).reshape(1,-1), 4))
5 print("Actual Class :", test_y[test_point_index])
6 indices = np.argsort(-clf.feature_importances_)
7 print("-"*50)
8 for i in indices:
9     if i<9:
10         print("Gene is important feature")
11     elif i<18:
12         print("Variation is important feature")
13     else:
14         print("Text is important feature")
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0281 0.2006 0.203 0.0857 0.0626 0.0906 0.2249 0.0676 0.0369]]

Actual Class : 7

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

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

4.7 Stack the models

4.7.1 testing with hyper parameter tuning

```

In [0]: 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=None,
5 # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, pow
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-1
14 #-----
15
16
17 # read more about support vector machines with linear kernels here http://scikit-learn.org/stable/modules/g
18 # -----
19 # default parameters
20 # SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
21 # cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_stat
22
23 # Some of methods of SVM()
24 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
25 # predict(X) Perform classification on samples in X.
26 # -----
27 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivati
28 # -----
29
30
31 # read more about support vector machines with linear kernels here http://scikit-learn.org/stable/modules/g
32 # -----
33 # default parameters
34 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_sp
35 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_
36 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_st
37 # class_weight=None)
38
39 # Some of methods of RandomForestClassifier()
40 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
41 # predict(X) Perform classification on samples in X.

```

```

42 # predict_proba (X) Perform classification on samples in X.
43
44 # some of attributes of RandomForestClassifier()
45 # feature_importances_ : array of shape = [n_features]
46 # The feature importances (the higher, the more important the feature).
47
48 # -----
49 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-the
50 # -----
51
52
53 clf1 = SGDClassifier(alpha=0.001, penalty='l2', loss='log', class_weight='balanced', random_state=0)
54 clf1.fit(train_x_onehotCoding, train_y)
55 sig_clf1 = CalibratedClassifierCV(clf1, method="sigmoid")
56
57 clf2 = SGDClassifier(alpha=1, penalty='l2', loss='hinge', class_weight='balanced', random_state=0)
58 clf2.fit(train_x_onehotCoding, train_y)
59 sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")
60
61
62 clf3 = MultinomialNB(alpha=0.001)
63 clf3.fit(train_x_onehotCoding, train_y)
64 sig_clf3 = CalibratedClassifierCV(clf3, method="sigmoid")
65
66 sig_clf1.fit(train_x_onehotCoding, train_y)
67 print("Logistic Regression : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf1.predict_proba(cv_x_onehotCoding)))
68 sig_clf2.fit(train_x_onehotCoding, train_y)
69 print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf2.predict_proba(cv_x_onehotCoding)))
70 sig_clf3.fit(train_x_onehotCoding, train_y)
71 print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predict_proba(cv_x_onehotCoding))))
72 print("-"*50)
73 alpha = [0.0001,0.001,0.01,0.1,1,10]
74 best_alpha = 999
75 for i in alpha:
76     lr = LogisticRegression(C=i)
77     sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probabilities=True)
78     sclf.fit(train_x_onehotCoding, train_y)
79     print("Stacking Classifier : for the value of alpha: %f Log Loss: %0.3f" % (i, log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))))
80     log_error = log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))
81     if best_alpha > log_error:
82         best_alpha = log_error

```

Logistic Regression : Log Loss: 1.24

```
Logistic Regression : Log Loss: 1.24
Support vector machines : Log Loss: 1.72
Naive Bayes : Log Loss: 1.37
-----
Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.179
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.049
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.577
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.224
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.366
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.690
```

4.7.2 testing the model with the best hyper parameters

```
In [0]: 1 lr = LogisticRegression(C=0.1)
2 scf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probas=True)
3 scf.fit(train_x_onehotCoding, train_y)
4
5 log_error = log_loss(train_y, scf.predict_proba(train_x_onehotCoding))
6 print("Log loss (train) on the stacking classifier :",log_error)
7
8 log_error = log_loss(cv_y, scf.predict_proba(cv_x_onehotCoding))
9 print("Log loss (CV) on the stacking classifier :",log_error)
10
11 log_error = log_loss(test_y, scf.predict_proba(test_x_onehotCoding))
12 print("Log loss (test) on the stacking classifier :",log_error)
13
14 print("Number of missclassified point :", np.count_nonzero((scf.predict(test_x_onehotCoding)- test_y))/tes
15 plot_confusion_matrix(test_y=test_y, predict_y=scf.predict(test_x_onehotCoding))
```

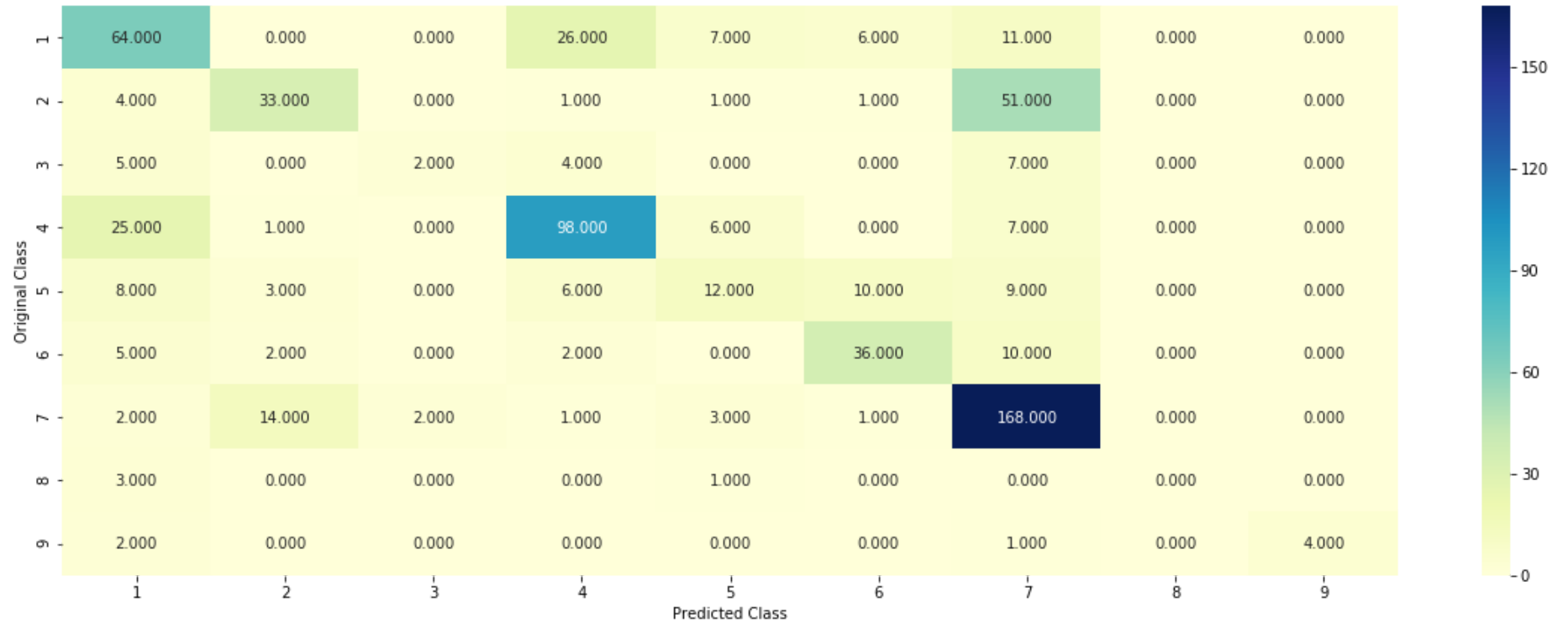
Log loss (train) on the stacking classifier : 0.6760284396805781

Log loss (CV) on the stacking classifier : 1.2243084610674686

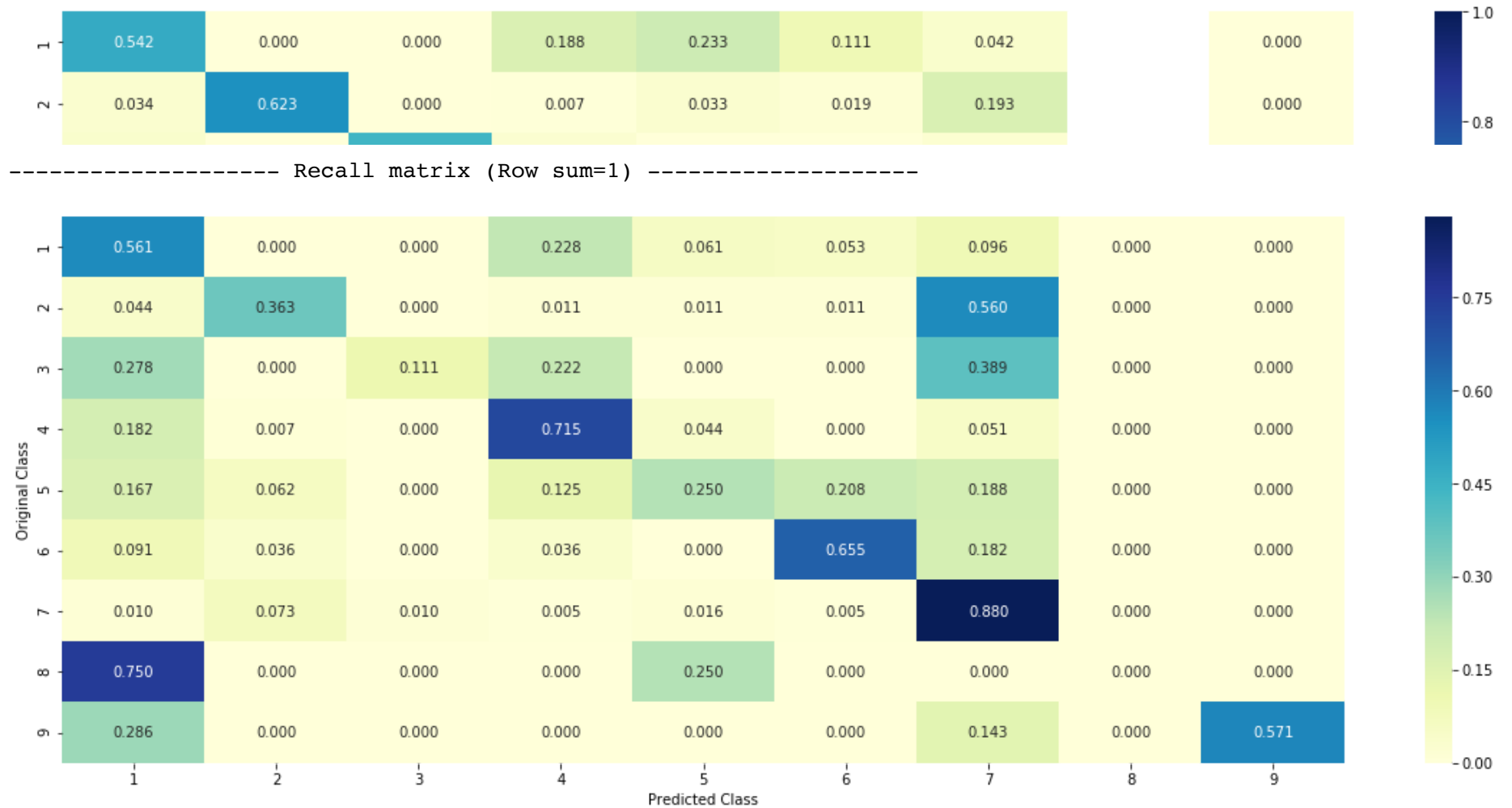
Log loss (test) on the stacking classifier : 1.1562525475350196

Number of missclassified point : 0.37293233082706767

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



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



4.7.3 Maximum Voting classifier

```
In [0]: 1 #Refer:http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingClassifier.html
2 from sklearn.ensemble import VotingClassifier
3 vclf = VotingClassifier(estimators=[('lr', sig_clf1), ('svc', sig_clf2), ('rf', sig_clf3)], voting='soft')
4 vclf.fit(train_x_onehotCoding, train_y)
5 print("Log loss (train) on the VotingClassifier :", log_loss(train_y, vclf.predict_proba(train_x_onehotCoding)))
6 print("Log loss (CV) on the VotingClassifier :", log_loss(cv_y, vclf.predict_proba(cv_x_onehotCoding)))
7 print("Log loss (test) on the VotingClassifier :", log_loss(test_y, vclf.predict_proba(test_x_onehotCoding)))
8 print("Number of missclassified point :", np.count_nonzero((vclf.predict(test_x_onehotCoding) - test_y))/test_y.size)
9 plot_confusion_matrix(test_y=test_y, predict_y=vclf.predict(test_x_onehotCoding))
```

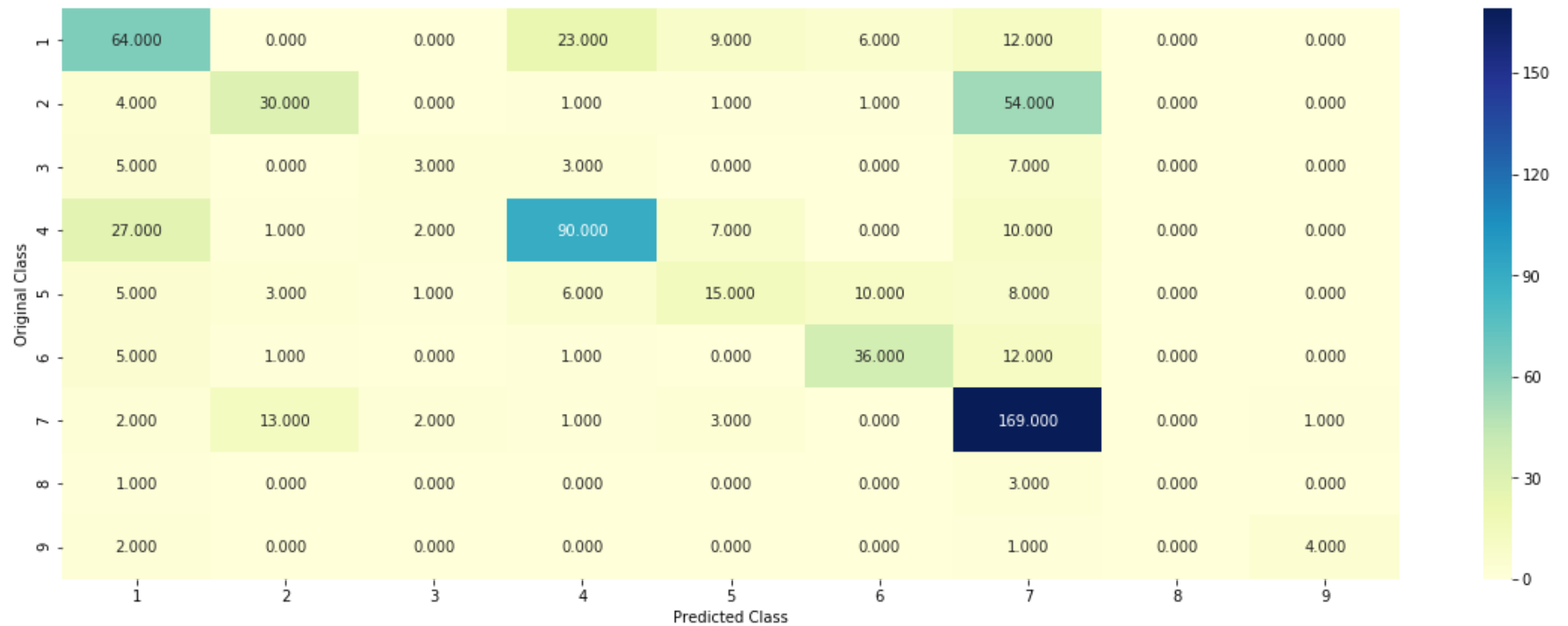
Log loss (train) on the VotingClassifier : 0.9407598679043604

Log loss (CV) on the VotingClassifier : 1.2835402100341697

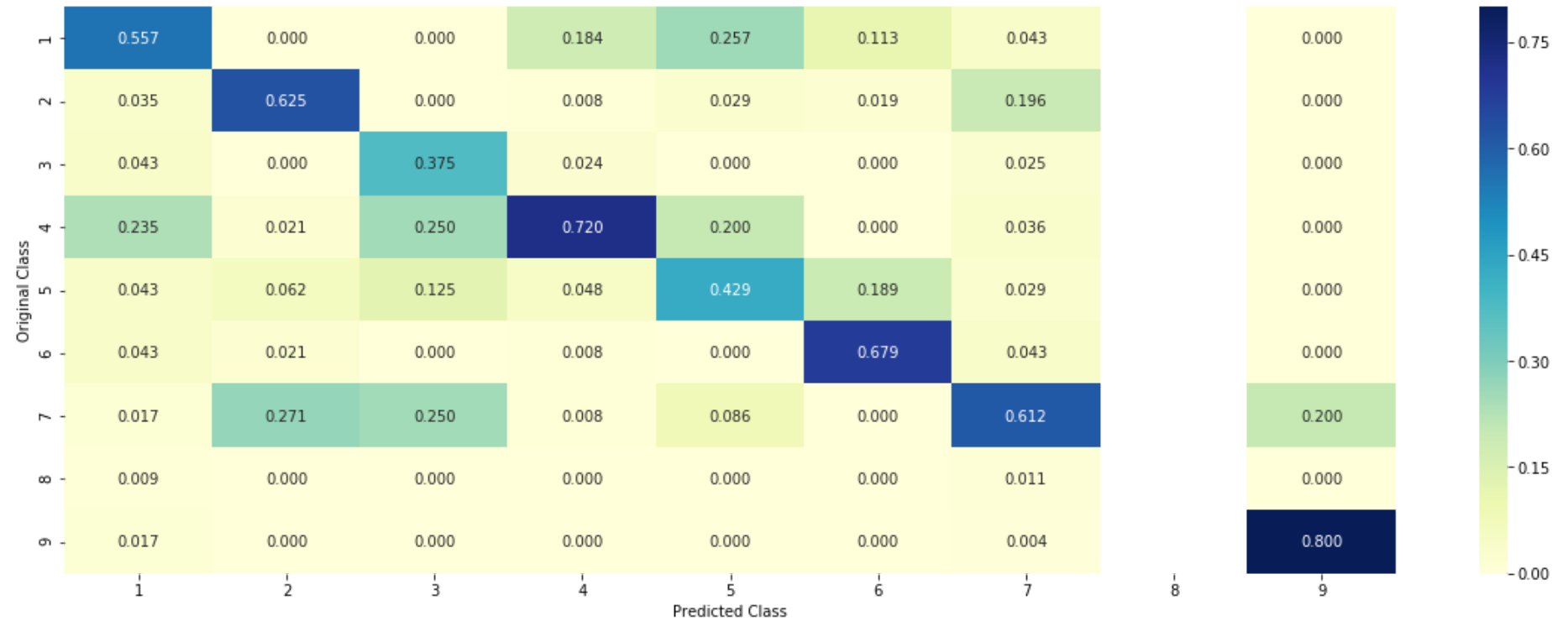
Log loss (test) on the VotingClassifier : 1.223278167176945

Number of missclassified point : 0.3819548872180451

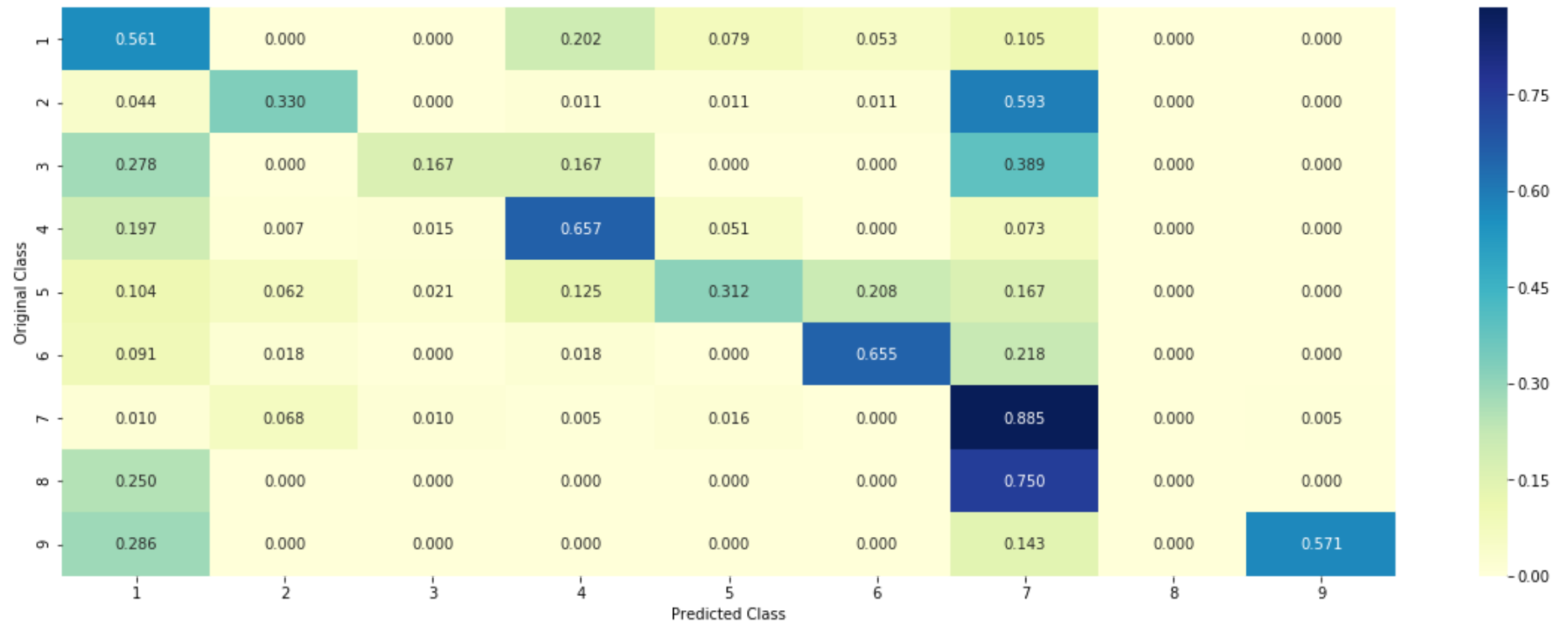
----- Confusion matrix -----



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



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



5. Assignments

1. Apply All the models with tf-idf features (Replace CountVectorizer with tfidfVectorizer and run the same cells)
2. Instead of using all the words in the dataset, use only the top 1000 words based of tf-idf values
3. Apply Logistic regression with CountVectorizer Features, including both unigrams and bigrams
4. Try any of the feature engineering techniques discussed in the course to reduce the CV and test log-loss to a value less than 1.0

```
In [99]: 1 feature_array = np.array(text_vectorizer.get_feature_names())
          2 tfidf_sorting = np.argsort(train_text_feature_onehotCoding.toarray()).flatten()[::-1]
          3
          4 n = 1000
          5 top_n = feature_array[tfidf_sorting][:n]
```

```
In [104]: 1 train_text_feature_tfidf_top_1000 = train_text_feature_onehotCoding[:, np.argsort(text_vectorizer.idf_)[:1000]]
```

```
In [106]: 1 train_text_feature_tfidf_top_1000 = normalize(train_text_feature_tfidf_top_1000)
```

```
In [107]: 1 train_text_feature_tfidf_top_1000
```

```
Out[107]: <2124x1000 sparse matrix of type '<class 'numpy.float64'>'
          with 1130862 stored elements in Compressed Sparse Row format>
```

```
In [109]: 1 test_text_feature_tfidf_top_1000 = text_vectorizer.transform(test_df['TEXT'])[:, np.argsort(text_vectorizer.idf_)[:1000]]
          2 test_text_feature_tfidf_top_1000 = normalize(test_text_feature_tfidf_top_1000)
          3
          4 cv_text_feature_tfidf_top_1000 = text_vectorizer.transform(cv_df['TEXT'])[:, np.argsort(text_vectorizer.idf_)[:1000]]
          5 cv_text_feature_tfidf_top_1000 = normalize(cv_text_feature_tfidf_top_1000)
```

```
In [110]: 1 test_text_feature_tfidf_top_1000
```

```
Out[110]: <665x1000 sparse matrix of type '<class 'numpy.float64'>'
          with 349459 stored elements in Compressed Sparse Row format>
```

```
In [126]: 1
          2 cv_text_feature_tfidf_top_1000
```

```
Out[126]: <532x1000 sparse matrix of type '<class 'numpy.float64'>'
          with 284687 stored elements in Compressed Sparse Row format>
```

```
In [112]: 1 # Stack gene, variation and text features horizontally
          2 train_x_onehotCoding_tfidf_top_1000 = hstack([train_gene_var_onehotCoding, train_text_feature_tfidf_top_1000])
          3
          4 test_x_onehotCoding_tfidf_top_1000 = hstack((test_gene_var_onehotCoding, test_text_feature_tfidf_top_1000))
          5
          6 cv_x_onehotCoding_tfidf_top_1000 = hstack((cv_gene_var_onehotCoding, cv_text_feature_tfidf_top_1000))
```

```
In [113]: 1 train_y = np.array(list(train_df['Class']))
          2 test_y=np.array(list(test_df['Class']))
          3 cv_y=np.array(list(cv_df['Class']))
```

```

In [114]: 1 alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100,1000]
          2 cv_log_error_array = []
          3 for i in alpha:
          4     print("for alpha =", i)
          5     clf = MultinomialNB(alpha=i)
          6     clf.fit(train_x_onehotCoding_tfidf_top_1000, train_y)
          7     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
          8     sig_clf.fit(train_x_onehotCoding_tfidf_top_1000, train_y)
          9     sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding_tfidf_top_1000)
         10     cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
         11     # to avoid rounding error while multiplying probabillites we use log-probability estimates
         12     print("Log Loss :",log_loss(cv_y, sig_clf_probs, eps=1e-15))
         13
         14 fig, ax = plt.subplots()
         15 ax.plot(np.log10(alpha), cv_log_error_array,c='g')
         16 for i, txt in enumerate(np.round(cv_log_error_array,3)):
         17     ax.annotate((alpha[i],str(txt)), (np.log10(alpha[i]),cv_log_error_array[i]))
         18 plt.grid()
         19 plt.xticks(np.log10(alpha))
         20 plt.title("Cross Validation Error for each alpha")
         21 plt.xlabel("Alpha i's")
         22 plt.ylabel("Error measure")
         23 plt.show()
         24
         25
         26 best_alpha = np.argmin(cv_log_error_array)
         27 clf = MultinomialNB(alpha=alpha[best_alpha])
         28 clf.fit(train_x_onehotCoding_tfidf_top_1000, train_y)
         29 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
         30 sig_clf.fit(train_x_onehotCoding_tfidf_top_1000, train_y)
         31
         32
         33 predict_y = sig_clf.predict_proba(train_x_onehotCoding_tfidf_top_1000)
         34 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_
         35 predict_y = sig_clf.predict_proba(cv_x_onehotCoding_tfidf_top_1000)
         36 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv,
         37 predict_y = sig_clf.predict_proba(test_x_onehotCoding_tfidf_top_1000)
         38 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y,
         39

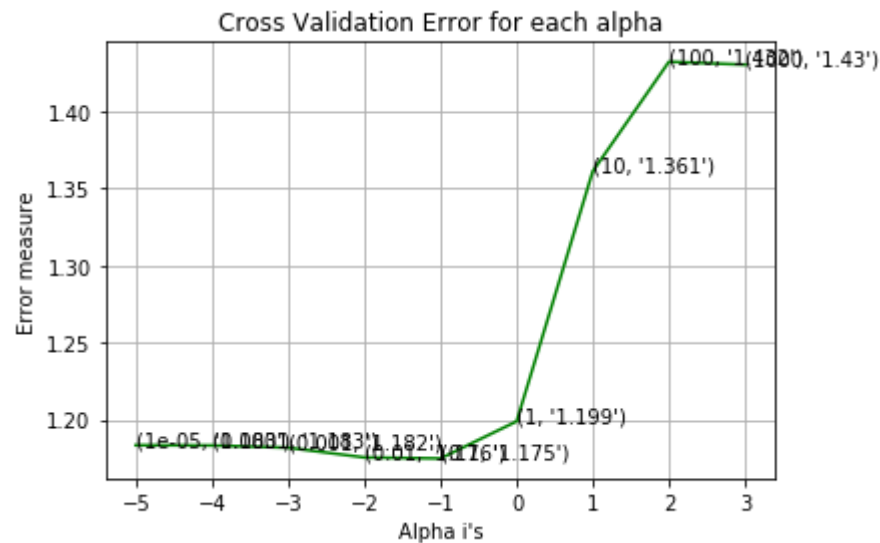
```

for alpha = 1e-05

```

Log Loss : 1.1834303887673983
for alpha = 0.0001
Log Loss : 1.183267843927018
for alpha = 0.001
Log Loss : 1.181819646655514
for alpha = 0.01
Log Loss : 1.175509324129162
for alpha = 0.1
Log Loss : 1.1748243948287305
for alpha = 1
Log Loss : 1.1988744136473353
for alpha = 10
Log Loss : 1.3609421484837234
for alpha = 100
Log Loss : 1.4319274865409584
for alpha = 1000
Log Loss : 1.430027611678733

```



```

For values of best alpha = 0.1 The train log loss is: 0.6900123099128517
For values of best alpha = 0.1 The cross validation log loss is: 1.1748243948287305
For values of best alpha = 0.1 The test log loss is: 1.173064487909269

```

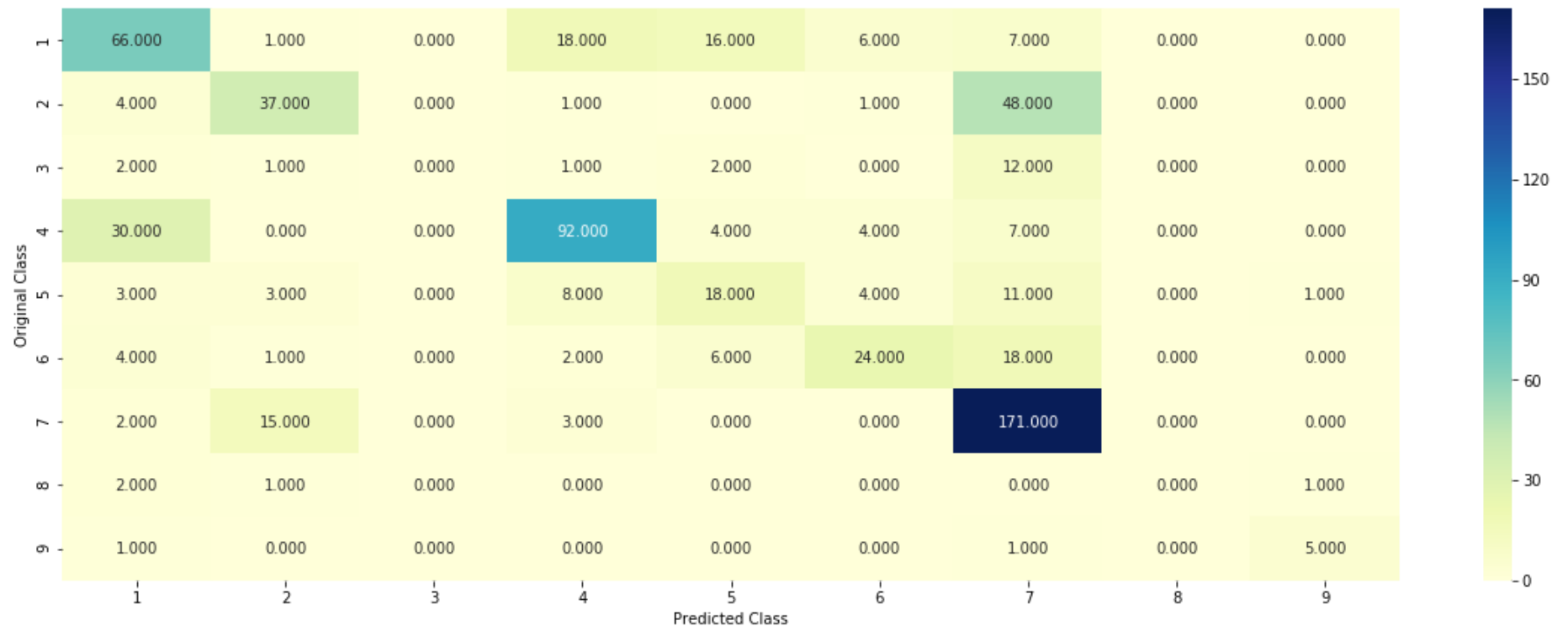


```
In [115]: 1 clf = MultinomialNB(alpha=alpha[best_alpha])
2 clf.fit(train_x_onehotCoding_tfidf_top_1000, train_y)
3 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
4 sig_clf.fit(train_x_onehotCoding_tfidf_top_1000, train_y)
5 sig_clf_probs = sig_clf.predict_proba(test_x_onehotCoding_tfidf_top_1000)
6 # to avoid rounding error while multiplying probabilities we use log-probability estimates
7 print("Log Loss :", log_loss(test_y, sig_clf_probs))
8 print("Number of missclassified point :", np.count_nonzero((sig_clf.predict(test_x_onehotCoding_tfidf_top_1
9 plot_confusion_matrix(test_y, sig_clf.predict(test_x_onehotCoding_tfidf_top_1000))
```

Log Loss : 1.173064487909269

Number of missclassified point : 0.37894736842105264

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



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



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



In [116]: 1 *# Just Like this perform on various models*

In [117]: 1 *# TASK 3*

```
In [118]: 1 # building a CountVectorizer with all the words that occurred minimum 3 times in train data
2 text_vectorizer = CountVectorizer(min_df=3, ngram_range=(1,2))
3 train_text_feature_onehotCoding_bigrams = 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) vect
8 train_text_fea_counts = train_text_feature_onehotCoding_bigrams.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 1-grams and 2-grams in train data :", len(train_text_features))
```

Total number of unique 1-grams and 2-grams in train data : 787757

```
In [120]: 1 cv_text_onehotCoding_bigrams=text_vectorizer.transform(cv_df['TEXT'])
2 test_text_onehotCoding_bigrams=text_vectorizer.transform(test_df['TEXT'])
```

```
In [121]: 1 train_x_onehotCoding_bigrams = hstack([train_gene_var_onehotCoding, train_text_feature_onehotCoding_bigrams]
2
3 test_x_onehotCoding_bigrams = hstack((test_gene_var_onehotCoding, test_text_onehotCoding_bigrams))
4
5 cv_x_onehotCoding_bigrams = hstack((cv_gene_var_onehotCoding, cv_text_onehotCoding_bigrams))
```

```

In [123]: 1 alpha = [10 ** x for x in range(0, 6)]
2 cv_log_error_array = []
3 for i in alpha:
4     print("for alpha =", i)
5     clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42, n_jobs=
6     clf.fit(train_x_onehotCoding_bigrams, train_y)
7     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
8     sig_clf.fit(train_x_onehotCoding_bigrams, train_y)
9     sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding_bigrams)
10    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
11    # to avoid rounding error while multiplying probabilities we use log-probability estimates
12    print("Log Loss :", log_loss(cv_y, sig_clf_probs))
13
14    fig, ax = plt.subplots()
15    ax.plot(alpha, cv_log_error_array, c='g')
16    for i, txt in enumerate(np.round(cv_log_error_array, 3)):
17        ax.annotate((alpha[i], str(txt)), (alpha[i], cv_log_error_array[i]))
18    plt.grid()
19    plt.title("Cross Validation Error for each alpha")
20    plt.xlabel("Alpha i's")
21    plt.ylabel("Error measure")
22    plt.show()
23
24
25    best_alpha = np.argmin(cv_log_error_array)
26    clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=
27    clf.fit(train_x_onehotCoding_bigrams, train_y)
28    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
29    sig_clf.fit(train_x_onehotCoding_bigrams, train_y)
30
31    predict_y = sig_clf.predict_proba(train_x_onehotCoding_bigrams)
32    print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_
33    predict_y = sig_clf.predict_proba(cv_x_onehotCoding_bigrams)
34    print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv,
35    predict_y = sig_clf.predict_proba(test_x_onehotCoding_bigrams)
36    print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y,

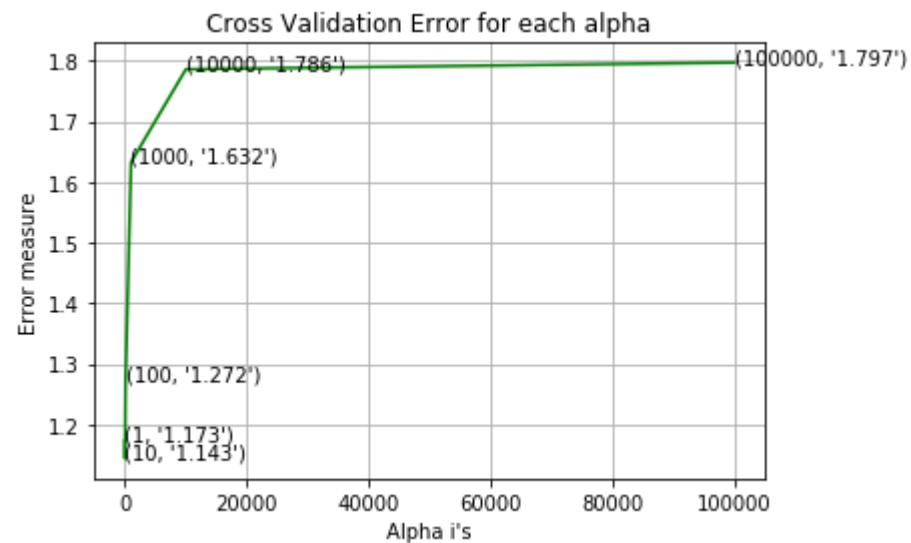
```

```

for alpha = 1
Log Loss : 1.1732603879951011
for alpha = 10
Log Loss : 1.143466339743986

```

```
for alpha = 100
Log Loss : 1.2717856739723545
for alpha = 1000
Log Loss : 1.6324919602788461
for alpha = 10000
Log Loss : 1.7858845090446755
for alpha = 100000
Log Loss : 1.7968406884272987
```



For values of best alpha = 10 The train log loss is: 0.8812540858048751

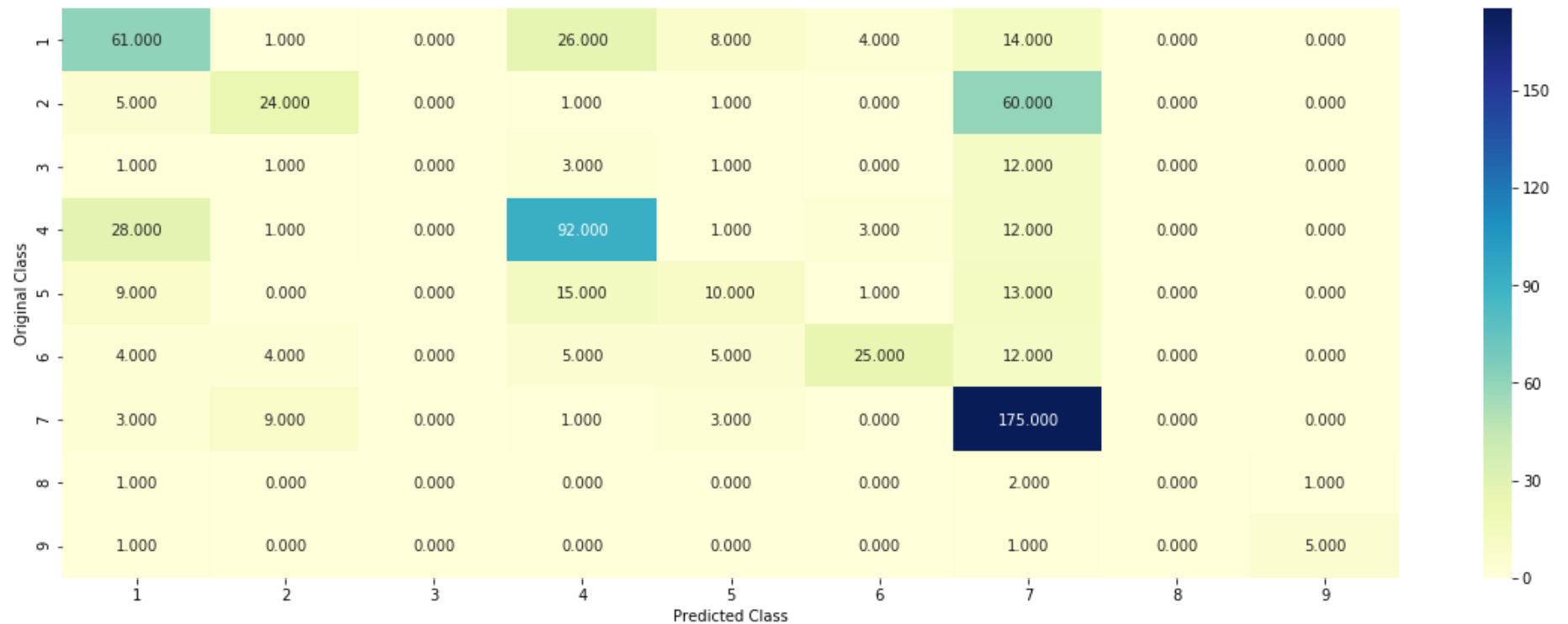
```
For values of best alpha = 10 The cross validation log loss is: 1.143466339743986  
For values of best alpha = 10 The test log loss is: 1.1868654202322484
```

```
In [125]: 1 clf=SGDClassifier(class_weight='balanced', random_state=42, loss='log', alpha=10, n_jobs=3)
2 clf.fit(train_x_onehotCoding_bigrams, y_train)
3 predict_and_plot_confusion_matrix(train_x_onehotCoding_bigrams, train_y, test_x_onehotCoding_bigrams, test_y)
```

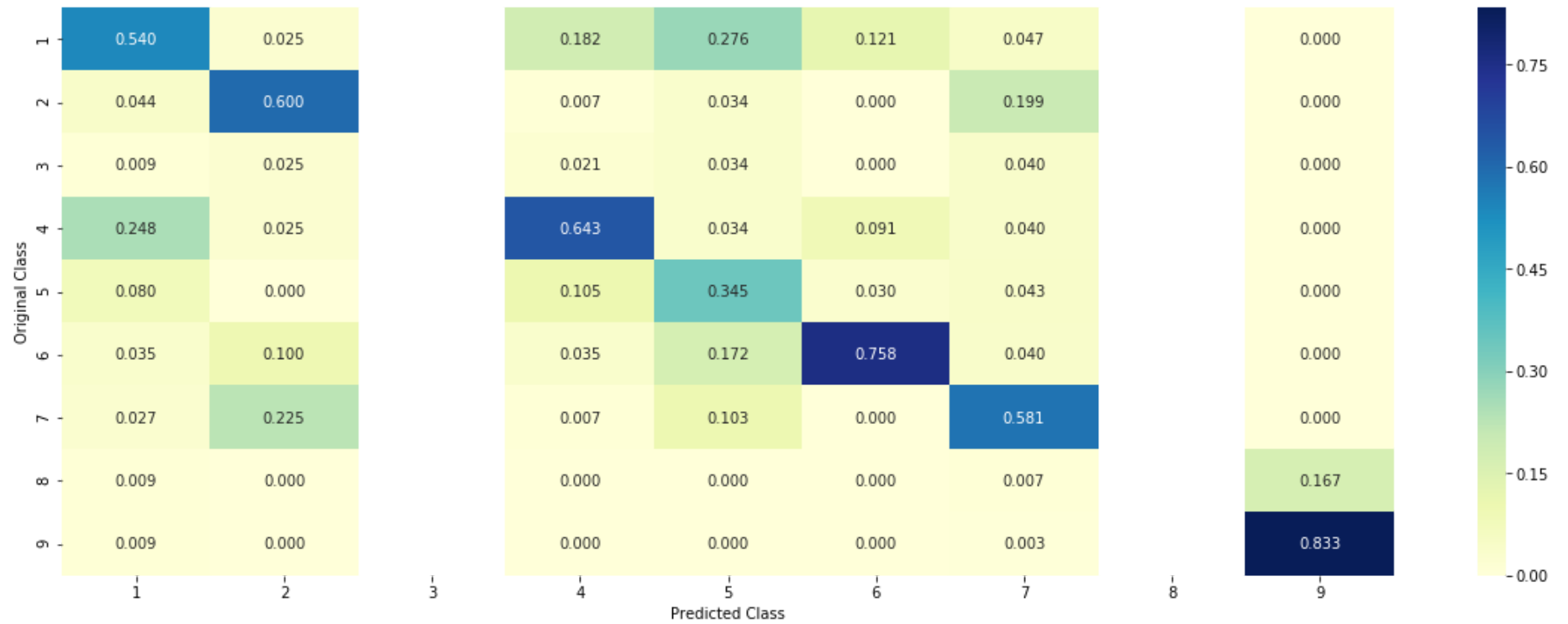
Log loss : 1.1868654202322484

Number of mis-classified points : 0.4105263157894737

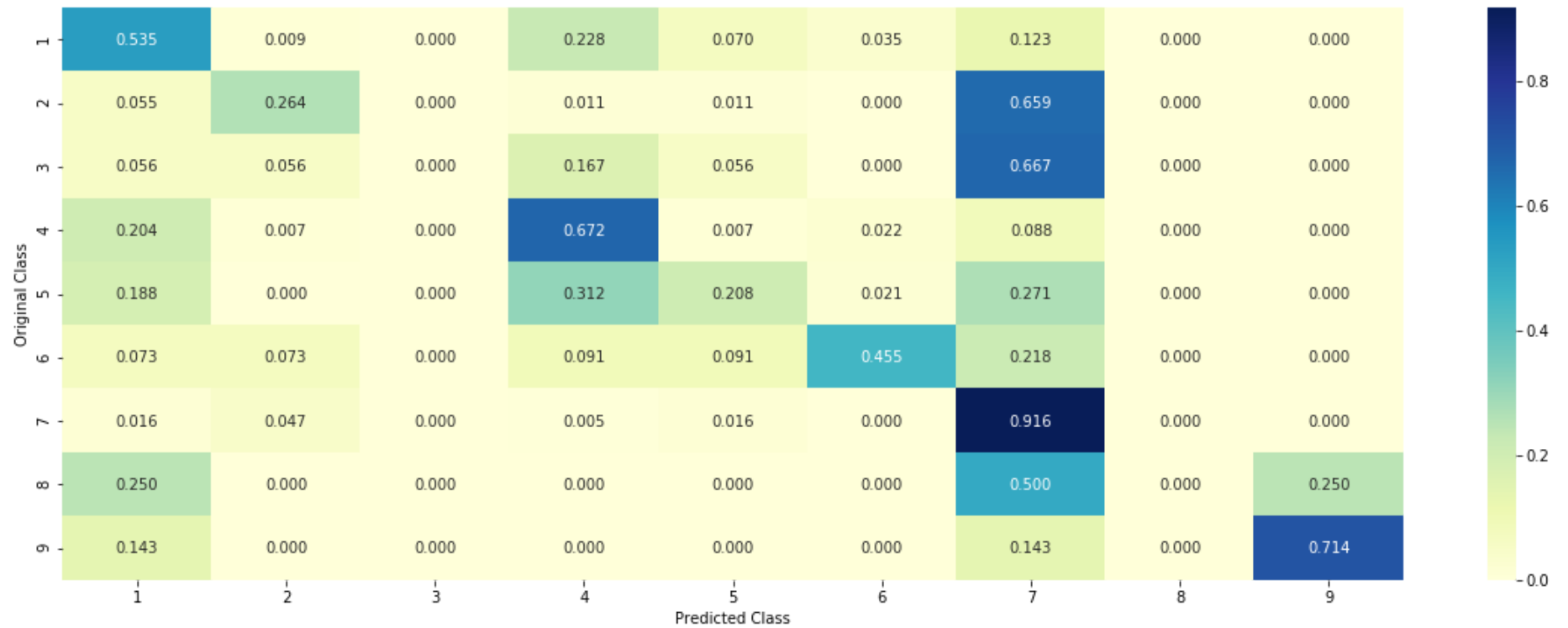
----- Confusion matrix -----



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



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



```
In [119]: 1 #Task 4
```

```
In [127]: 1 import dill
          2 dill.dump_session('cancer.db')
```

```
In [ ]: 1
```