# Personalized cancer diagnosis

## 1. Business Problem

### 1.1. Description

Source: https://www.kaggle.com/c/msk-redefining-cancer-treatment/ (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 (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=qxXRKVompI8 (https://www.youtube.com/watch?v=qxXRKVompI8)

### 1.3. Real-world/Business objectives and constraints.

- No low-latency requirement.
- Interpretability is important.
- Errors can be very costly.
- Probability of a data-point belonging to each class is needed.

## 2. Machine Learning Problem Formulation

### 2.1. Data

#### 2.1.1. Data Overview

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

#### 2.1.2. Example Data Point

***training_variants***

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

***training_text***

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

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

### 2.2.1. Type of Machine Learning Problem

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

### 2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

### 2.2.3. Machine Learing 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 probabilites => 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

# Exploratory Data Analysis ¶

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

from mlxtend.classifier import StackingClassifier

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
```

## Reading Data

### Reading Gene and Variation Data

```python
data = pd.read_csv('training_variants.zip',compression = 'zip')
print('Number of data points : ', data.shape[0])
print('Number of features : ', data.shape[1])
print('Features : ', data.columns.values)
data.head()
```

```
Number of data points :  3321
Number of features :  4
Features :  ['ID' 'Gene' 'Variation' 'Class']
```

Out[11]:

| | 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

### Reading Text Data

```
In [0]:  # note the seprator in this file
         data_text =pd.read_csv("training_text.zip",compression='zip',sep="\|\|",engine="python",names=["ID","TEXT"],skiprows=1)
         print('Number of data points : ', data_text.shape[0])
         print('Number of features : ', data_text.shape[1])
         print('Features : ', data_text.columns.values)
         data_text.head()

         Number of data points :  3321
         Number of features :  2
         Features :  ['ID' 'TEXT']
```

Out[12]:

| | 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... |

## Preprocessing of text

```
In [0]:  import nltk
         nltk.download('stopwords')

         [nltk_data] Downloading package stopwords to /root/nltk_data...
         [nltk_data]   Unzipping corpora/stopwords.zip.
```

Out[15]:  True

```
In [0]:  # loading stop words from nltk library
         stop_words = set(stopwords.words('english'))

         def nlp_preprocessing(total_text, index, column):
             if type(total_text) is not int:
                 string = ""
                 # replace every special char with space
                 total_text = re.sub('[^a-zA-Z\n]', ' ', total_text)
                 # replace multiple spaces with single space
                 total_text = re.sub('\s+',' ', total_text)
                 # converting all the chars into lower-case.
                 total_text = total_text.lower()

                 for word in total_text.split():
                 # if the word is a not a stop word then retain that word from the data
                     if not word in stop_words:
                         string += word + " "

                 data_text[column][index] = string
```

```
In [0]:  #text processing stage.
         start_time = time.clock()
         for index, row in data_text.iterrows():
             if type(row['TEXT']) is str:
                 nlp_preprocessing(row['TEXT'], index, 'TEXT')
             else:
                 print("there is no text description for id:",index)
         print('Time took for preprocessing the text :',time.clock() - start_time, "seconds")

         there is no text description for id: 1109
         there is no text description for id: 1277
         there is no text description for id: 1407
         there is no text description for id: 1639
         there is no text description for id: 2755
         Time took for preprocessing the text : 301.533856 seconds
```

```
In [0]:  #merging both gene_variations and text data based on ID
         result = pd.merge(data, data_text,on='ID', how='left')
         result.head()
```

Out[18]:

| | 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 [0]:  result[result.isnull().any(axis=1)]
```

Out[19]:

| | 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 [0]:  result.loc[result['TEXT'].isnull(),'TEXT'] = result['Gene'] +' '+result['Variation']
```

```
In [0]:  result[result['ID']==1109]
```

Out[21]:

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

```
In [0]:  variation_counts.head()
```

Out[25]:
```
Truncating Mutations    93
Deletion                74
Amplification           71
Fusions                 34
Overexpression           6
Name: Variation, dtype: int64
```

```
In [0]:  result.drop(columns=['Unnamed: 0'],inplace=True)
```

## Feature Engineering Part I (Gene and Variation)

**1)** Upon closer examination of the Variations, it is found that some variations are just combination of two genes with a suffix 'Fusion' in the end.

**2)** Some Variations contain an Asterisk term in the end.

**3)** A simple feature engineering is applied which creates a new column in the dataframe 'IsFusion' and 'IsAsterisk' which displays 1 if Fusion or Asterisk is contained and 0 otherwise.

```python
isFusion = [] #List to keep track of Fusions
cnt=0 #Keeps count of number of Variations with fusions
for feature in result.Variation.values:
    if feature[-6:]=='Fusion': #Append 1 if last 6 strings of the variation word consists of 'Fusion'
        isFusion.append(1)
        cnt+=1
    else:
        isFusion.append(0)
print('Number of Variations which are a Fusion are',cnt)

cnt=0 #Keeps count of number of Variations with Asterisk
isAsterisk = [] #List to keep track of Asterisk Variations
for feature in result.Variation.values:
    if feature[-1]=='*': #Append 1 if last string of the variation word consists of '*'
        isAsterisk.append(1)
        cnt+=1
    else:
        isAsterisk.append(0)
print('Number of Variations containing asterisk at the end are',cnt)
```

```
Number of Variations which are a Fusion are 148
Number of Variations containing asterisk at the end are 56
```

```python
result['IsFusion']=isFusion #Creates a new column in the dataframe with the feature
result['IsAsterisk']=isAsterisk
result.drop(columns=['Unnamed: 0'],inplace=True)
```

```python
result.head()
```

Out[30]:

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

## Train, Test and Cross validation split (64:20:16)

```python
y_true = result['Class'].values
result.Gene      = result.Gene.str.replace('\s+', '_')
result.Variation = result.Variation.str.replace('\s+', '_')

# split the data into test and train by maintaining same distribution of output varaible 'y_true' [stratify=y_true]
X_train, test_df, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_size=0.2)
# split the train data into train and cross validation by maintaining same distribution of output varaible 'y_train' [stratify=y_train]
train_df, cv_df, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)
```

We split the data into train, test and cross validation data sets, preserving the ratio of class distribution in the original data set

```python
print('Number of data points in train data:', train_df.shape[0])
print('Number of data points in test data:', test_df.shape[0])
print('Number of data points in cross validation data:', cv_df.shape[0])
```

```
Number of data points in train data: 2124
Number of data points in test data: 665
Number of data points in cross validation data: 532
```

### Distribution of y_i's in Train, Test and Cross Validation datasets

```
In [0]:  # it returns a dict, keys as class labels and values as the number of data points in that class
         train_class_distribution = train_df['Class'].value_counts().sortlevel()
         test_class_distribution = test_df['Class'].value_counts().sortlevel()
         cv_class_distribution = cv_df['Class'].value_counts().sortlevel()

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

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


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

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

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

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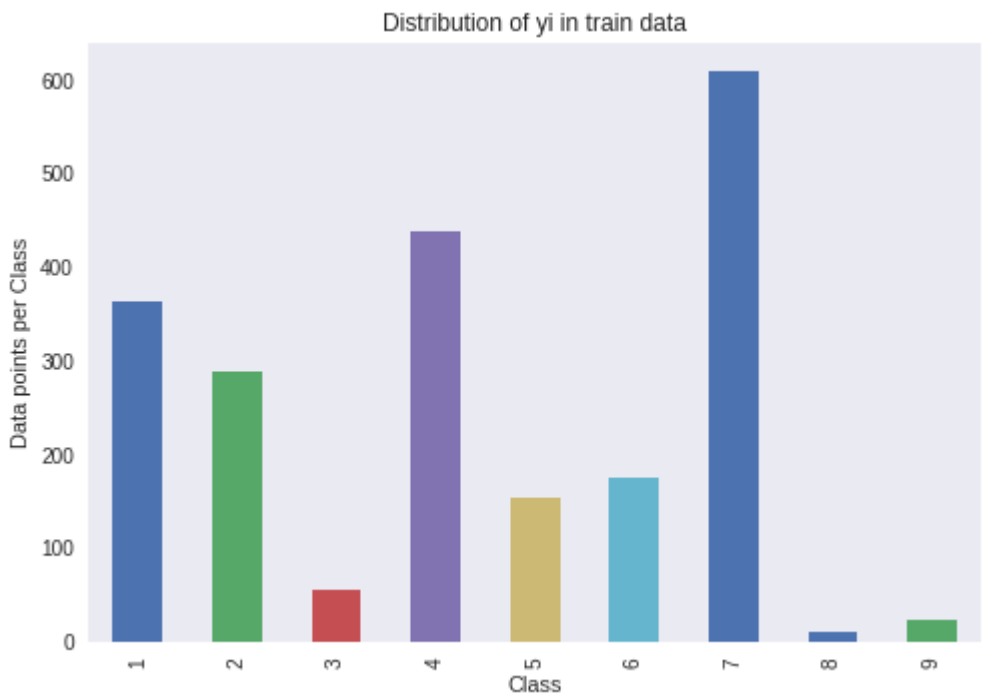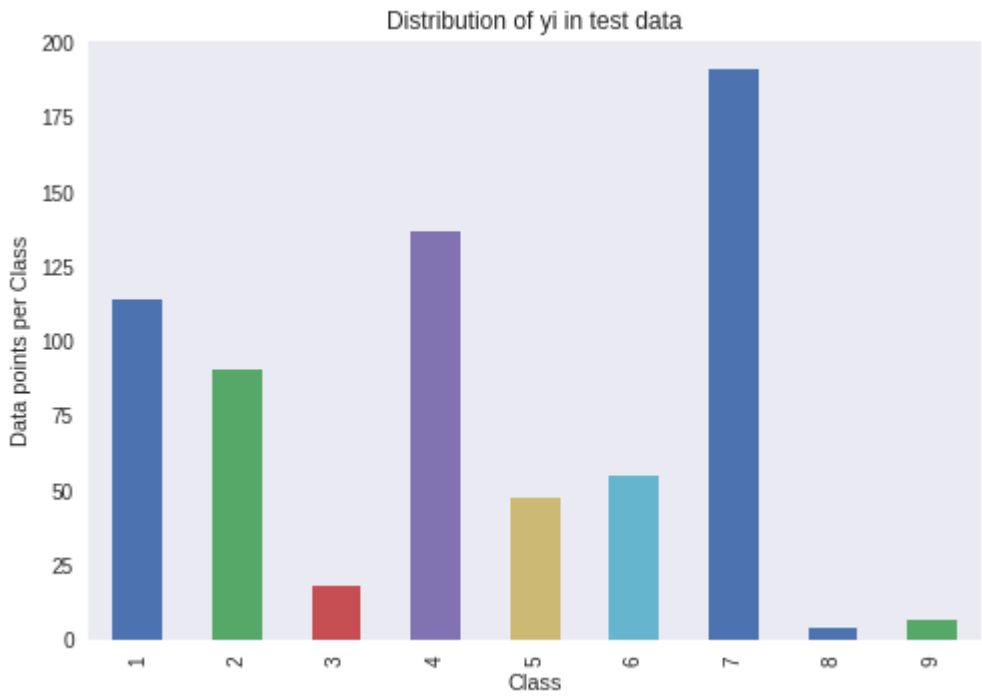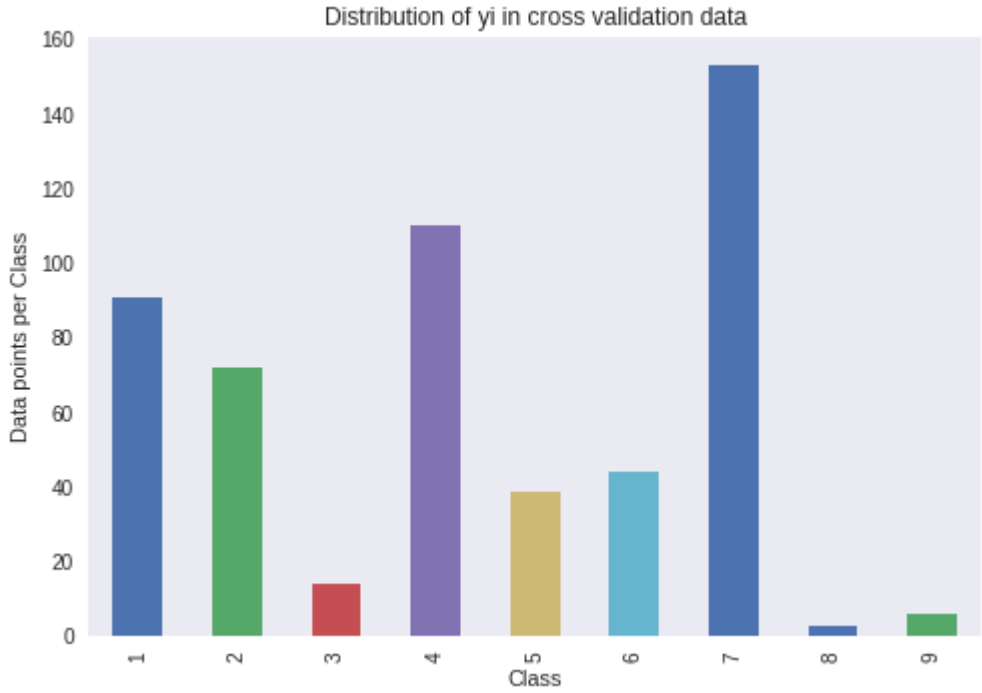


Distribution of yi in train data

```
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 %)
--------------------------------------------------------------------------------
```



Distribution of yi in test data

```
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 %)
--------------------------------------------------------------------------------
```

Distribution of yi in cross validation data

```
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 %)
```

## Prediction using a 'Random' Model

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

```python
In [0]:  # This function plots the confusion matrices given y_i, y_i_hat.
         def plot_confusion_matrix(test_y, predict_y):
             C = confusion_matrix(test_y, predict_y)
             # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j

             A =(((C.T)/(C.sum(axis=1))).T)
             #divid each element of the confusion matrix with the sum of elements in that column

             # C = [[1, 2],
             #      [3, 4]]
             # C.T = [[1, 3],
             #        [2, 4]]
             # C.sum(axis = 1)  axis=0 corresonds to columns and axis=1 corresponds to rows in two diamensional array
             # C.sum(axix =1) = [[3, 7]]
             # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
             #                            [2/3, 4/7]]

             # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
             #                              [3/7, 4/7]]
             # sum of row elements = 1

             B =(C/C.sum(axis=0))
             #divid each element of the confusion matrix with the sum of elements in that row
             # C = [[1, 2],
             #      [3, 4]]
             # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to rows in two diamensional array
             # C.sum(axix =0) = [[4, 6]]
             # (C/C.sum(axis=0)) = [[1/4, 2/6],
             #                      [3/4, 4/6]]

             labels = [1,2,3,4,5,6,7,8,9]
             # representing A in heatmap format
             print("-"*20, "Confusion matrix", "-"*20)
             plt.figure(figsize=(20,7))
             sns.heatmap(C, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
             plt.xlabel('Predicted Class')
             plt.ylabel('Original Class')
             plt.show()

             print("-"*20, "Precision matrix (Columm Sum=1)", "-"*20)
             plt.figure(figsize=(20,7))
             sns.heatmap(B, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
             plt.xlabel('Predicted Class')
             plt.ylabel('Original Class')
             plt.show()

             # representing B in heatmap format
             print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
             plt.figure(figsize=(20,7))
             sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
             plt.xlabel('Predicted Class')
             plt.ylabel('Original Class')
             plt.show()
```
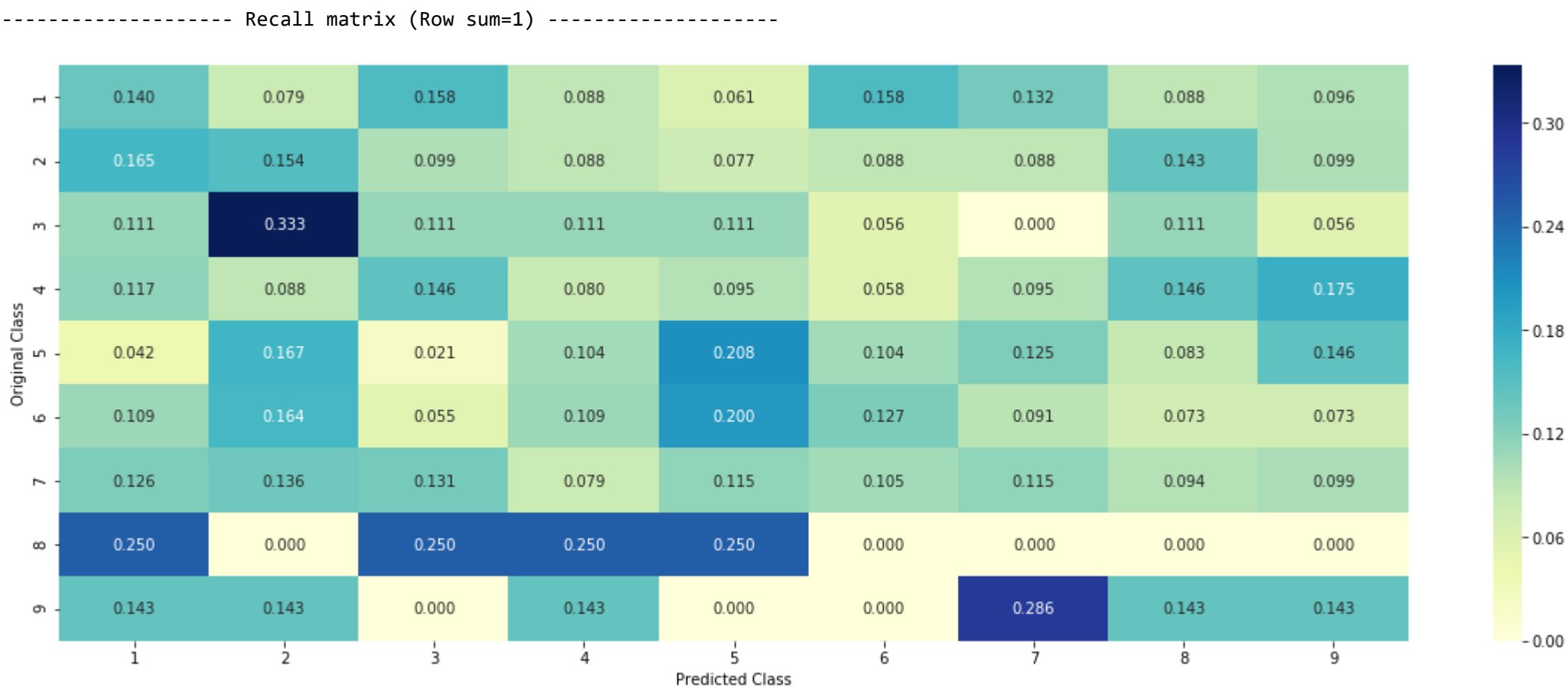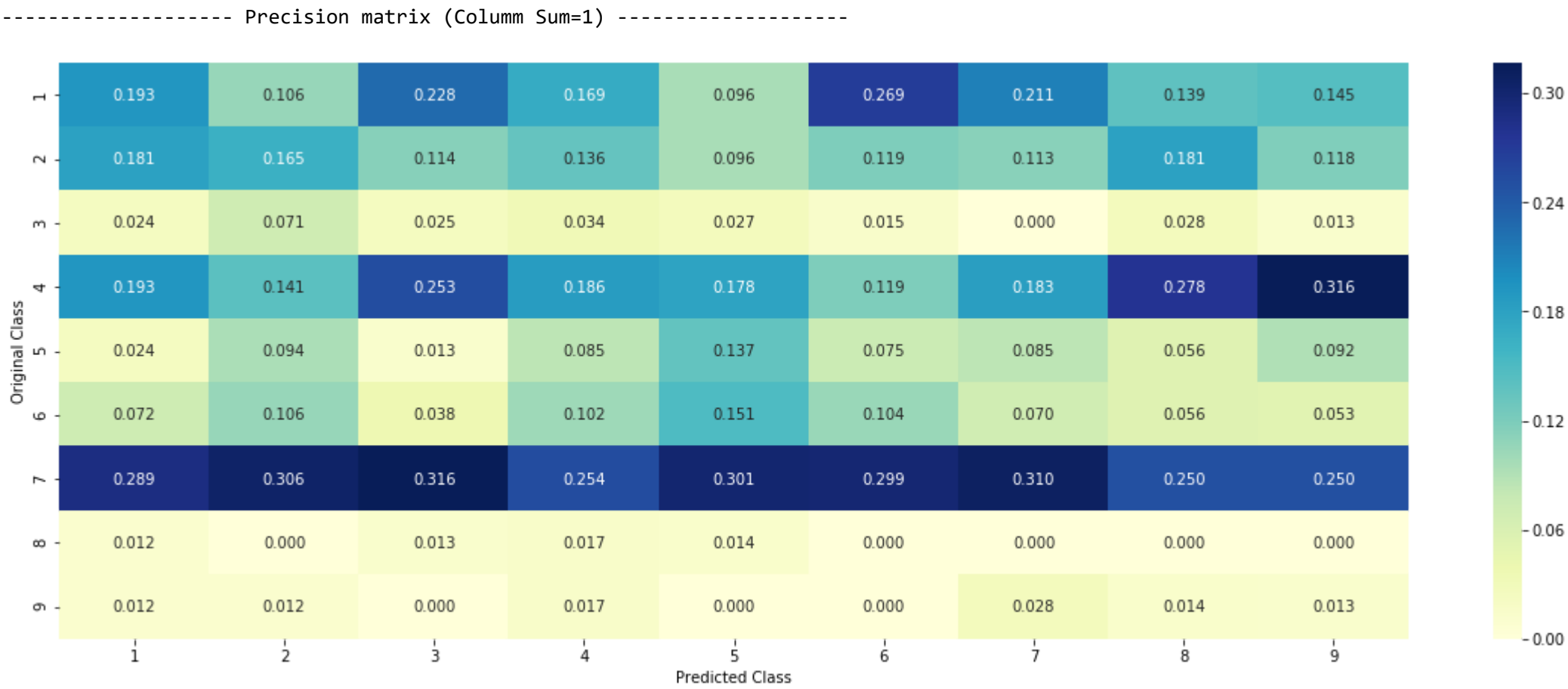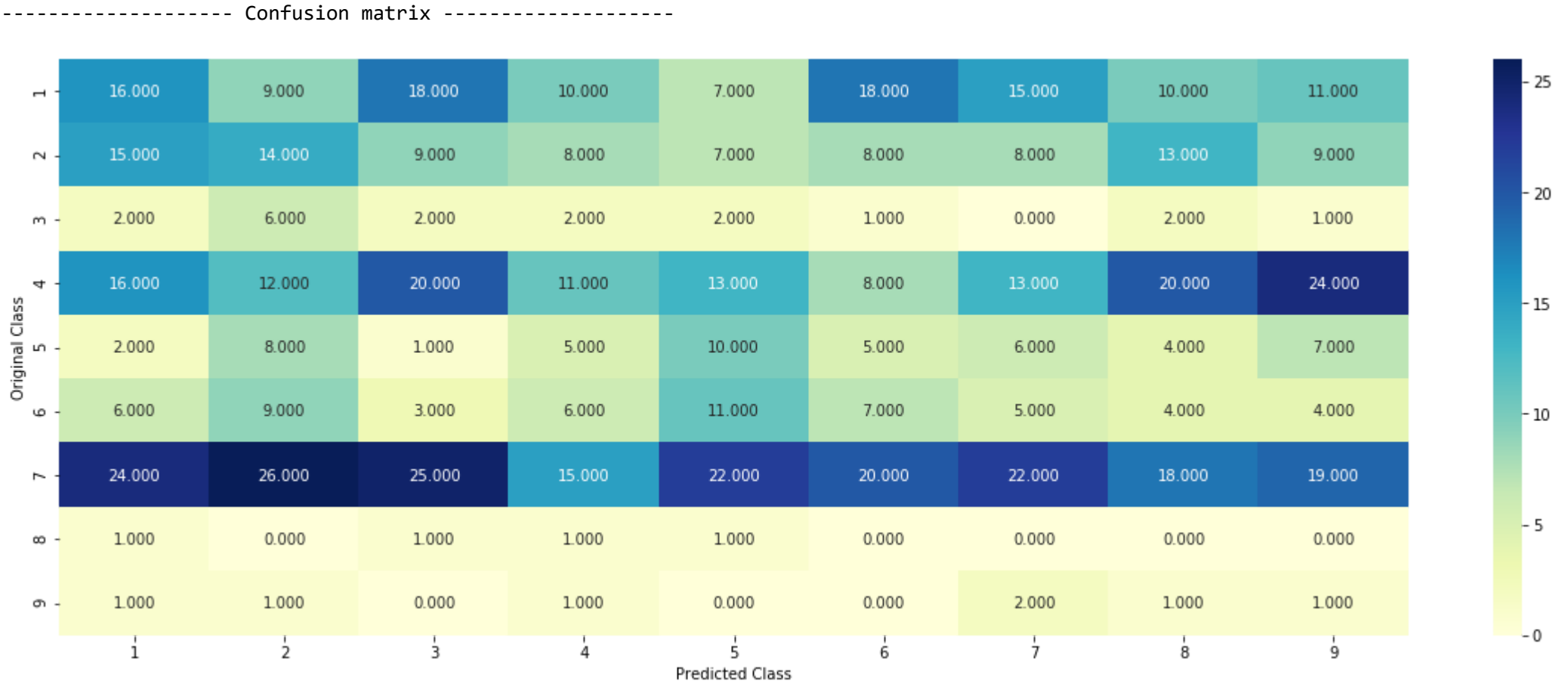
```
In [0]:  # we need to generate 9 numbers and the sum of numbers should be 1
         # one solution is to genarate 9 numbers and divide each of the numbers by their sum
         # ref: https://stackoverflow.com/a/18662466/4084039
         test_data_len = test_df.shape[0]
         cv_data_len = cv_df.shape[0]

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

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

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

```
Log loss on Cross Validation Data using Random Model 2.508503263105886
Log loss on Test Data using Random Model 2.4673429249348153
-------------------- Confusion matrix --------------------
```



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



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



# Univariate Analysis

```python
In [0]:  # code for response coding with Laplace smoothing.
         # alpha : used for Laplace smoothing
         # feature: ['gene', 'variation']
         # df: ['train_df', 'test_df', 'cv_df']
         # algorithm
         # ----------
         # Consider all unique values and the number of occurances of given feature in train data dataframe
         # build a vector (1*9) , the first element = (number of times it occured in class1 + 10*alpha / number of time it occurred in total data+90*alpha)
         # gv_dict is like a look up table, for every gene it store a (1*9) representation of it
         # for a value of feature in df:
         # if it is in train data:
         # we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
         # if it is not there is train:
         # we add [1/9, 1/9, 1/9, 1/9,1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
         # return 'gv_fea'
         # ----------------------

         # get_gv_fea_dict: Get Gene varaition Feature Dict
         def get_gv_fea_dict(alpha, feature, df):
             # value_count: it contains a dict like
             # print(train_df['Gene'].value_counts())
             # output:
             #        {BRCA1      174
             #         TP53       106
             #         EGFR        86
             #         BRCA2       75
             #         PTEN        69
             #         KIT         61
             #         BRAF        60
             #         ERBB2       47
             #         PDGFRA      46
             #          ...}
             # print(train_df['Variation'].value_counts())
             # output:
             # {
             # Truncating_Mutations                     63
             # Deletion                                 43
             # Amplification                            43
             # Fusions                                  22
             # Overexpression                            3
             # E17K                                      3
             # Q61L                                      3
             # S222D                                     2
             # P130S                                     2
             # ...
             # }
             value_count = train_df[feature].value_counts()

             # gv_dict : Gene Variation Dict, which contains the probability array for each gene/variation
             gv_dict = dict()

             # denominator will contain the number of time that particular feature occured in whole data
             for i, denominator in value_count.items():
                 # vec will contain (p(yi==1/Gi) probability of gene/variation belongs to perticular class
                 # vec is 9 diamensional vector
                 vec = []
                 for k in range(1,10):
                     # print(train_df.loc[(train_df['Class']==1) & (train_df['Gene']=='BRCA1')])
                     #           ID   Gene            Variation  Class
                     # 2470  2470  BRCA1               S1715C      1
                     # 2486  2486  BRCA1               S1841R      1
                     # 2614  2614  BRCA1                  M1R      1
                     # 2432  2432  BRCA1               L1657P      1
                     # 2567  2567  BRCA1               T1685A      1
                     # 2583  2583  BRCA1               E1660G      1
                     # 2634  2634  BRCA1               W1718L      1
                     # cls_cnt.shape[0] will return the number of rows

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

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

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

         # Get Gene variation feature
         def get_gv_feature(alpha, feature, df):
             # print(gv_dict)
             #     {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.068181818181818177, 0.13636363636363635, 0.25, 0.19318181818181818, 0.03787878787878788, 0.03787878787878788, 0.03787
             #      'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.27040816326530615, 0.061224489795918366, 0.066326530612244902, 0.051020408163265307, 0.05102040
             #      'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.068181818181818177, 0.068181818181818177, 0.0625, 0.34659090909090912, 0.0625, 0.056818181818181816],
             #      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.078787878787878782, 0.13939393939393939, 0.34545454545454546, 0.060606060606060608, 0.060606060
             #      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.46540880503144655, 0.07547169811320754, 0.062893081761006289, 0.069182389937106917, 0.0628936
             #      'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295, 0.066225165562913912, 0.066225165562913912, 0.27152317880794702, 0.066225165
             #      'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333333333333334, 0.073333333333333334, 0.093333333333333338, 0.080000000000000002, 0.29999999999999999, 0.06666666
             #      ...
             #      }
             gv_dict = get_gv_fea_dict(alpha, feature, df)
             # value_count is similar in get_gv_fea_dict
             value_count = train_df[feature].value_counts()

             # gv_fea: Gene_variation feature, it will contain the feature for each feature value in the data
             gv_fea = []
             # for every feature values in the given data frame we will check if it is there in the train data then we will add the feature to gv_fea
             # if not we will add [1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9] to gv_fea
             for index, row in df.iterrows():
                 if row[feature] in dict(value_count).keys():
                     gv_fea.append(gv_dict[row[feature]])
                 else:
                     gv_fea.append([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
         #            gv_fea.append([-1,-1,-1,-1,-1,-1,-1,-1,-1])
             return gv_fea
```

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

- (numerator + 10*alpha) / (denominator + 90*alpha)

## Univariate Analysis on Gene Feature

**Q1.** Gene, What type of feature it is ?

**Ans.** Gene is a categorical variable

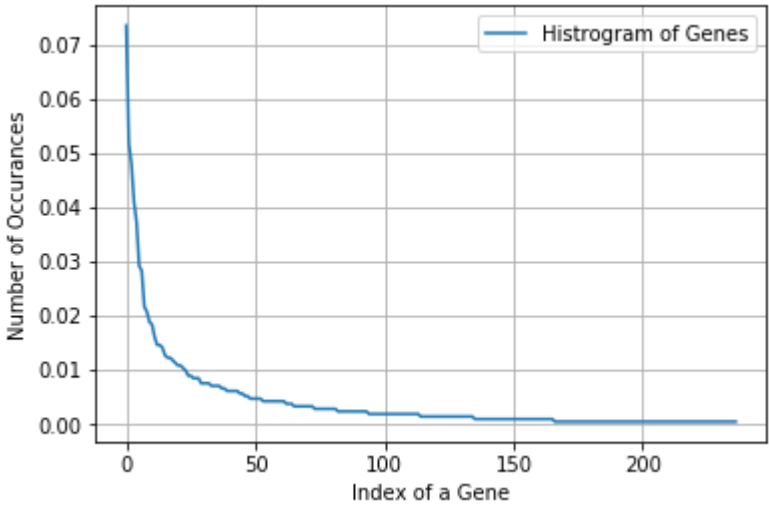**Q2.** How many categories are there and How they are distributed?

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

```
Number of Unique Genes : 235
BRCA1      163
TP53       108
BRCA2       91
EGFR        86
PTEN        77
BRAF        67
KIT         53
ALK         45
PIK3CA      40
PDGFRA      40
Name: Gene, dtype: int64
```
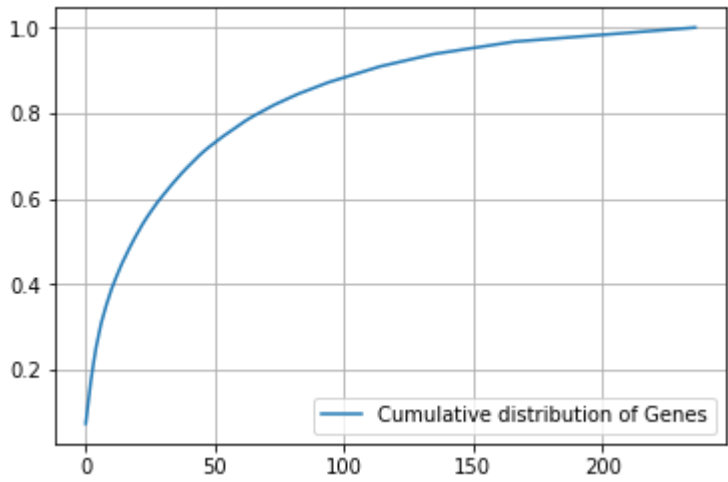
```
In [0]: print("Ans: There are", unique_genes.shape[0] ,"different categories of genes in the train data, and they are distibuted as follows",)
```

```
Ans: There are 238 different categories of genes in the train data, and they are distibuted as follows
```

```
In [0]: s = sum(unique_genes.values);
        h = unique_genes.values/s;
        plt.plot(h, label="Histrogram of Genes")
        plt.xlabel('Index of a Gene')
        plt.ylabel('Number of Occurances')
        plt.legend()
        plt.grid(linestyle='-')
        plt.show()
```



```
In [0]: c = np.cumsum(h)
        plt.plot(c,label='Cumulative distribution of Genes')
        plt.grid(linestyle='-')
        plt.legend()
        plt.show()
```



**Q3.** How to featurize this Gene feature ?

**Ans.** There are two ways we can featurize this variable

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 [0]: #response-coding of the Gene feature
        # alpha is used for laplace smoothing
        alpha = 1
        # train gene feature
        train_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", train_df))
        # test gene feature
        test_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", test_df))
        # cross validation gene feature
        cv_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", cv_df))
```

```
In [0]: print("train_gene_feature_responseCoding is converted feature using respone coding method. The shape of gene feature:", train_gene_feature_responseCoding.shape)
```

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

```
In [0]: # one-hot encoding of Gene feature.
        gene_vectorizer = CountVectorizer()
        train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])
        test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])
        cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

```
In [0]: train_df['Gene'].head()
```

```
Out[15]: 672       CDKN2A
         2072        TET2
         1908      SMARCA4
         641       CDKN1B
         1693        PMS2
         Name: Gene, dtype: object
```

```
In [0]: print("train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature:", train_gene_feature_onehotCoding.shape)
```

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 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 [0]:  alpha = [10 ** x for x in range(-5, 1)] # hyperparam for SGD classifier.

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

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

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

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

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

         predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
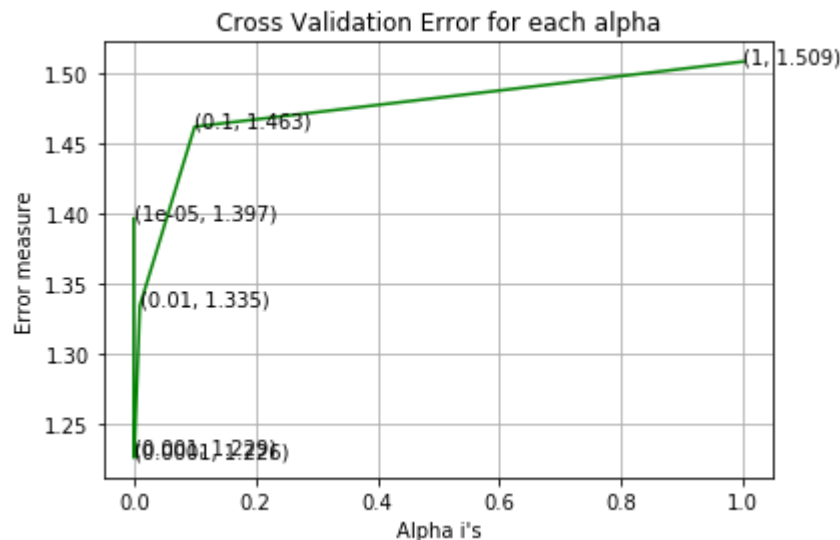
```
For values of alpha =  1e-05 The log loss is: 1.3969593188908682
For values of alpha =  0.0001 The log loss is: 1.226372068366399
For values of alpha =  0.001 The log loss is: 1.229043148119161
For values of alpha =  0.01 The log loss is: 1.3347319401796107
For values of alpha =  0.1 The log loss is: 1.4625194743962024
For values of alpha =  1 The log loss is: 1.5089351225309489
```



```
For values of best alpha =  0.0001 The train log loss is: 1.0768679625883426
For values of best alpha =  0.0001 The cross validation log loss is: 1.226372068366399
For values of best alpha =  0.0001 The test log loss is: 1.1952311902238153
```

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

         test_coverage=test_df[test_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]
         cv_coverage=cv_df[cv_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]

         print('Ans\n1. In test data',test_coverage, 'out of',test_df.shape[0], ":",(test_coverage/test_df.shape[0])*100)
         print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0],":" ,(cv_coverage/cv_df.shape[0])*100)
```

```
Q6. How many data points in Test and CV datasets are covered by the  235  genes in train dataset?
Ans
1. In test data 649 out of 665 : 97.59398496240601
2. In cross validation data 516 out of  532 : 96.99248120300751
```

## Univariate Analysis on Variation Feature

**Q7.** Variation, What type of feature is it ?

**Ans.** Variation is a categorical variable

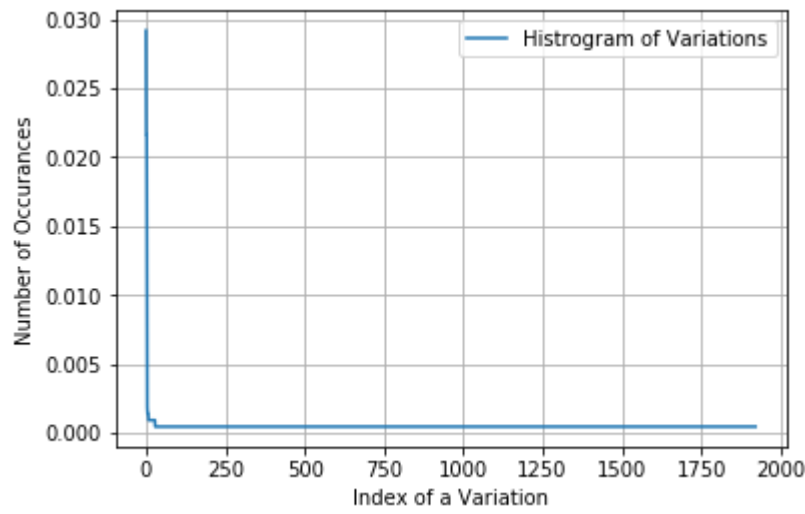**Q8.** How many categories are there?

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

```
Number of Unique Variations : 1930
Truncating_Mutations    63
Amplification           52
Deletion                38
Fusions                 17
Overexpression           4
E17K                     3
T58I                     3
Q61R                     3
Q22K                     2
Q61K                     2
Name: Variation, dtype: int64
```

```
In [0]: print("Ans: There are", unique_variations.shape[0] ,"different categories of variations in the train data, and they are distibuted as follows",)
```
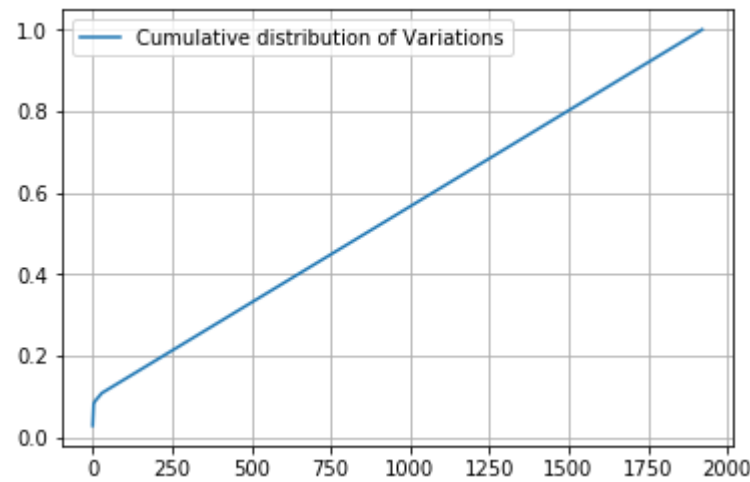
Ans: There are 1930 different categories of variations in the train data, and they are distibuted as follows

```
In [0]: s = sum(unique_variations.values);
        h = unique_variations.values/s;
        plt.plot(h, label="Histrogram of Variations")
        plt.xlabel('Index of a Variation')
        plt.ylabel('Number of Occurances')
        plt.legend()
        plt.grid(linestyle='-')
        plt.show()
```



```
In [0]: c = np.cumsum(h)
        print(c)
        plt.plot(c,label='Cumulative distribution of Variations')
        plt.grid(linestyle='-')
        plt.legend()
        plt.show()
```

[0.02919021 0.05084746 0.07250471 ... 0.99905838 0.99952919 1.        ]



**Q9.** How to featurize this Variation feature ?

**Ans.** There are two ways we can featurize this variable

1. One hot Encoding
2. Response coding

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

```
In [0]: # alpha is used for laplace smoothing
        alpha = 1
        # train gene feature
        train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", train_df))
        # test gene feature
        test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", test_df))
        # cross validation gene feature
        cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", cv_df))
```

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

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

```
In [0]: # one-hot encoding of variation feature.
        variation_vectorizer = CountVectorizer()
        train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(train_df['Variation'])
        test_variation_feature_onehotCoding = variation_vectorizer.transform(test_df['Variation'])
        cv_variation_feature_onehotCoding = variation_vectorizer.transform(cv_df['Variation'])
```

```
In [0]: print("train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape of Variation feature:", train_variation_feature_onehotCoding.shape)
```

train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape of Variation feature: (2124, 1961)

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

Let's build a model just like the earlier!

```
In [0]:    alpha = [10 ** x for x in range(-5, 1)]

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

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

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


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

               sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
               sig_clf.fit(train_variation_feature_onehotCoding, y_train)
               predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)

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

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


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

           predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
           print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
           predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
           print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
           predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding)
           print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
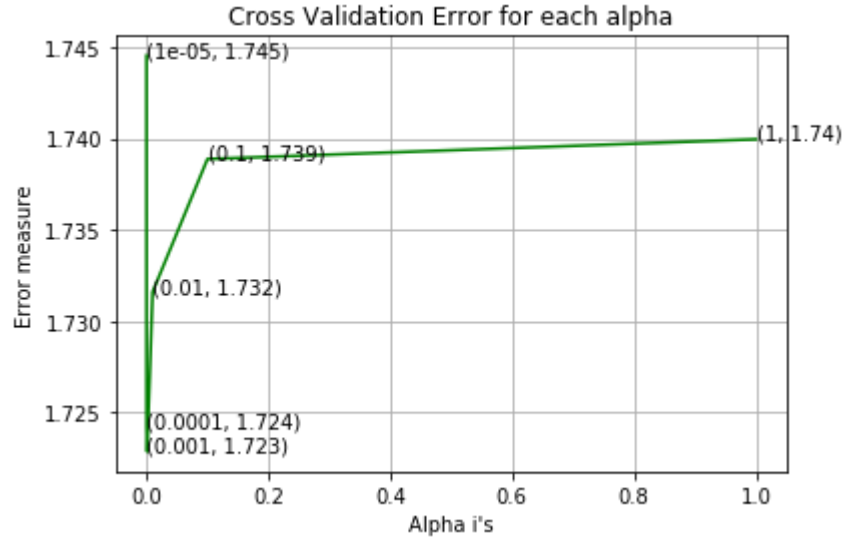
```
For values of alpha =  1e-05 The log loss is: 1.7445442414024976
For values of alpha =  0.0001 The log loss is: 1.724169807225462
For values of alpha =  0.001 The log loss is: 1.722833344920799
For values of alpha =  0.01 The log loss is: 1.7315868954086109
For values of alpha =  0.1 The log loss is: 1.7388743241154072
For values of alpha =  1 The log loss is: 1.7399580414394316
```



```
For values of best alpha =  0.001 The train log loss is: 1.073330998150363
For values of best alpha =  0.001 The cross validation log loss is: 1.722833344920799
For values of best alpha =  0.001 The test log loss is: 1.6978620691350326
```
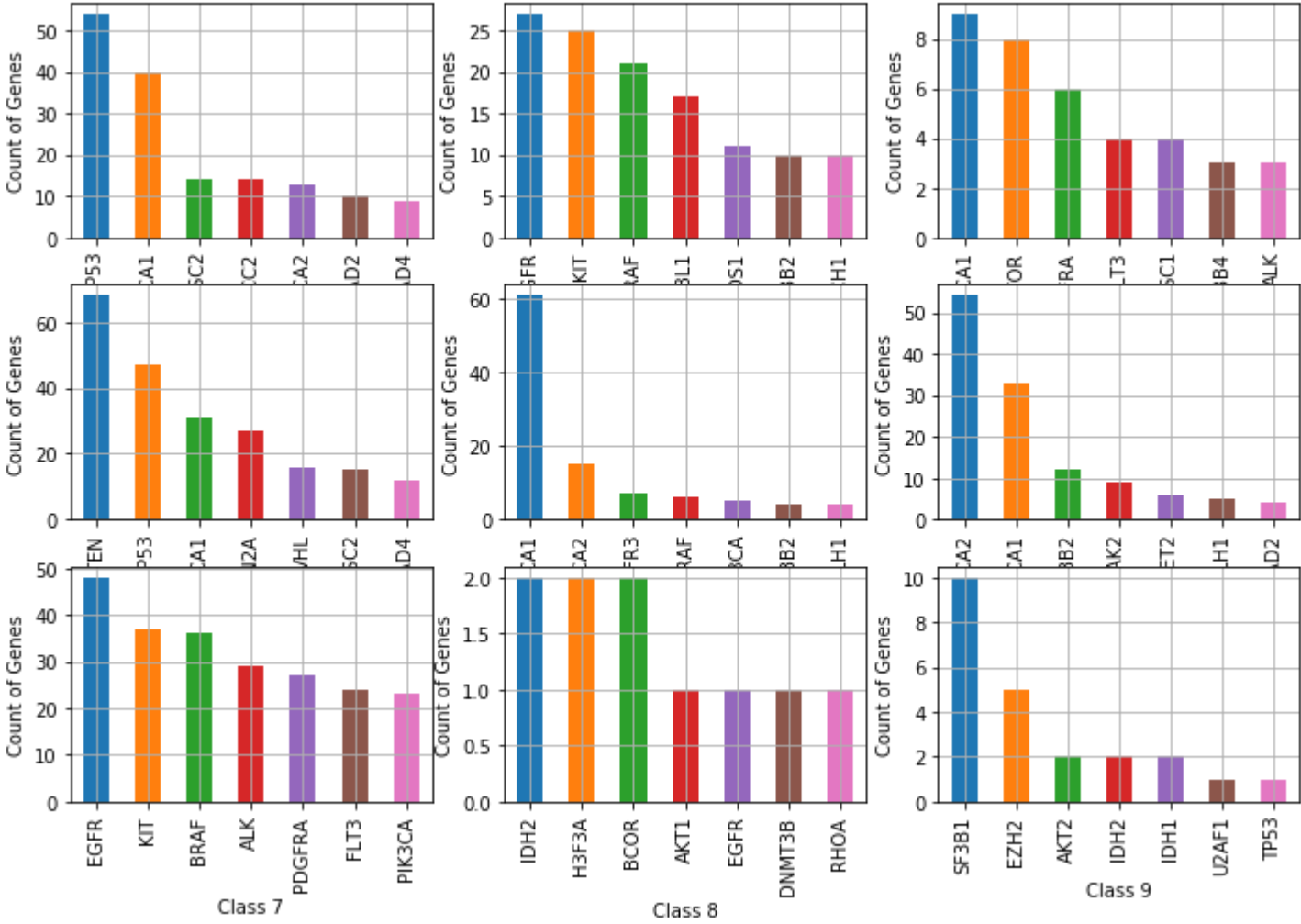
**Q11.** Is the Variation feature stable across all the data sets (Test, Train, Cross validation)?

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

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

```
Q12. How many data points are covered by total  1930  genes in test and cross validation data sets?
Ans
1. In test data 62 out of 665 : 9.323308270676693
2. In cross validation data 65 out of  532 : 12.218045112781954
```

```
In [0]:  #Plots most occuring Genes across all nine classes
         c=331
         plt.figure(figsize=(12,8))
         for i in range(1,10):
             plt.subplot(c)
             train_df[train_df.Class==i].Gene.value_counts()[:7].plot(kind='bar')
             plt.xlabel('Class '+str(i))
             plt.ylabel('Count of Genes')
             plt.grid(linestyle='-')
             c+=1
```



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

```
In [0]:  # cls_text is a data frame
         # for every row in data fram consider the 'TEXT'
         # split the words by space
         # make a dict with those words
         # increment its count whenever we see that word

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

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

```
In [0]:  # building a TFIDFVectorizer with all the words that occured minimum 3 times in train data
         text_vectorizer = TfidfVectorizer(min_df=3,ngram_range=(1,1),max_features=2000)
         train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])
         # getting all the feature names (words)
         train_text_features= text_vectorizer.get_feature_names()

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

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


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

```
Total number of unique words in train data : 2000
```

```
In [0]:  dict_list = []
         # dict_list =[] contains 9 dictionaries each corresponds to a class
         for i in range(1,10):
             cls_text = train_df[train_df['Class']==i]
             # build a word dict based on the words in that class
             dict_list.append(extract_dictionary_paddle(cls_text))
             # append it to dict_list

         # dict_list[i] is build on i'th  class text data
         # total_dict is buid on whole training text data
         total_dict = extract_dictionary_paddle(train_df)


         confuse_array = []
         for i in train_text_features:
             ratios = []
             max_val = -1
             for j in range(0,9):
                 ratios.append((dict_list[j][i]+10 )/(total_dict[i]+90))
             confuse_array.append(ratios)
         confuse_array = np.array(confuse_array)
```

```
In [0]:  #response coding of text features
         train_text_feature_responseCoding  = get_text_responsecoding(train_df)
         test_text_feature_responseCoding   = get_text_responsecoding(test_df)
         cv_text_feature_responseCoding     = get_text_responsecoding(cv_df)
```

```
In [0]:  # https://stackoverflow.com/a/16202486
         # we convert each row values such that they sum to 1
         train_text_feature_responseCoding = (train_text_feature_responseCoding.T/train_text_feature_responseCoding.sum(axis=1)).T
         test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.sum(axis=1)).T
         cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=1)).T
```

```
In [0]:  # Normalize every feature
         train_text_feature_onehotCoding = normalize(train_text_feature_onehotCoding, axis=0)

         # We use the same vectorizer that was trained on train data
         test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
         # Normalize every feature
         test_text_feature_onehotCoding = normalize(test_text_feature_onehotCoding, axis=0)

         # We use the same vectorizer that was trained on train data
         cv_text_feature_onehotCoding = text_vectorizer.transform(cv_df['TEXT'])
         # Normalize every feature
         cv_text_feature_onehotCoding = normalize(cv_text_feature_onehotCoding, axis=0)
```

```
In [0]:  #https://stackoverflow.com/a/2258273/4084039
         sorted_text_fea_dict = dict(sorted(text_fea_dict.items(), key=lambda x: x[1] , reverse=True))
         sorted_text_occur = np.array(list(sorted_text_fea_dict.values()))
```

```
In [0]:  # Train a Logistic regression+Calibration model using text features whicha re on-hot encoded
         alpha = [10 ** x for x in range(-5, 1)]

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

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

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

             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_text_feature_onehotCoding, y_train)
             predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
             cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
             print('For values of alpha = ', i, "The log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

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


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

         predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(test_text_feature_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
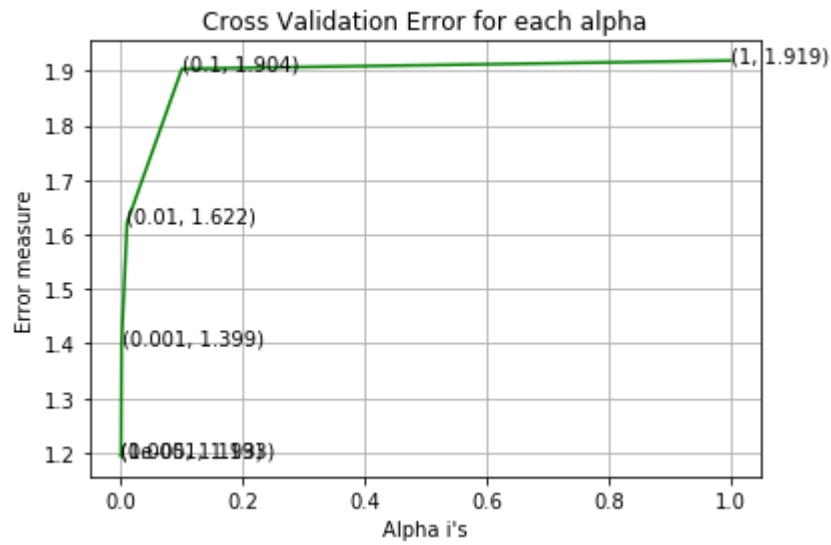
```
For values of alpha =  1e-05 The log loss is: 1.1928781998387548
For values of alpha =  0.0001 The log loss is: 1.1926089184737765
For values of alpha =  0.001 The log loss is: 1.3989495831833352
For values of alpha =  0.01 The log loss is: 1.6220182677326027
For values of alpha =  0.1 The log loss is: 1.903798534554901
For values of alpha =  1 The log loss is: 1.9188319665768077
```



```
For values of best alpha =  0.0001 The train log loss is: 0.7212436582114516
For values of best alpha =  0.0001 The cross validation log loss is: 1.1926089184737765
For values of best alpha =  0.0001 The test log loss is: 1.1158259505127084
```

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

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

```
In [0]:  def get_intersec_text(df):
             df_text_vec = CountVectorizer(min_df=3)
             df_text_fea = df_text_vec.fit_transform(df['TEXT'])
             df_text_features = df_text_vec.get_feature_names()

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

```
In [0]:  len1,len2 = get_intersec_text(test_df)
         print(np.round((len2/len1)*100, 3), "% of words of test data appeared in train data")
         len1,len2 = get_intersec_text(cv_df)
         print(np.round((len2/len1)*100, 3), "% of words of Cross Validation appeared in train data")
```

```
9.524 % of words of test data appeared in train data
10.474 % of words of Cross Validation appeared in train data
```

## Feature Engineering Part II (Text)

```
In [0]:  word_imp_train = [] #Creates a list to store ratio of count of important words in text to the total number of words in that text
         text_len_train = [] #Creates a list to store length of words in the text
         for text in train_df.TEXT.values:
             cnt=0 #Counts the number of times a word which is important according to IDF values appears in the text
             for word in text.split():
                 if word in tfidf_features:
                     cnt+=1
             word_imp_train.append(cnt/len(text.split()))
             text_len_train.append(len(text.split()))

         norm_text_len_train = [(value-min(text_len_train))/(max(text_len_train)-min(text_len_train)) for value in text_len_train] #Perform data normalization on length of words
```

```
In [0]:  word_imp_test = []
         text_len_test = []
         for text in test_df.TEXT.values:
             cnt=0
             for word in text.split():
                 if word in tfidf_features:
                     cnt+=1
             word_imp_test.append(cnt/len(text.split()))
             text_len_test.append(len(text.split()))

         norm_text_len_test = [(value-min(text_len_test))/(max(text_len_test)-min(text_len_test)) for value in text_len_test]
```

```
In [0]:  word_imp_cv = []
         text_len_cv = []
         for text in cv_df.TEXT.values:
             cnt=0
             for word in text.split():
                 if word in tfidf_features:
                     cnt+=1
             word_imp_cv.append(cnt/len(text.split()))
             text_len_cv.append(len(text.split()))

         norm_text_len_cv = [(value-min(text_len_cv))/(max(text_len_cv)-min(text_len_cv)) for value in text_len_cv]
```

```
In [0]:  train_df['norm_text_len_train'] = norm_text_len_train #Adds the normalized words length to the dataframe as a new column
         test_df['norm_text_len_test'] = norm_text_len_test
         cv_df['norm_text_len_cv'] = norm_text_len_cv

         train_df['word_imp_train'] = word_imp_train #Adds the ratio of important words to the dataframe as a new column
         cv_df['word_imp_cv'] = word_imp_cv
         test_df['word_imp_test'] = word_imp_test
```

```
In [0]:  train_df.head()
```

Out[62]:

| | ID | Gene | Variation | Class | TEXT | IsFusion | IsAsterisk | norm_text_len_train | word_imp_train |
|---|---|---|---|---|---|---|---|---|---|
| **672** | 672 | CDKN2A | R80L | 4 | background point mutations tumor suppressor ge... | 0 | 0 | 0.065126 | 0.752126 |
| **2072** | 2072 | TET2 | H1904R | 1 | tet proteins oxidize methylcytosine mc dna pla... | 0 | 0 | 0.097207 | 0.678134 |
| **1908** | 1908 | SMARCA4 | Truncating_Mutations | 1 | small cell carcinoma ovary hypercalcemic type ... | 0 | 0 | 0.226001 | 0.691002 |
| **641** | 641 | CDKN1B | Truncating_Mutations | 1 | cdkn b gene encodes cyclin dependent kinase in... | 0 | 0 | 0.427988 | 0.672485 |
| **1693** | 1693 | PMS2 | G207E | 1 | hereditary nonpolyposis colorectal cancer hnpc... | 0 | 0 | 0.053078 | 0.719474 |

```
In [0]:  features_train = train_df.drop(columns=['ID','Gene','Variation','Class','TEXT']) #Gets seperate dataframe with only feature engineering values
         features_test = test_df.drop(columns=['ID','Gene','Variation','Class','TEXT'])
         features_cv = cv_df.drop(columns=['ID','Gene','Variation','Class','TEXT'])
```

```
In [0]:  features_train_mat = features_train.as_matrix() #Convert it into a matrix
         features_test_mat = features_test.as_matrix()
         features_cv_mat = features_cv.as_matrix()
```

## Stacking the three types of features

```
In [0]:  # merging gene, variance and text features

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

         train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding,train_variation_feature_onehotCoding,features_train_mat))
         test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding,test_variation_feature_onehotCoding,features_test_mat))
         cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding,cv_variation_feature_onehotCoding,features_cv_mat))

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

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

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

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

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

```
In [0]: print("One hot encoding features :")
        print("(number of data points * number of features) in train data = ", train_x_onehotCoding.shape)
        print("(number of data points * number of features) in test data = ", test_x_onehotCoding.shape)
        print("(number of data points * number of features) in cross validation data =", cv_x_onehotCoding.shape)
```

```
One hot encoding features :
(number of data points * number of features) in train data =  (2124, 4200)
(number of data points * number of features) in test data =  (665, 4200)
(number of data points * number of features) in cross validation data = (532, 4200)
```

```
In [0]: print(" Response encoding features :")
        print("(number of data points * number of features) in train data = ", train_x_responseCoding.shape)
        print("(number of data points * number of features) in test data = ", test_x_responseCoding.shape)
        print("(number of data points * number of features) in cross validation data =", cv_x_responseCoding.shape)
```

```
 Response encoding features :
(number of data points * number of features) in train data =  (2124, 31)
(number of data points * number of features) in test data =  (665, 31)
(number of data points * number of features) in cross validation data = (532, 31)
```

## Machine learning model

```
In [0]: #Data preparation for ML models.

        #Misc. functions for ML models

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

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

```
In [0]: def report_log_loss(train_x, train_y, test_x, test_y,  clf):
            clf.fit(train_x, train_y)
            sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
            sig_clf.fit(train_x, train_y)
            sig_clf_probs = sig_clf.predict_proba(test_x)
            return log_loss(test_y, sig_clf_probs, eps=1e-15)
```

```
In [0]: # this function will be used just for naive bayes
        # for the given indices, we will print the name of the features
        # and we will check whether the feature present in the test point text or not
        def get_impfeature_names(indices, text, gene, var, no_features):
            gene_count_vec = CountVectorizer()
            var_count_vec = CountVectorizer()
            text_count_vec = CountVectorizer(min_df=3)

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

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

            word_present = 0
            for i,v in enumerate(indices):
                if (v < fea1_len):
                    word = gene_vec.get_feature_names()[v]
                    yes_no = True if word == gene else False
                    if yes_no:
                        word_present += 1
                        print(i, "Gene feature [{}] present in test data point [{}]".format(word,yes_no))
                elif (v < fea1_len+fea2_len):
                    word = var_vec.get_feature_names()[v-(fea1_len)]
                    yes_no = True if word == var else False
                    if yes_no:
                        word_present += 1
                        print(i, "variation feature [{}] present in test data point [{}]".format(word,yes_no))
                else:
                    word = text_vec.get_feature_names()[v-(fea1_len+fea2_len)]
                    yes_no = True if word in text.split() else False
                    if yes_no:
                        word_present += 1
                        print(i, "Text feature [{}] present in test data point [{}]".format(word,yes_no))

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

### Naive Bayes (Base Line Model)

### Hyper parameter tuning

```python
In [0]:  # find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html
         # -------------------------
         # default paramters
         # sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)

         # some of methods of MultinomialNB()
         # fit(X, y[, sample_weight])    Fit Naive Bayes classifier according to X, y
         # predict(X)    Perform classification on an array of test vectors X.
         # predict_log_proba(X)  Return log-probability estimates for the test vector X.


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


         alpha = [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100,1000]
         cv_log_error_array = []
         for i in alpha:
             print("for alpha =", i)
             clf = MultinomialNB(alpha=i)
             clf.fit(train_x_onehotCoding, train_y)
             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_x_onehotCoding, train_y)
             sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
             cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
             # to avoid rounding error while multiplying probabilites we use log-probability estimates
             print("Log Loss :",log_loss(cv_y, sig_clf_probs))

         fig, ax = plt.subplots()
         ax.plot(np.log10(alpha), cv_log_error_array,c='g')
         for i, txt in enumerate(np.round(cv_log_error_array,3)):
             ax.annotate((alpha[i],str(txt)), (np.log10(alpha[i]),cv_log_error_array[i]))
         plt.grid(linestyle='-')
         plt.xticks(np.log10(alpha))
         plt.title("Cross Validation Error for each alpha")
         plt.xlabel("Alpha i's")
         plt.ylabel("Error measure")
         plt.show()


         best_alpha = np.argmin(cv_log_error_array)
         clf = MultinomialNB(alpha=alpha[best_alpha])
         clf.fit(train_x_onehotCoding, train_y)
         sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
         sig_clf.fit(train_x_onehotCoding, train_y)


         predict_y = sig_clf.predict_proba(train_x_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(test_x_onehotCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
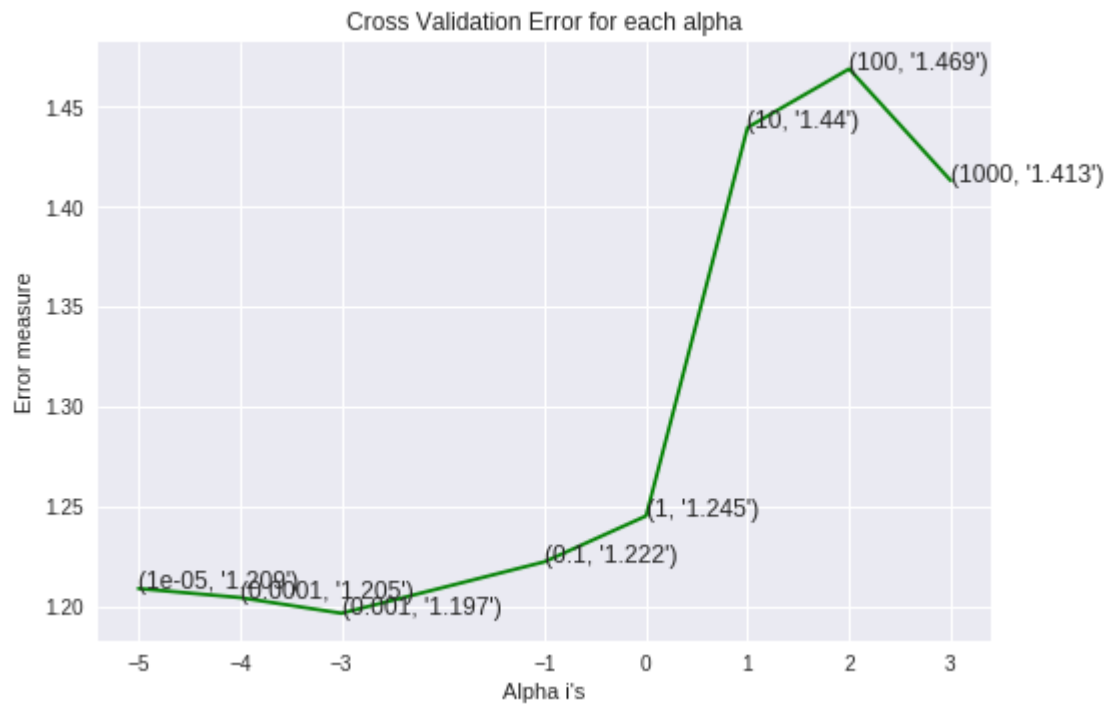
```
for alpha = 1e-05
Log Loss : 1.2090383130894389
for alpha = 0.0001
Log Loss : 1.204523303460323
for alpha = 0.001
Log Loss : 1.1966739622927822
for alpha = 0.1
Log Loss : 1.2224104337118806
for alpha = 1
Log Loss : 1.2453292576605728
for alpha = 10
Log Loss : 1.4395333154414918
for alpha = 100
Log Loss : 1.4688070412050847
for alpha = 1000
Log Loss : 1.412987365119048
```



```
For values of best alpha =  0.001 The train log loss is: 0.5757360550574471
For values of best alpha =  0.001 The cross validation log loss is: 1.1966739622927822
For values of best alpha =  0.001 The test log loss is: 1.1791910545786317
```

### Testing the model with best hyper paramters

```
In [0]: # find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html
        # -------------------------
        # default parmers
        # sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)

        # some of methods of MultinomialNB()
        # fit(X, y[, sample_weight])     Fit Naive Bayes classifier according to X, y
        # predict(X)     Perform classification on an array of test vectors X.
        # predict_log_proba(X)  Return log-probability estimates for the test vector X.


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

        clf = MultinomialNB(alpha=alpha[best_alpha])
        clf.fit(train_x_onehotCoding, train_y)
        sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
        sig_clf.fit(train_x_onehotCoding, train_y)
        sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
        # to avoid rounding error while multiplying probabilites we use log-probability estimates
        print("Log Loss :",log_loss(cv_y, sig_clf_probs))
        print("Number of missclassified point :", np.count_nonzero((sig_clf.predict(cv_x_onehotCoding)- cv_y))/cv_y.shape[0])
        plot_confusion_matrix(cv_y, sig_clf.predict(cv_x_onehotCoding.toarray()))
```
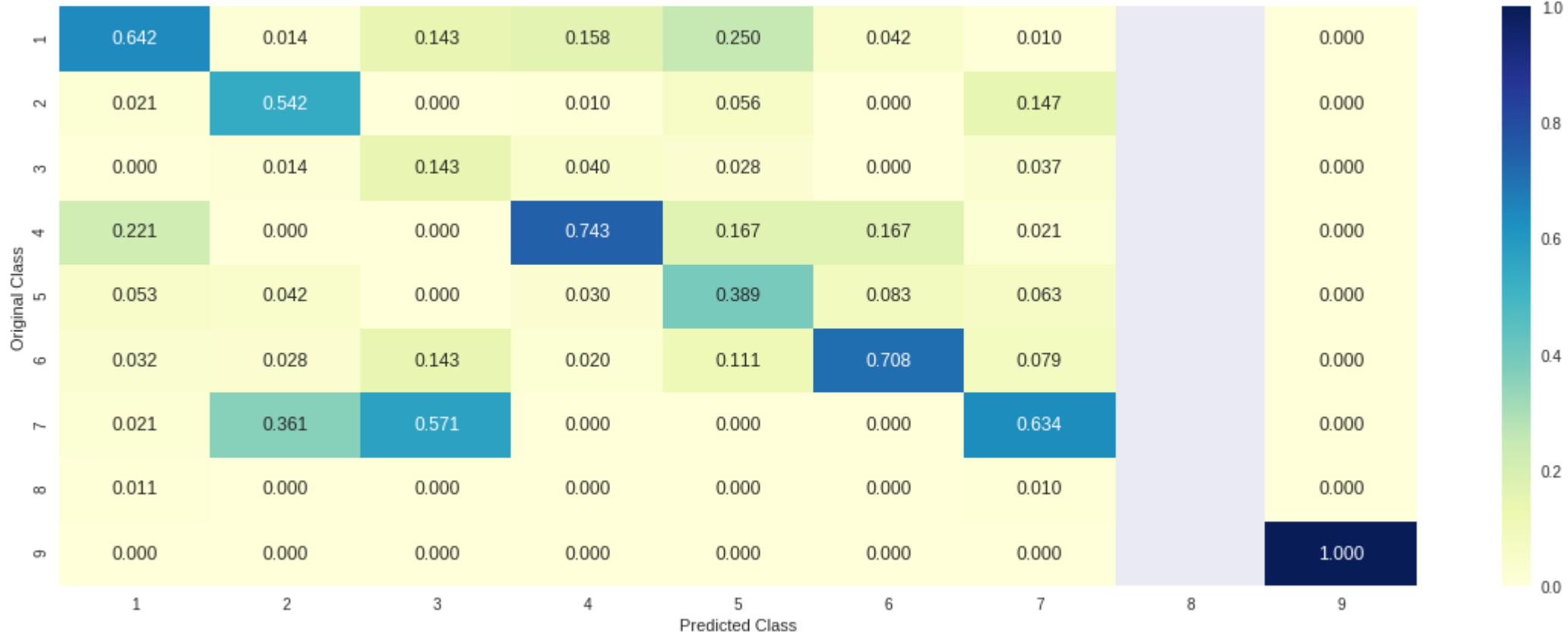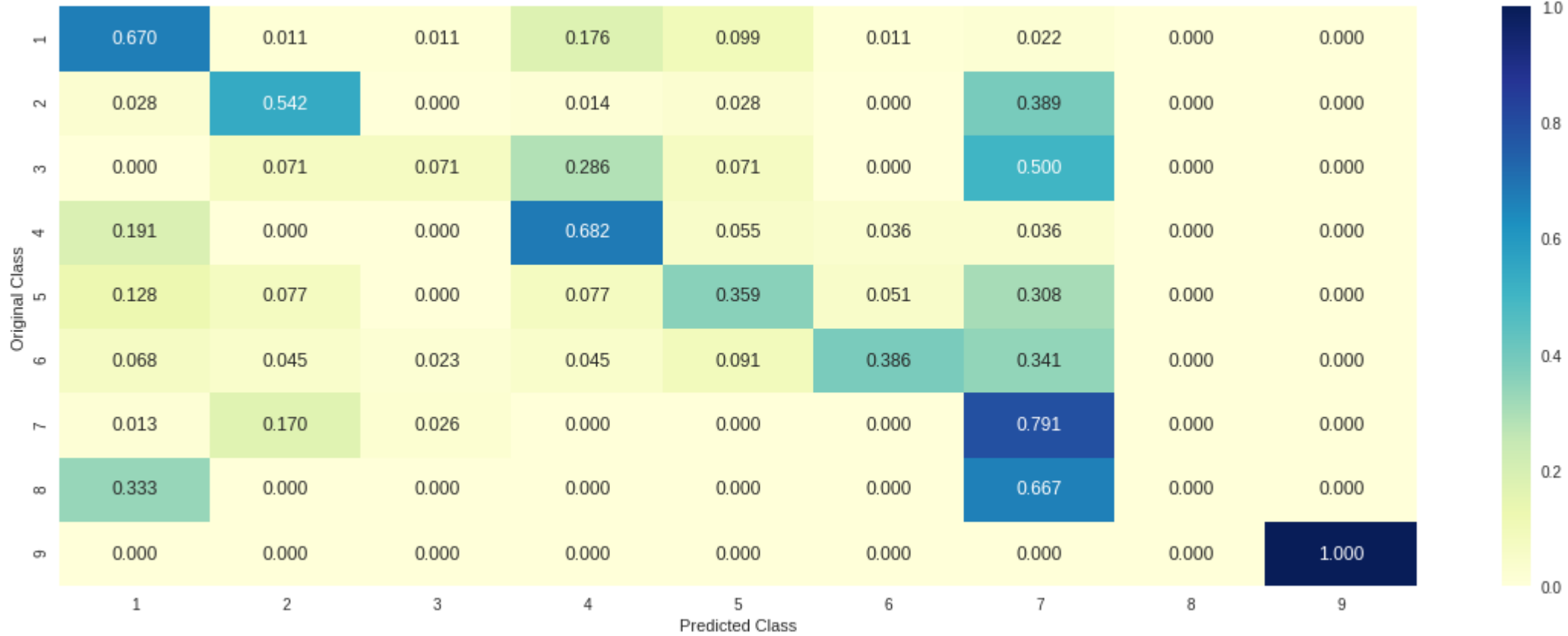
```
Log Loss : 1.1966739622927822
Number of missclassified point : 0.37218045112781956
------------------- Confusion matrix --------------------
```



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



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



**Feature Importance, Correctly classified point**

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.057  0.0456 0.0119 0.7313 0.0366 0.0326 0.0782 0.0032 0.0036]]
Actual Class : 4
-------------------------------------------------
52 Text feature [ala] present in test data point [True]
63 Text feature [abrogated] present in test data point [True]
73 Text feature [asds] present in test data point [True]
74 Text feature [act] present in test data point [True]
Out of the top  100  features  4 are present in query point
```

**Feature Importance, Incorrectly classified point**

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

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0634 0.0735 0.0132 0.0697 0.0409 0.0431 0.6888 0.0035 0.0039]]
Actual Class : 7
-------------------------------------------------
29 Text feature [according] present in test data point [True]
Out of the top  100  features  1 are present in query point
```

# K Nearest Neighbour Classification

## Hyper parameter tuning

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.057  0.0456 0.0119 0.7313 0.0366 0.0326 0.0782 0.0032 0.0036]]
Actual Class : 4
```

```
In [0]:  # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
         # -------------------------
         # default parameter
         # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
         # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

         # methods of
         # fit(X, y) : Fit the model using X as training data and y as target values
         # predict(X):Predict the class labels for the provided data
         # predict_proba(X):Return probability estimates for the test data X.


         alpha = [5, 11, 15, 21, 31, 41, 51, 99]
         cv_log_error_array = []
         for i in alpha:
             print("for alpha =", i)
             clf = KNeighborsClassifier(n_neighbors=i)
             clf.fit(train_x_responseCoding, train_y)
             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_x_responseCoding, train_y)
             sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
             cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
             # to avoid rounding error while multiplying probabilites we use log-probability estimates
             print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


         best_alpha = np.argmin(cv_log_error_array)
         clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
         clf.fit(train_x_responseCoding, train_y)
         sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
         sig_clf.fit(train_x_responseCoding, train_y)

         predict_y = sig_clf.predict_proba(train_x_responseCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(cv_x_responseCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(test_x_responseCoding)
         print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
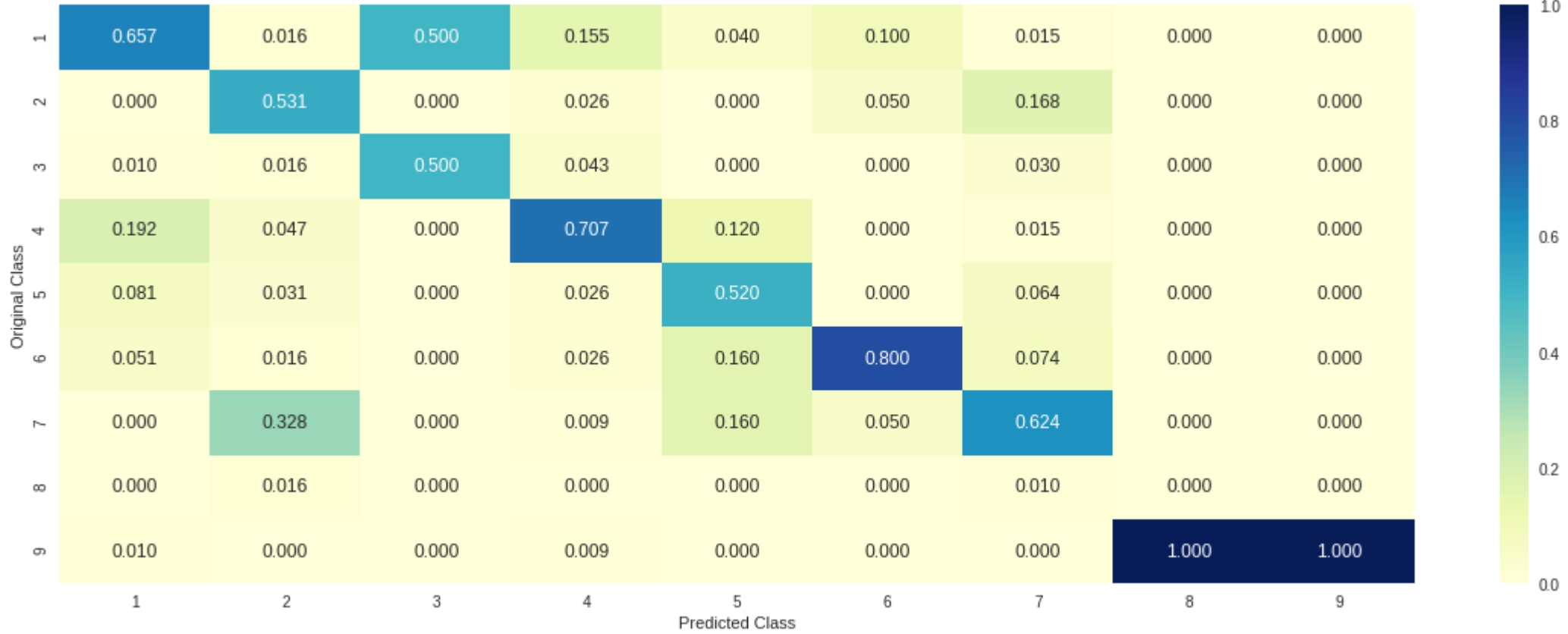
```
for alpha = 5
Log Loss : 1.1291712211953413
for alpha = 11
Log Loss : 1.097591660232571
for alpha = 15
Log Loss : 1.0845563095161594
for alpha = 21
Log Loss : 1.0941827074111
for alpha = 31
Log Loss : 1.1182732265845867
for alpha = 41
Log Loss : 1.1250115332851671
for alpha = 51
Log Loss : 1.14315472896924981
for alpha = 99
Log Loss : 1.1525176829718227
```



```
For values of best alpha =  15 The train log loss is: 0.8670664602899972
For values of best alpha =  15 The cross validation log loss is: 1.0845563095161594
For values of best alpha =  15 The test log loss is: 1.0285146235289795
```

**Testing the model with best hyper paramters**

```
In [0]: # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
        # -------------------------
        # default parameter
        # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
        # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

        # methods of
        # fit(X, y) : Fit the model using X as training data and y as target values
        # predict(X):Predict the class labels for the provided data
        # predict_proba(X):Return probability estimates for the test data X.

        clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
        predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseCoding, cv_y, clf)
```
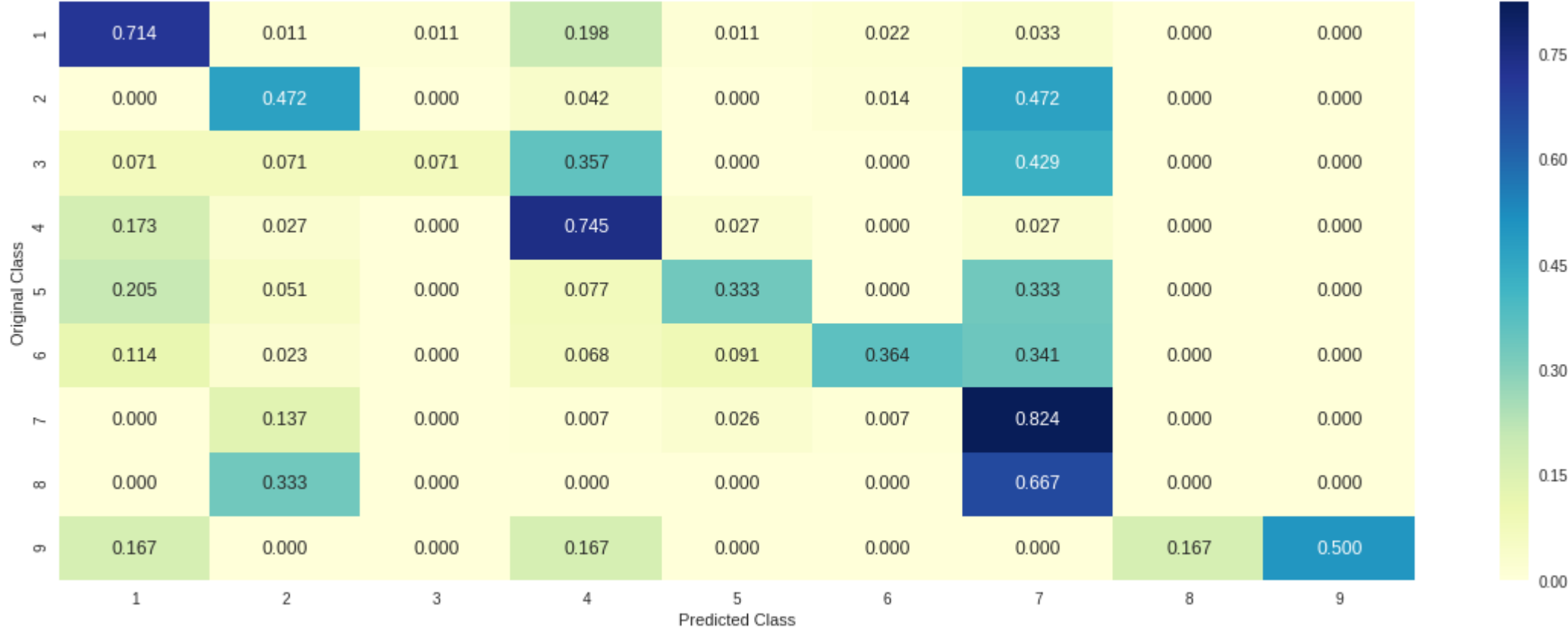
```
Log loss : 1.0845563095161594
Number of mis-classified points : 0.3609022556390977
------------------- Confusion matrix -------------------
```



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



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



### Sample Query point -1

```
In [0]: clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
        clf.fit(train_x_responseCoding, train_y)
        sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
        sig_clf.fit(train_x_responseCoding, train_y)

        test_point_index = 1
        predicted_cls = sig_clf.predict(test_x_responseCoding[0].reshape(1,-1))
        print("Predicted Class :", predicted_cls[0])
        print("Actual Class :", test_y[test_point_index])
        neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
        print("The ",alpha[best_alpha]," nearest neighbours of the test points belongs to classes",train_y[neighbors[1][0]])
        print("Fequency of nearest points :",Counter(train_y[neighbors[1][0]]))
```

```
Predicted Class : 2
Actual Class : 4
The  15  nearest neighbours of the test points belongs to classes [4 4 4 4 4 4 3 4 3 4 4 4 4 4 4]
Fequency of nearest points : Counter({4: 13, 3: 2})
```

### Sample Query Point-2

```
In [0]:   clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
          clf.fit(train_x_responseCoding, train_y)
          sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
          sig_clf.fit(train_x_responseCoding, train_y)

          test_point_index = 100

          predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
          print("Predicted Class :", predicted_cls[0])
          print("Actual Class :", test_y[test_point_index])
          neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
          print("the k value for knn is",alpha[best_alpha],"and the nearest neighbours of the test points belongs to classes",train_y[neighbors[1][0]])
          print("Fequency 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 [7 7 7 7 7 7 7 7 7 6 6 6 6 6]
Fequency of nearest points : Counter({7: 9, 6: 6})
```

## Logistic Regression

### Count Vectorizer with unigrams and bigrams

We apply the same text preprocessing steps used for TFIDF vectorizer applied above, with only changing TFIDFVectorizer to CountVectorizer.

```
In [0]:   # Building a CountVectorizer with all the words that occured minimum 3 times in train data
          text_count_vectorizer = CountVectorizer(min_df=3,ngram_range=(1,2))
          train_text_count_feature_onehotCoding = text_count_vectorizer.fit_transform(train_df['TEXT'])
          # Getting all the feature names (words)
          train_text_count_features= text_count_vectorizer.get_feature_names()

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

          # Zip(list(text_features),text_fea_counts) will zip a word with its number of times it occured
          text_count_fea_dict = dict(zip(list(train_text__count_features),train_text_count_fea_counts))


          #print("Total number of unique words in train data :", len(train_text_count_features))
          #Output>> Total number of unique words in train data : 675913

          #Response coding of text features
          train_text_feature_responseCoding  = get_text_responsecoding(train_df)
          test_text_feature_responseCoding   = get_text_responsecoding(test_df)
          cv_text_feature_responseCoding   = get_text_responsecoding(cv_df)

          # https://stackoverflow.com/a/16202486
          # We convert each row values such that they sum to 1
          train_text_feature_responseCoding = (train_text_feature_responseCoding.T/train_text_feature_responseCoding.sum(axis=1)).T
          test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.sum(axis=1)).T
          cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=1)).T

          # Normalize every feature
          train_text_count_feature_onehotCoding = normalize(train_text_count_feature_onehotCoding, axis=0)

          # We use the same vectorizer that was trained on train data
          test_text_count_feature_onehotCoding = text_count_vectorizer.transform(test_df['TEXT'])
          # Normalize every feature
          test_text_count_feature_onehotCoding = normalize(test_text_count_feature_onehotCoding, axis=0)

          # We use the same vectorizer that was trained on train data
          cv_text_count_feature_onehotCoding = text_count_vectorizer.transform(cv_df['TEXT'])
          # Normalize every feature
          cv_text_count_feature_onehotCoding = normalize(cv_text_feature_onehotCoding, axis=0)

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

          #We use the same feature engineering used for TFIDFVectorization.

          #Apply stacking of all the vectorizations

          train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding,train_variation_feature_onehotCoding,features_train_mat))
          test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding,test_variation_feature_onehotCoding,features_test_mat))
          cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding,cv_variation_feature_onehotCoding,features_cv_mat))

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

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

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


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

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

```
In [0]:   print("One hot encoding features :")
          print("(number of data points * number of features) in train data = ", train_x_onehotCoding.shape)
          print("(number of data points * number of features) in test data = ", test_x_onehotCoding.shape)
          print("(number of data points * number of features) in cross validation data =", cv_x_onehotCoding.shape)
```

```
One hot encoding features :
(number of data points * number of features) in train data =  (2124, 678116)
(number of data points * number of features) in test data =  (665, 678116)
(number of data points * number of features) in cross validation data = (532, 678116)
```

```
In [0]:   print(" Response encoding features :")
          print("(number of data points * number of features) in train data = ", train_x_responseCoding.shape)
          print("(number of data points * number of features) in test data = ", test_x_responseCoding.shape)
          print("(number of data points * number of features) in cross validation data =", cv_x_responseCoding.shape)
```

```
 Response encoding features :
(number of data points * number of features) in train data =  (2124, 31)
(number of data points * number of features) in test data =  (665, 31)
(number of data points * number of features) in cross validation data = (532, 31)
```

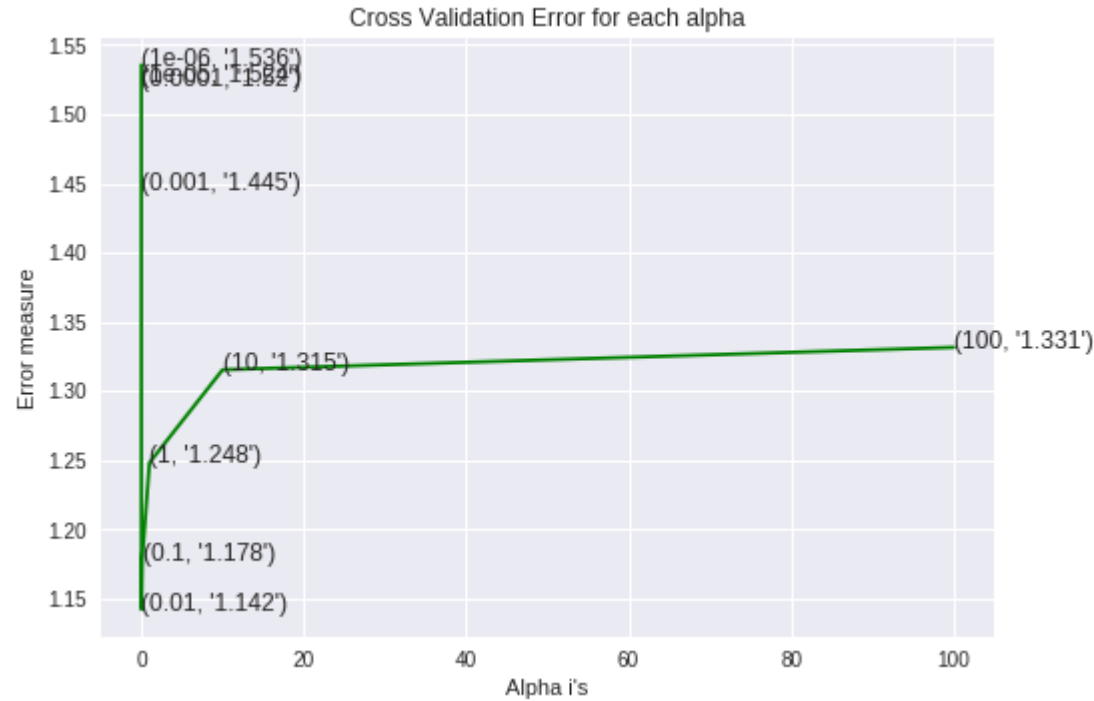**With class balancing**

```
In [0]:  alpha = [10 ** x for x in range(-6, 3)]
         cv_log_error_array = []
         for i in alpha:
             print("for alpha =", i)
             clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
             clf.fit(train_x_onehotCoding, train_y)
             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_x_onehotCoding, train_y)
             sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
             cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
             # to avoid rounding error while multiplying probabilites we use log-probability estimates
             print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


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

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

```
for alpha = 1e-06
Log Loss : 1.5356390696485183
for alpha = 1e-05
Log Loss : 1.5242174163062283
for alpha = 0.0001
Log Loss : 1.520182991130175
for alpha = 0.001
Log Loss : 1.4451895040372738
for alpha = 0.01
Log Loss : 1.142001569594803
for alpha = 0.1
Log Loss : 1.1779294794516042
for alpha = 1
Log Loss : 1.2480104960428926
for alpha = 10
Log Loss : 1.315214860533208
for alpha = 100
Log Loss : 1.3313901431741268
```



```
For values of best alpha =  0.01 The train log loss is: 0.7983328360516818
For values of best alpha =  0.01 The cross validation log loss is: 1.142001569594803
For values of best alpha =  0.01 The test log loss is: 1.1686332768474184
```

```
In [0]: clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
        predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```
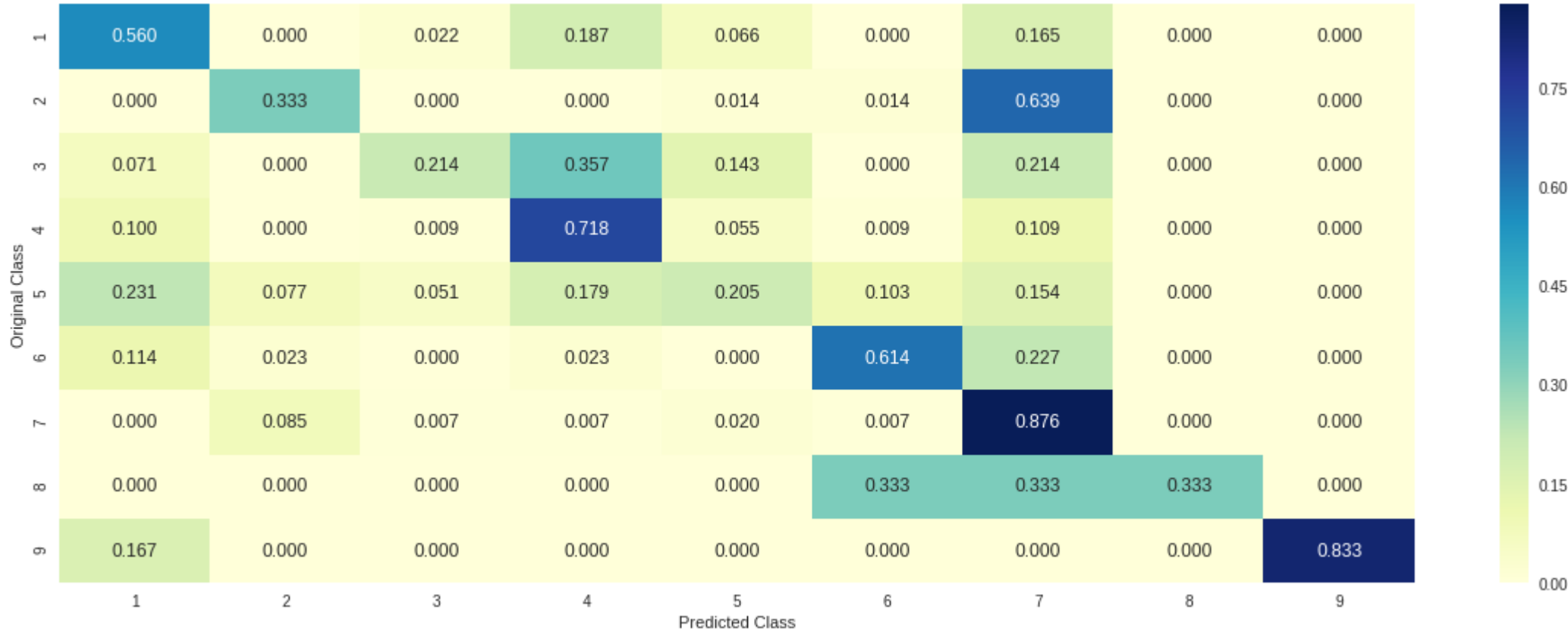
```
Log loss : 1.142001569594803
Number of mis-classified points : 0.37593984962406013
------------------- Confusion matrix --------------------
```



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



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



**Without class balancing**
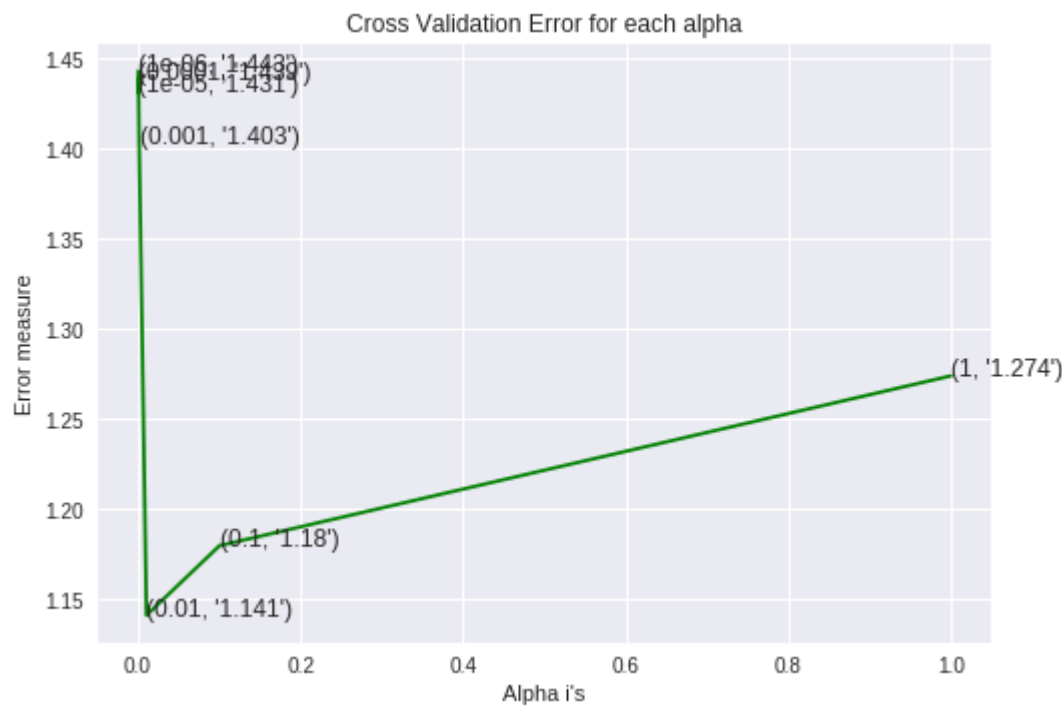
```
In [0]:  alpha = [10 ** x for x in range(-6, 1)]
         cv_log_error_array = []
         for i in alpha:
             print("for alpha =", i)
             clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
             clf.fit(train_x_onehotCoding, train_y)
             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_x_onehotCoding, train_y)
             sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
             cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
             print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


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

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

```
for alpha = 1e-06
Log Loss : 1.4430964093076402
for alpha = 1e-05
Log Loss : 1.4312635126043278
for alpha = 0.0001
Log Loss : 1.4386446359772678
for alpha = 0.001
Log Loss : 1.403095124703155
for alpha = 0.01
Log Loss : 1.1412324489930532
for alpha = 0.1
Log Loss : 1.1798462518590513
for alpha = 1
Log Loss : 1.2739800487558217
```



```
For values of best alpha =  0.01 The train log loss is: 0.786167189725999
For values of best alpha =  0.01 The cross validation log loss is: 1.1412324489930532
For values of best alpha =  0.01 The test log loss is: 1.1866184314042312
```

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

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


        clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
        predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```
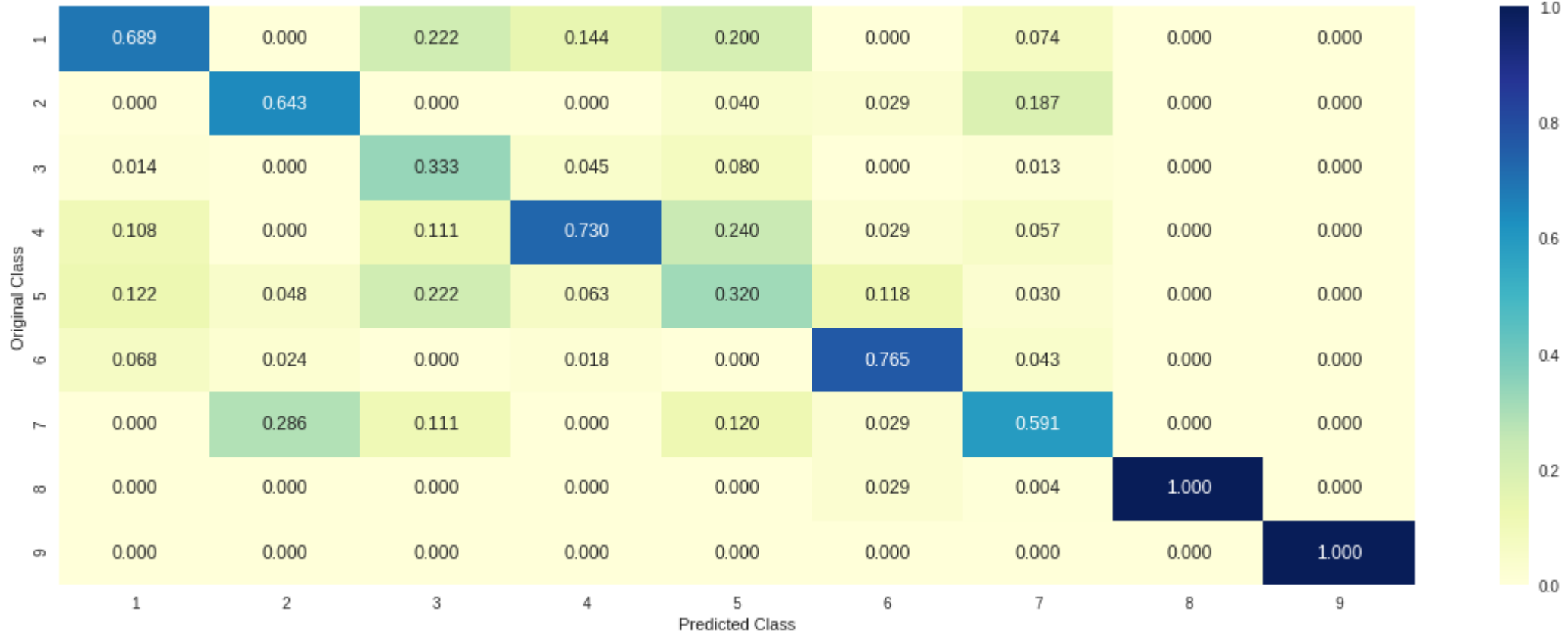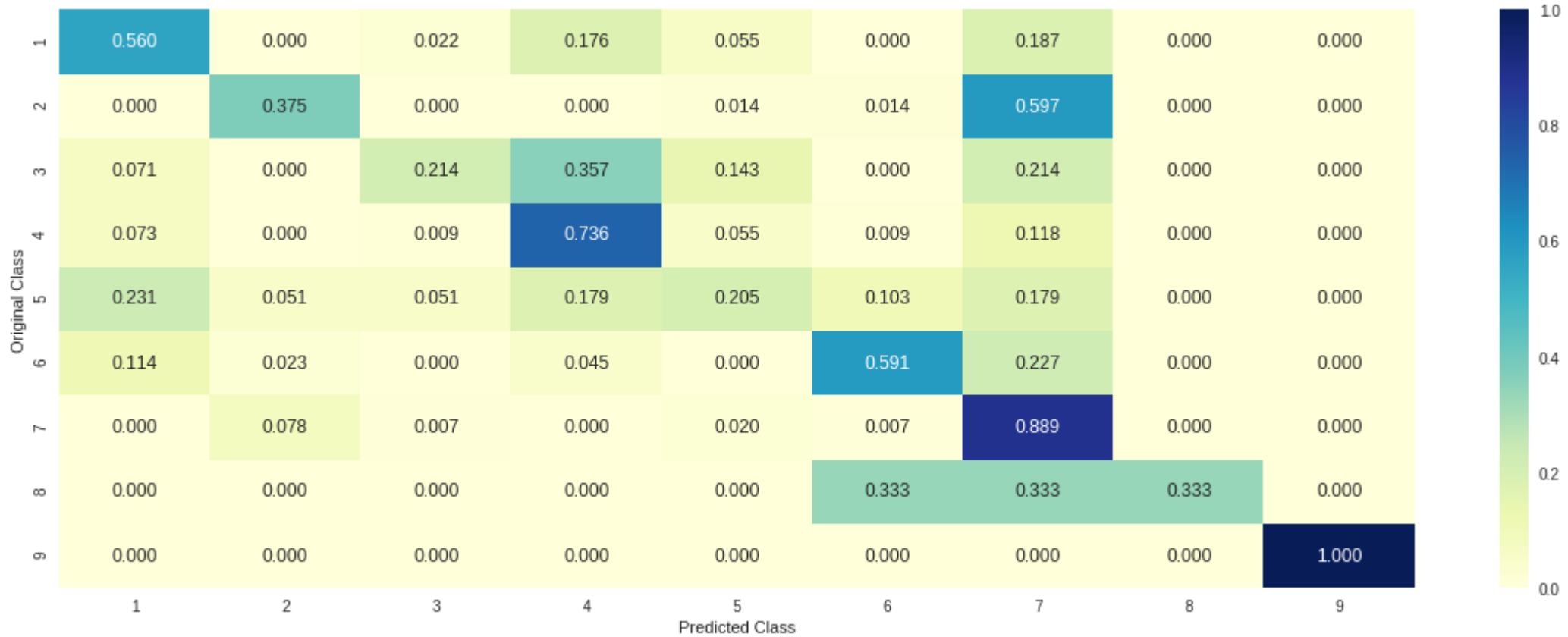
```
Log loss : 1.1412324489930532
Number of mis-classified points : 0.36278195488721804
------------------- Confusion matrix -------------------
```



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



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



**TFIDF vectorizer with unigrams and bigrams**

In [0]:
```python
# building a CountVectorizer with all the words that occured minimum 3 times in train data
text_bigr_vectorizer = TfidfVectorizer(min_df=3,ngram_range=(1,2),max_features=2000)
train_text_bigr_feature_onehotCoding = text_bigr_vectorizer.fit_transform(train_df['TEXT'])
# getting all the feature names (words)
train_text_bigr_features= text_bigr_vectorizer.get_feature_names()

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

# zip(list(text_features),text_fea_counts) will zip a word with its number of times it occured
text_bigr_fea_dict = dict(zip(list(train_text_bigr_features),train_text_bigr_fea_counts))


#print("Total number of unique words in train data :", len(train_text_features))
#Total number of unique words in train data : 2000

#response coding of text features
train_text_feature_responseCoding  = get_text_responsecoding(train_df)
test_text_feature_responseCoding  = get_text_responsecoding(test_df)
cv_text_feature_responseCoding  = get_text_responsecoding(cv_df)
# https://stackoverflow.com/a/16202486
# we convert each row values such that they sum to 1
train_text_feature_responseCoding = (train_text_feature_responseCoding.T/train_text_feature_responseCoding.sum(axis=1)).T
test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.sum(axis=1)).T
cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=1)).T

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

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

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

#https://stackoverflow.com/a/2258273/4084039
sorted_text_bigr_fea_dict = dict(sorted(text_bigr_fea_dict.items(), key=lambda x: x[1] , reverse=True))
sorted_text_bigr_occur = np.array(list(sorted_text_bigr_fea_dict.values()))
```

In [0]:
```python
word_imp_train_bigr = []
text_len_train_bigr = []
for text in train_df.TEXT.values:
    cnt=0
    for word in text.split():
        if word in train_text_bigr_features:
            cnt+=1
    word_imp_train_bigr.append(cnt/len(text.split()))
    text_len_train_bigr.append(len(text.split()))

norm_text_len_train_bigr = [(value-min(text_len_train_bigr))/(max(text_len_train_bigr)-min(text_len_train_bigr)) for value in text_len_train_bigr]

word_imp_test_bigr = []
text_len_test_bigr = []
for text in test_df.TEXT.values:
    cnt=0
    for word in text.split():
        if word in train_text_bigr_features:
            cnt+=1
    word_imp_test_bigr.append(cnt/len(text.split()))
    text_len_test_bigr.append(len(text.split()))

norm_text_len_test_bigr = [(value-min(text_len_test_bigr))/(max(text_len_test_bigr)-min(text_len_test_bigr)) for value in text_len_test_bigr]

word_imp_cv_bigr = []
text_len_cv_bigr = []
for text in cv_df.TEXT.values:
    cnt=0
    for word in text.split():
        if word in train_text_bigr_features:
            cnt+=1
    word_imp_cv_bigr.append(cnt/len(text.split()))
    text_len_cv_bigr.append(len(text.split()))

norm_text_len_cv_bigr = [(value-min(text_len_cv_bigr))/(max(text_len_cv_bigr)-min(text_len_cv_bigr)) for value in text_len_cv_bigr]
```

In [0]:
```python
train_df['norm_text_len_train_bigr'] = norm_text_len_train_bigr
test_df['norm_text_len_test_bigr'] = norm_text_len_test_bigr
cv_df['norm_text_len_cv_bigr'] = norm_text_len_cv_bigr

train_df['word_imp_train_bigr'] = word_imp_train_bigr
cv_df['word_imp_cv_bigr'] = word_imp_cv_bigr
test_df['word_imp_test_bigr'] = word_imp_test_bigr

features_train_bigr = train_df.drop(columns=['ID','Gene','Variation','Class','TEXT','norm_text_len_train','word_imp_train'])
features_test_bigr = test_df.drop(columns=['ID','Gene','Variation','Class','TEXT','norm_text_len_test','word_imp_test'])
features_cv_bigr = cv_df.drop(columns=['ID','Gene','Variation','Class','TEXT','norm_text_len_cv','word_imp_cv'])

features_train_mat_bigr = features_train_bigr.as_matrix()
features_test_mat_bigr = features_test_bigr.as_matrix()
features_cv_mat_bigr = features_cv_bigr.as_matrix()
```

In [0]:
```python
train_gene_var_onehotCoding_bigr = hstack((train_gene_feature_onehotCoding,train_variation_feature_onehotCoding,features_train_mat_bigr))
test_gene_var_onehotCoding_bigr = hstack((test_gene_feature_onehotCoding,test_variation_feature_onehotCoding,features_test_mat_bigr))
cv_gene_var_onehotCoding_bigr = hstack((cv_gene_feature_onehotCoding,cv_variation_feature_onehotCoding,features_cv_mat_bigr))

train_x_onehotCoding_bigr = hstack((train_gene_var_onehotCoding_bigr, train_text_bigr_feature_onehotCoding)).tocsr()
train_y_bigr = np.array(list(train_df['Class']))

test_x_onehotCoding_bigr = hstack((test_gene_var_onehotCoding_bigr, test_text_bigr_feature_onehotCoding)).tocsr()
test_y_bigr = np.array(list(test_df['Class']))

cv_x_onehotCoding_bigr = hstack((cv_gene_var_onehotCoding_bigr, cv_text_bigr_feature_onehotCoding)).tocsr()
cv_y_bigr = np.array(list(cv_df['Class']))


train_gene_var_responseCoding_bigr = np.hstack((train_gene_feature_responseCoding,train_variation_feature_responseCoding,features_train_mat_bigr))
test_gene_var_responseCoding_bigr = np.hstack((test_gene_feature_responseCoding,test_variation_feature_responseCoding,features_test_mat_bigr))
cv_gene_var_responseCoding_bigr = np.hstack((cv_gene_feature_responseCoding,cv_variation_feature_responseCoding,features_cv_mat_bigr))

train_x_responseCoding_bigr = np.hstack((train_gene_var_responseCoding_bigr, train_text_feature_responseCoding))
test_x_responseCoding_bigr = np.hstack((test_gene_var_responseCoding_bigr, test_text_feature_responseCoding))
cv_x_responseCoding_bigr = np.hstack((cv_gene_var_responseCoding_bigr, cv_text_feature_responseCoding))
```

```
In [0]: print("One hot encoding features :")
        print("(number of data points * number of features) in train data = ", train_x_onehotCoding_bigr.shape)
        print("(number of data points * number of features) in test data = ", test_x_onehotCoding_bigr.shape)
        print("(number of data points * number of features) in cross validation data =", cv_x_onehotCoding_bigr.shape)
```

```
One hot encoding features :
(number of data points * number of features) in train data =  (2124, 4200)
(number of data points * number of features) in test data =  (665, 4200)
(number of data points * number of features) in cross validation data = (532, 4200)
```

```
In [0]: print(" Response encoding features :")
        print("(number of data points * number of features) in train data = ", train_x_responseCoding_bigr.shape)
        print("(number of data points * number of features) in test data = ", test_x_responseCoding_bigr.shape)
        print("(number of data points * number of features) in cross validation data =", cv_x_responseCoding_bigr.shape)
```

```
 Response encoding features :
(number of data points * number of features) in train data =  (2124, 31)
(number of data points * number of features) in test data =  (665, 31)
(number of data points * number of features) in cross validation data = (532, 31)
```

**With class balancing**

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

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

        alpha = [10 ** x for x in range(-6, 3)]
        cv_log_error_array = []
        for i in alpha:
            print("for alpha =", i)
            clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
            clf.fit(train_x_onehotCoding_bigr, train_y_bigr)
            sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
            sig_clf.fit(train_x_onehotCoding_bigr, train_y_bigr)
            sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding_bigr)
            cv_log_error_array.append(log_loss(cv_y_bigr, sig_clf_probs, labels=clf.classes_, eps=1e-15))
            # to avoid rounding error while multiplying probabilites we use log-probability estimates
            print("Log Loss :",log_loss(cv_y_bigr, sig_clf_probs))

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

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

        predict_y = sig_clf.predict_proba(train_x_onehotCoding_bigr)
        print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(train_y_bigr, predict_y, labels=clf.classes_, eps=1e-15))
        predict_y = sig_clf.predict_proba(cv_x_onehotCoding_bigr)
        print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(cv_y_bigr, predict_y, labels=clf.classes_, eps=1e-15))
        predict_y = sig_clf.predict_proba(test_x_onehotCoding_bigr)
        print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(test_y_bigr, predict_y, labels=clf.classes_, eps=1e-15))
```
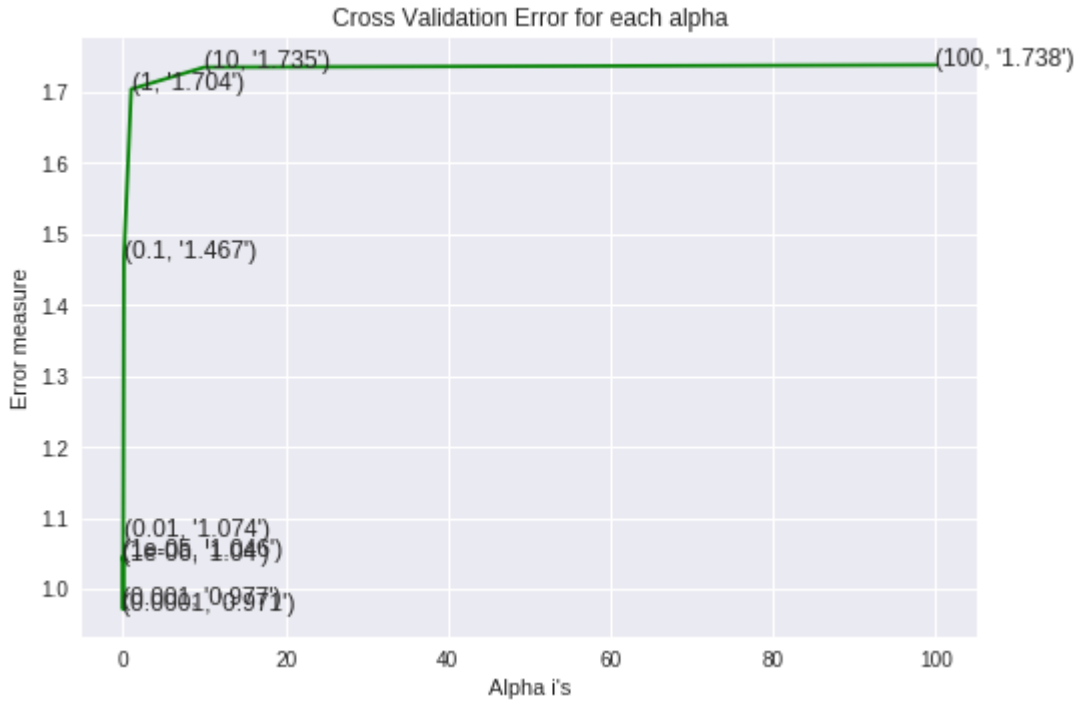
```
for alpha = 1e-06
Log Loss : 1.0403833802514242
for alpha = 1e-05
Log Loss : 1.0459218513585802
for alpha = 0.0001
Log Loss : 0.9711490158812979
for alpha = 0.001
Log Loss : 0.977472372445084
for alpha = 0.01
Log Loss : 1.0739322872420711
for alpha = 0.1
Log Loss : 1.46714259261282256
for alpha = 1
Log Loss : 1.70426609379258
for alpha = 10
Log Loss : 1.7350603404550524
for alpha = 100
Log Loss : 1.7383333015448899
```



Cross Validation Error for each alpha

```
For values of best alpha =  0.0001 The train log loss is: 0.42969881515295244
For values of best alpha =  0.0001 The cross validation log loss is: 0.9711490158812979
For values of best alpha =  0.0001 The test log loss is: 0.9162068466388409
```

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

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

         clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
         predict_and_plot_confusion_matrix(train_x_onehotCoding_bigr, train_y_bigr, cv_x_onehotCoding_bigr, cv_y_bigr, clf)
```
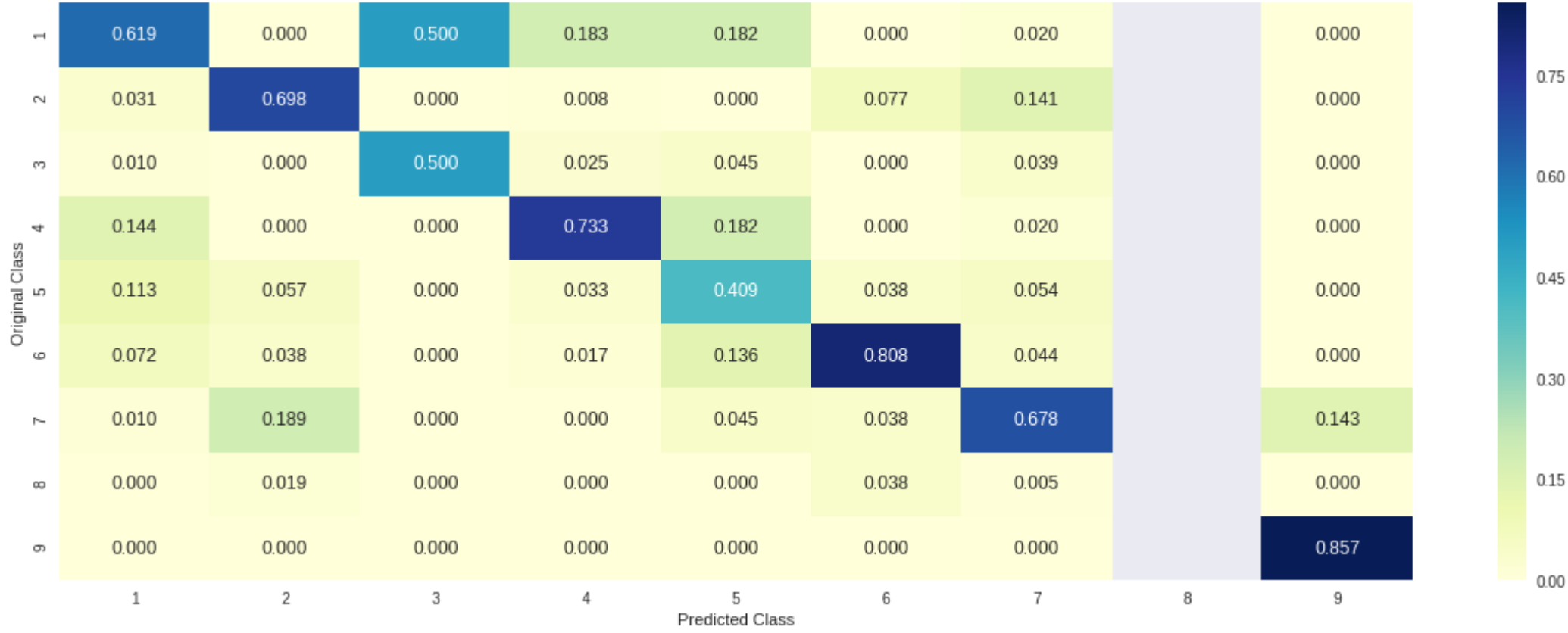
```
Log loss : 0.9711490158812979
Number of mis-classified points : 0.32142857142857145
------------------- Confusion matrix -------------------
```

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 60.000 | 0.000 | 1.000 | 22.000 | 4.000 | 0.000 | 4.000 | 0.000 | 0.000 |
| 2 | 3.000 | 37.000 | 0.000 | 1.000 | 0.000 | 2.000 | 29.000 | 0.000 | 0.000 |
| 3 | 1.000 | 0.000 | 1.000 | 3.000 | 1.000 | 0.000 | 8.000 | 0.000 | 0.000 |
| 4 | 14.000 | 0.000 | 0.000 | 88.000 | 4.000 | 0.000 | 4.000 | 0.000 | 0.000 |
| 5 | 11.000 | 3.000 | 0.000 | 4.000 | 9.000 | 1.000 | 11.000 | 0.000 | 0.000 |
| 6 | 7.000 | 2.000 | 0.000 | 2.000 | 3.000 | 21.000 | 9.000 | 0.000 | 0.000 |
| 7 | 1.000 | 10.000 | 0.000 | 0.000 | 1.000 | 1.000 | 139.000 | 0.000 | 1.000 |
| 8 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 1.000 | 1.000 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 6.000 |

Predicted Class

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

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.619 | 0.000 | 0.500 | 0.183 | 0.182 | 0.000 | 0.020 | | 0.000 |
| 2 | 0.031 | 0.698 | 0.000 | 0.008 | 0.000 | 0.077 | 0.141 | | 0.000 |
| 3 | 0.010 | 0.000 | 0.500 | 0.025 | 0.045 | 0.000 | 0.039 | | 0.000 |
| 4 | 0.144 | 0.000 | 0.000 | 0.733 | 0.182 | 0.000 | 0.020 | | 0.000 |
| 5 | 0.113 | 0.057 | 0.000 | 0.033 | 0.409 | 0.038 | 0.054 | | 0.000 |
| 6 | 0.072 | 0.038 | 0.000 | 0.017 | 0.136 | 0.808 | 0.044 | | 0.000 |
| 7 | 0.010 | 0.189 | 0.000 | 0.000 | 0.045 | 0.038 | 0.678 | | 0.143 |
| 8 | 0.000 | 0.019 | 0.000 | 0.000 | 0.000 | 0.038 | 0.005 | | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | | 0.857 |

Predicted Class

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

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.659 | 0.000 | 0.011 | 0.242 | 0.044 | 0.000 | 0.044 | 0.000 | 0.000 |
| 2 | 0.042 | 0.514 | 0.000 | 0.014 | 0.000 | 0.028 | 0.403 | 0.000 | 0.000 |
| 3 | 0.071 | 0.000 | 0.071 | 0.214 | 0.071 | 0.000 | 0.571 | 0.000 | 0.000 |
| 4 | 0.127 | 0.000 | 0.000 | 0.800 | 0.036 | 0.000 | 0.036 | 0.000 | 0.000 |
| 5 | 0.282 | 0.077 | 0.000 | 0.103 | 0.231 | 0.026 | 0.282 | 0.000 | 0.000 |
| 6 | 0.159 | 0.045 | 0.000 | 0.045 | 0.068 | 0.477 | 0.205 | 0.000 | 0.000 |
| 7 | 0.007 | 0.065 | 0.000 | 0.000 | 0.007 | 0.007 | 0.908 | 0.000 | 0.007 |
| 8 | 0.000 | 0.333 | 0.000 | 0.000 | 0.000 | 0.333 | 0.333 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 |

Predicted Class

*Feature Importance*

```
In [0]:  def get_imp_feature_names(text, indices, removed_ind = []):
             word_present = 0
             tabulte_list = []
             incresingorder_ind = 0
             for i in indices:
                 if i < train_gene_feature_onehotCoding.shape[1]:
                     tabulte_list.append([incresingorder_ind, "Gene", "Yes"])
                 elif i< 18:
                     tabulte_list.append([incresingorder_ind,"Variation", "Yes"])
                 if ((i > 17) & (i not in removed_ind)) :
                     word = train_text_features[i]
                     yes_no = True if word in text.split() else False
                     if yes_no:
                         word_present += 1
                         tabulte_list.append([incresingorder_ind,train_text_features[i], yes_no])
                 incresingorder_ind += 1
             print(word_present, "most importent features are present in our query point")
             print("-"*50)
             print("The features that are most importent of the ",predicted_cls[0]," class:")
             print (tabulate(tabulte_list, headers=["Index",'Feature name', 'Present or Not']))
```

*Correctly Classified point*

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.025  0.0068 0.0323 0.9053 0.0107 0.0035 0.0132 0.0013 0.0019]]
Actual Class : 4
--------------------------------------------------
102 Text feature [abnormalities] present in test data point [True]
175 Text feature [allowed] present in test data point [True]
274 Text feature [amino] present in test data point [True]
366 Text feature [activity] present in test data point [True]
462 Text feature [act] present in test data point [True]
491 Text feature [ala] present in test data point [True]
498 Text feature [along] present in test data point [True]
Out of the top  500  features  7 are present in query point
```

***Incorrectly Classified point***

In [0]:
```python
test_point_index = 100
no_feature = 500
stop=False
while stop==False:
    predicted_cls = sig_clf.predict(test_x_onehotCoding_bigr[test_point_index])
    if int(predicted_cls[0])!=int(test_y_bigr[test_point_index]):
        print("Predicted Class :", predicted_cls[0])
        print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding_bigr[test_point_index]),4))
        print("Actual Class :", test_y_bigr[test_point_index])
        indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
        print("-"*50)
        get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
        stop=True
    else:
        test_point_index+=2
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0305 0.0045 0.3231 0.3519 0.2493 0.0326 0.0041 0.0014 0.0025]]
Actual Class : 3
--------------------------------------------------
175 Text feature [allowed] present in test data point [True]
274 Text feature [amino] present in test data point [True]
326 Text feature [agency] present in test data point [True]
366 Text feature [activity] present in test data point [True]
383 Text feature [appears] present in test data point [True]
429 Text feature [affi] present in test data point [True]
442 Text feature [aberrant] present in test data point [True]
466 Text feature [advantage] present in test data point [True]
498 Text feature [along] present in test data point [True]
Out of the top  500  features  9 are present in query point
```

**Without Class balancing**

In [0]:
```python
alpha = [10 ** x for x in range(-6, 1)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_onehotCoding_bigr, train_y_bigr)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding_bigr, train_y_bigr)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding_bigr)
    cv_log_error_array.append(log_loss(cv_y_bigr, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    print("Log Loss :",log_loss(cv_y_bigr, sig_clf_probs))

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


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

predict_y = sig_clf.predict_proba(train_x_onehotCoding_bigr)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(train_y_bigr, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding_bigr)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(cv_y_bigr, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_onehotCoding_bigr)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(test_y_bigr, predict_y, labels=clf.classes_, eps=1e-15))
```
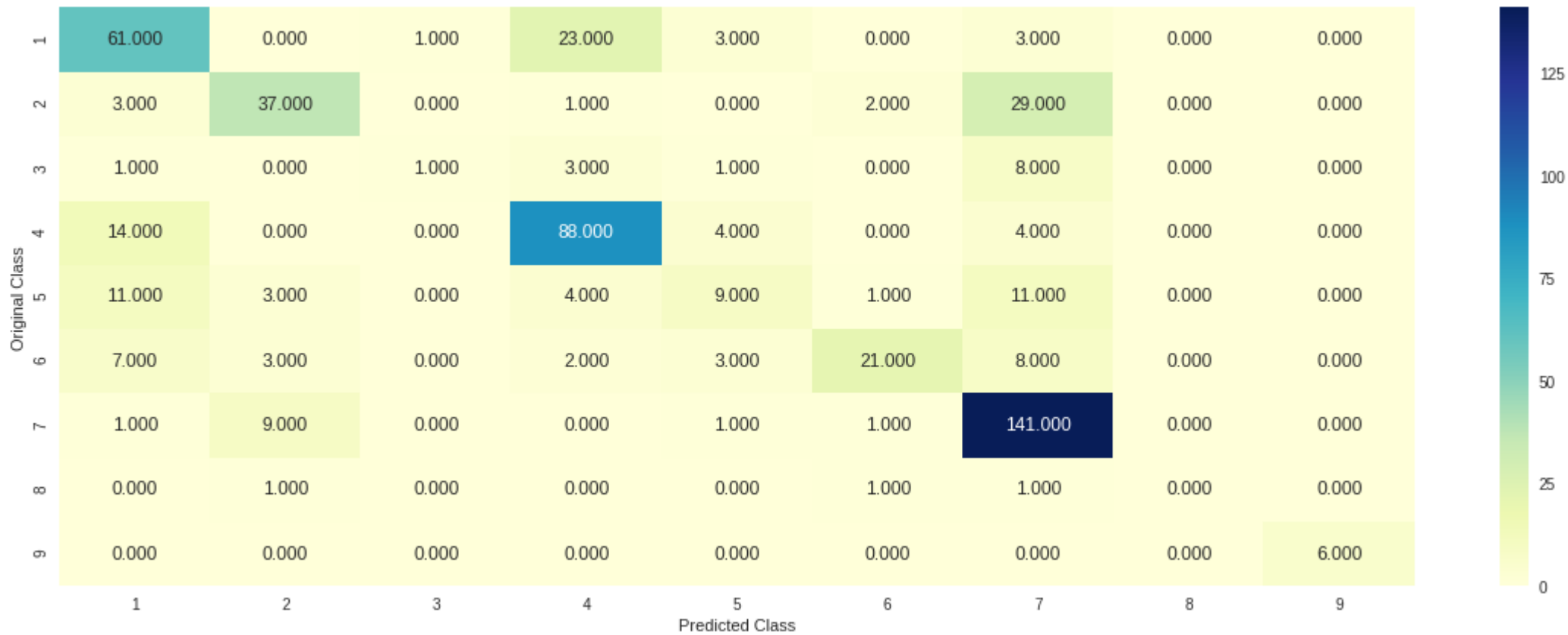
```
for alpha = 1e-06
Log Loss : 1.0629723555633865
for alpha = 1e-05
Log Loss : 1.0392640143675949
for alpha = 0.0001
Log Loss : 0.9728394794352108
for alpha = 0.001
Log Loss : 0.9949999044477307
for alpha = 0.01
Log Loss : 1.1483871506939078
for alpha = 0.1
Log Loss : 1.3857341041129565
for alpha = 1
Log Loss : 1.6799815567394123
```



```
For values of best alpha =  0.0001 The train log loss is: 0.4235676858021788
For values of best alpha =  0.0001 The cross validation log loss is: 0.9728394794352108
For values of best alpha =  0.0001 The test log loss is: 0.9201311286199884
```

```
In [0]:  clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
         predict_and_plot_confusion_matrix(train_x_onehotCoding_bigr, train_y_bigr, cv_x_onehotCoding_bigr, cv_y_bigr, clf)
```
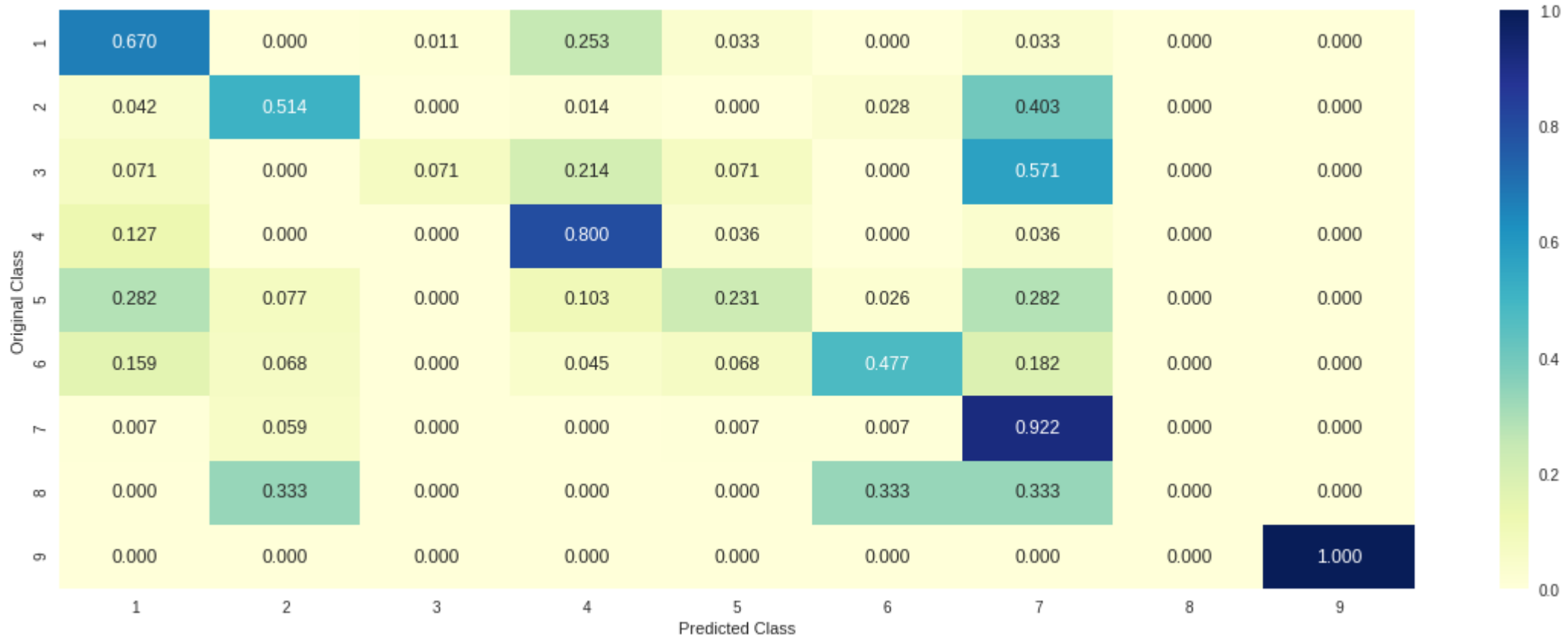
Log loss : 0.9728394794352108
Number of mis-classified points : 0.3157894736842105
------------------- Confusion matrix -------------------



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



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



***Feature Importance, Correctly Classified point***

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

Predicted Class : 4
Predicted Class Probabilities: [[0.025  0.0068 0.0323 0.9053 0.0107 0.0035 0.0132 0.0013 0.0019]]
Actual Class : 4
------------------------------------------------
66 Text feature [abnormalities] present in test data point [True]
144 Text feature [allowed] present in test data point [True]
322 Text feature [amino] present in test data point [True]
378 Text feature [activity] present in test data point [True]
392 Text feature [act] present in test data point [True]
423 Text feature [along] present in test data point [True]
Out of the top  500  features  6 are present in query point

***Feature Importance, Inorrectly Classified point***

```
In [0]:  test_point_index = 100
         no_feature = 500
         stop=False
         while stop==False:
             predicted_cls = sig_clf.predict(test_x_onehotCoding_bigr[test_point_index])
             if int(predicted_cls[0])!=int(test_y_bigr[test_point_index]):
                 print("Predicted Class :", predicted_cls[0])
                 print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding_bigr[test_point_index]),4))
                 print("Actual Class :", test_y_bigr[test_point_index])
                 indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
                 print("-"*50)
                 get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
                 stop=True
             else:
                 test_point_index+=2
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0304 0.0052 0.2449 0.4089 0.2653 0.0363 0.0064 0.0008 0.0019]]
Actual Class : 3
--------------------------------------------------
144 Text feature [allowed] present in test data point [True]
322 Text feature [amino] present in test data point [True]
378 Text feature [activity] present in test data point [True]
423 Text feature [along] present in test data point [True]
440 Text feature [appears] present in test data point [True]
478 Text feature [affi] present in test data point [True]
492 Text feature [agency] present in test data point [True]
Out of the top  500  features  7 are present in query point
```

## TFIDF vectorizer with unigrams

### With SMOTE class balancing

```
In [0]:  from imblearn.over_sampling import SMOTE

         print('The shape of train data before SMOTE: {}'.format(train_x_onehotCoding.shape))
         print("Number of labels before SMOTE: {}\n".format(train_y.shape[0]))

         sm = SMOTE()
         train_x_onehotCoding_smote, train_y_smote = sm.fit_sample(train_x_onehotCoding, train_y)

         print('The shape of train data after SMOTE: {}'.format(train_x_onehotCoding_smote.shape))
         print("Number of labels after SMOTE: {}\n".format(train_y_smote.shape[0]))
```

```
The shape of train data before SMOTE: (2124, 4200)
Number of labels before SMOTE: 2124

The shape of train data after SMOTE: (5481, 4200)
Number of labels after SMOTE: 5481
```
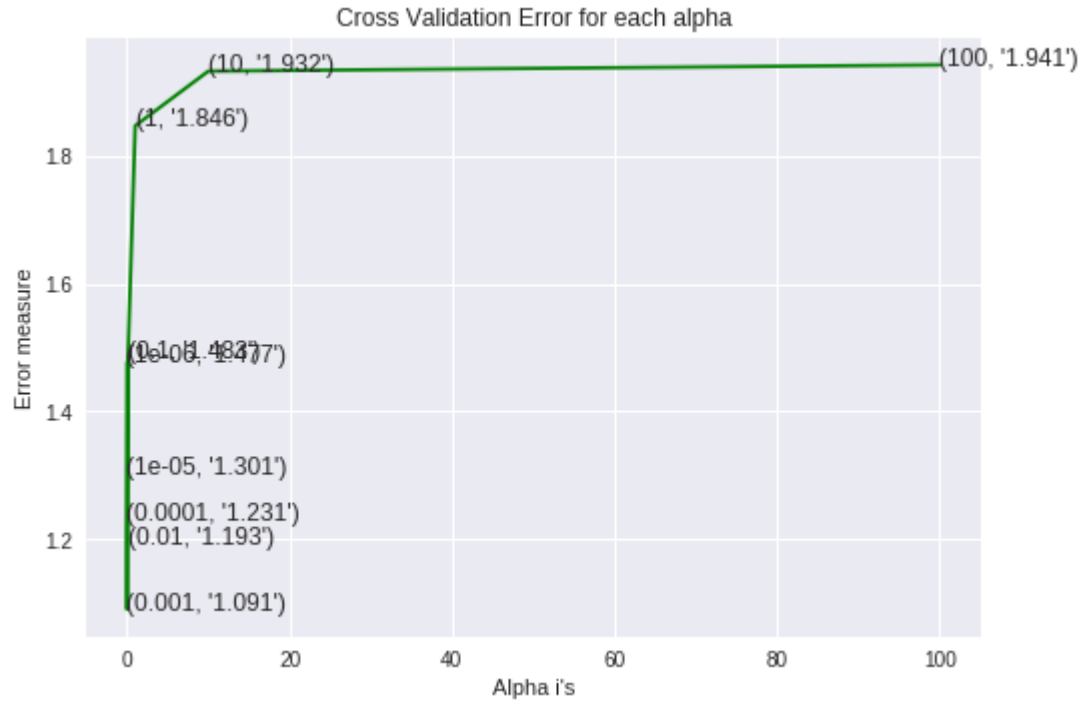
```python
In [0]:  alpha = [10 ** x for x in range(-6, 3)]
         cv_log_error_array = []
         for i in alpha:
             print("for alpha =", i)
             clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
             clf.fit(train_x_onehotCoding_smote, train_y_smote)
             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_x_onehotCoding_smote, train_y_smote)
             sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
             cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
             # to avoid rounding error while multiplying probabilites we use log-probability estimates
             print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


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

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

```
for alpha = 1e-06
Log Loss : 1.4765683060534416
for alpha = 1e-05
Log Loss : 1.3011147656514048
for alpha = 0.0001
Log Loss : 1.2305798010605404
for alpha = 0.001
Log Loss : 1.0908656925773599
for alpha = 0.01
Log Loss : 1.1933725958162464
for alpha = 0.1
Log Loss : 1.4832216534078249
for alpha = 1
Log Loss : 1.846312207019241
for alpha = 10
Log Loss : 1.9316700613929556
for alpha = 100
Log Loss : 1.9414226105485346
```
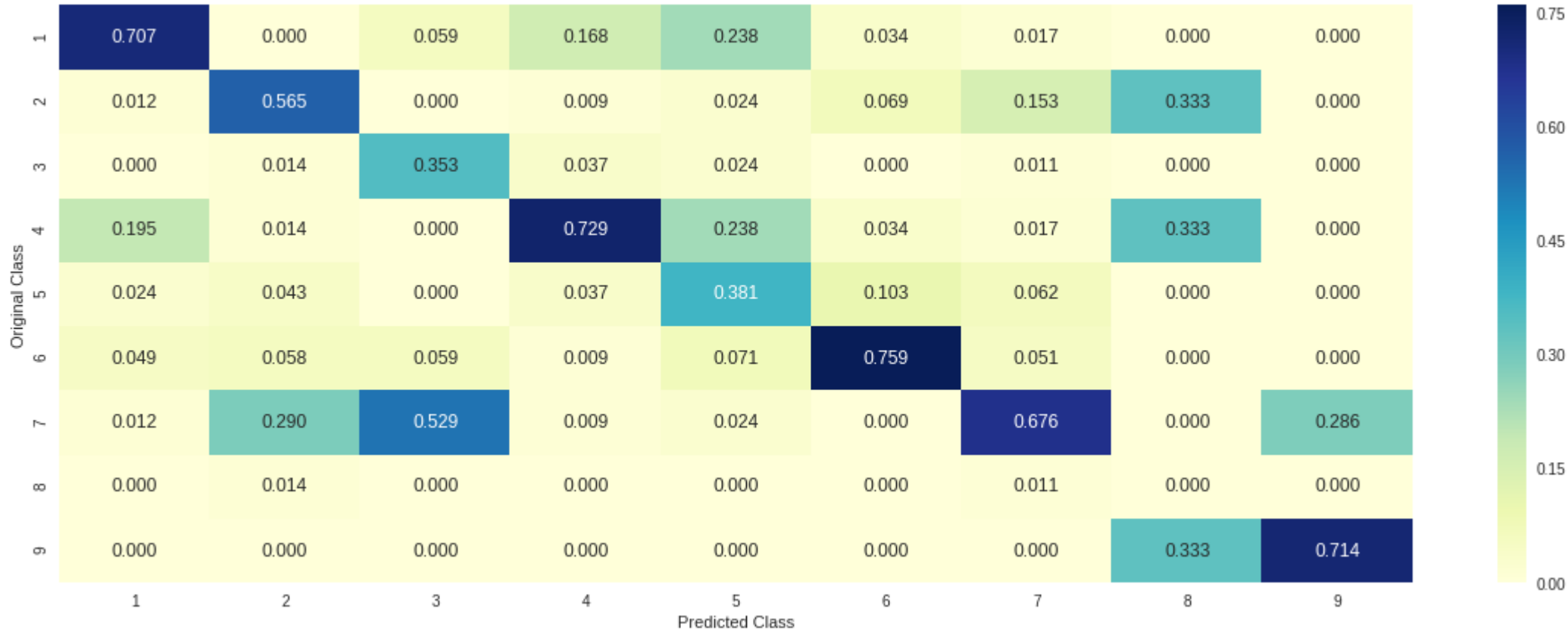


```
For values of best alpha =  0.001 The train log loss is: 0.4275748844582352
For values of best alpha =  0.001 The cross validation log loss is: 1.0908656925773599
For values of best alpha =  0.001 The test log loss is: 1.039852371229378
```

```
In [0]: clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
        predict_and_plot_confusion_matrix(train_x_onehotCoding_smote, train_y_smote, cv_x_onehotCoding, cv_y, clf)
```
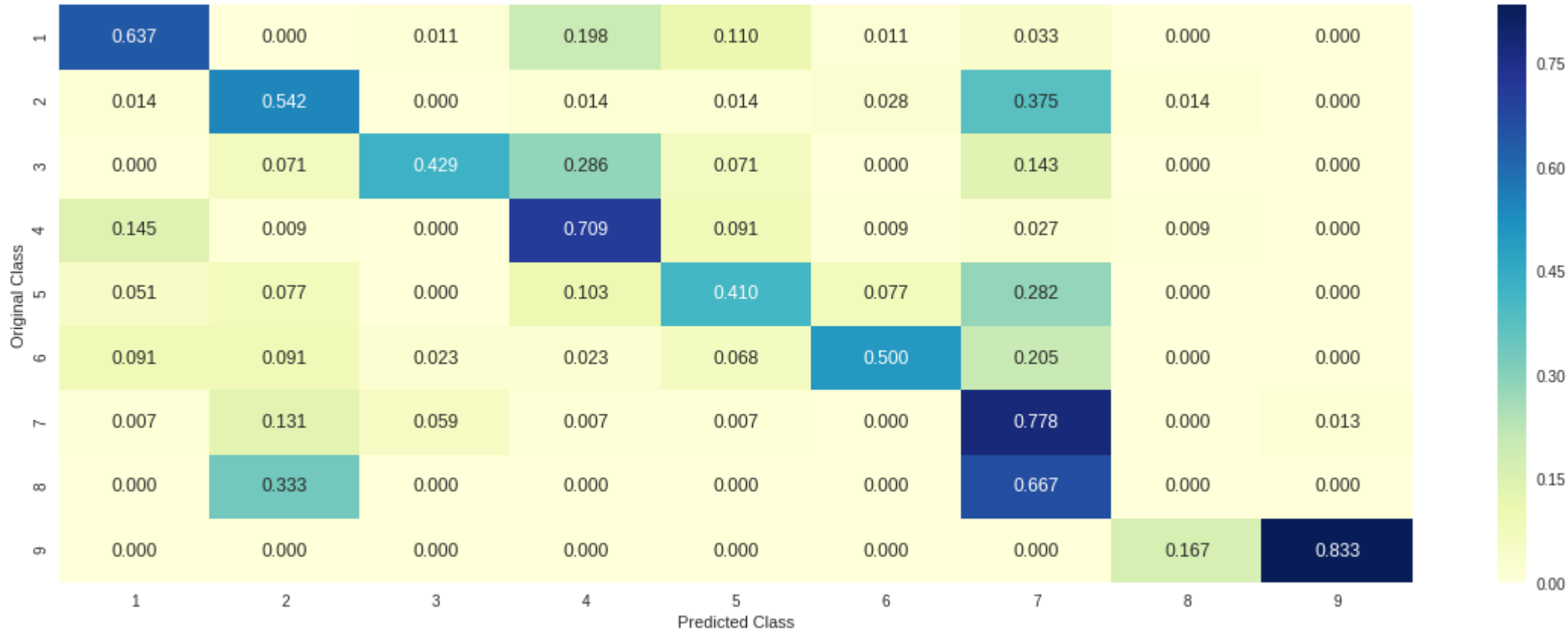
Log loss : 1.0908656925773599
Number of mis-classified points : 0.35526315789473684
------------------- Confusion matrix --------------------



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



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



**With class balancing**

In [0]:
```python
alpha = [10 ** x for x in range(-6, 3)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability estimates
    print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


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

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
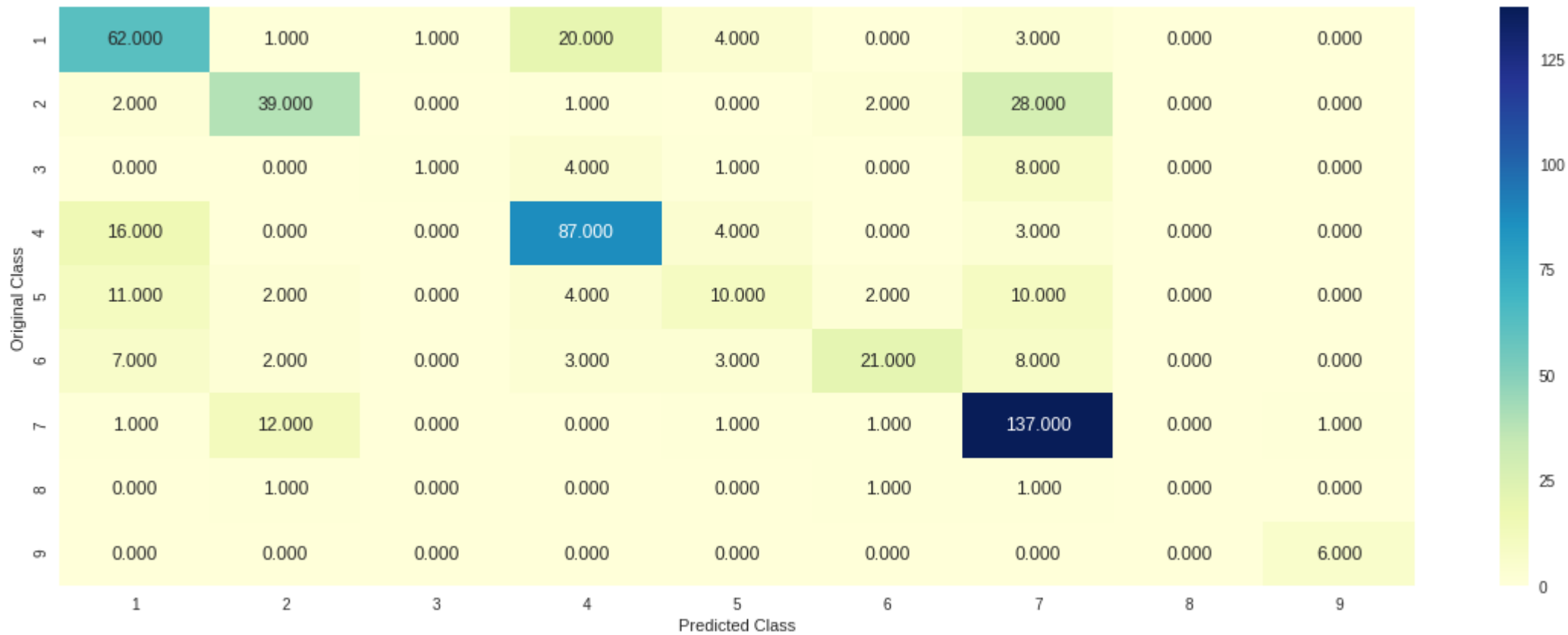
```
for alpha = 1e-06
Log Loss : 1.046283310416022
for alpha = 1e-05
Log Loss : 1.0253212289837434
for alpha = 0.0001
Log Loss : 0.9596416696455822
for alpha = 0.001
Log Loss : 0.9669542539240258
for alpha = 0.01
Log Loss : 1.07018010210269
for alpha = 0.1
Log Loss : 1.4631042646169317
for alpha = 1
Log Loss : 1.7089667768497054
for alpha = 10
Log Loss : 1.7422319736589027
for alpha = 100
Log Loss : 1.74575541645583
```
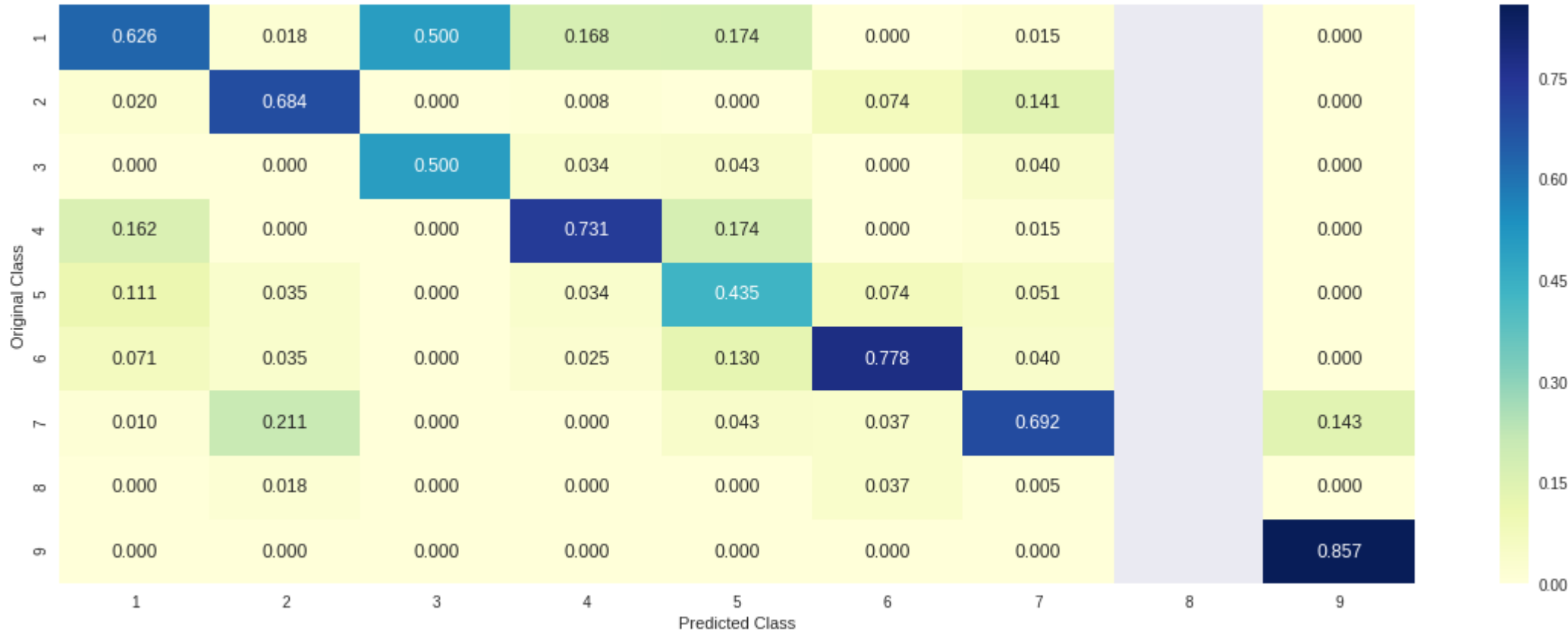


```
For values of best alpha =  0.0001 The train log loss is: 0.42873316349519613
For values of best alpha =  0.0001 The cross validation log loss is: 0.9596416696455822
For values of best alpha =  0.0001 The test log loss is: 0.922006523624239
```

```
In [0]: clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
        predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```
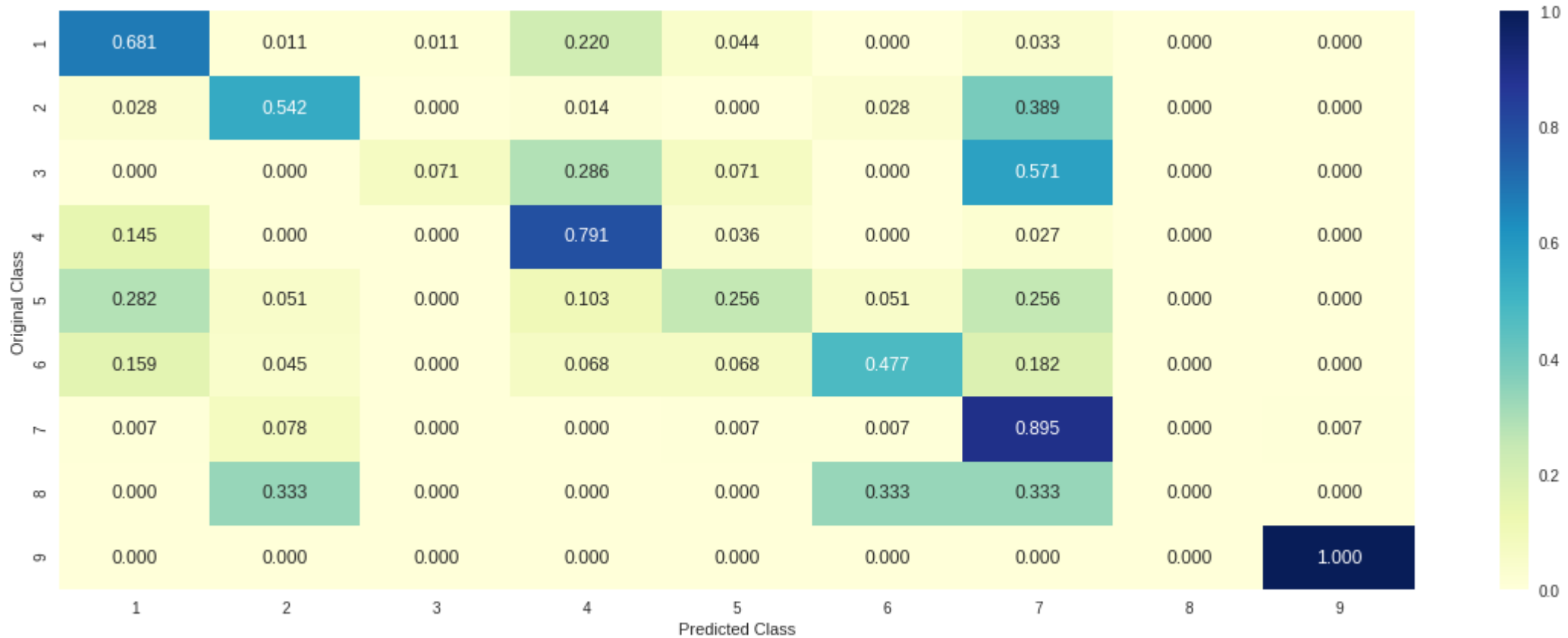
Log loss : 0.9596416696455822
Number of mis-classified points : 0.3176691729323308
------------------- Confusion matrix -------------------



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



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



### Feature Importance

```
In [0]: def get_imp_feature_names(text, indices, removed_ind = []):
            word_present = 0
            tabulte_list = []
            incresingorder_ind = 0
            for i in indices:
                if i < train_gene_feature_onehotCoding.shape[1]:
                    tabulte_list.append([incresingorder_ind, "Gene", "Yes"])
                elif i< 18:
                    tabulte_list.append([incresingorder_ind,"Variation", "Yes"])
                if ((i > 17) & (i not in removed_ind)) :
                    word = train_text_features[i]
                    yes_no = True if word in text.split() else False
                    if yes_no:
                        word_present += 1
                        tabulte_list.append([incresingorder_ind,train_text_features[i], yes_no])
                incresingorder_ind += 1
            print(word_present, "most important features are present in our query point")
            print("-"*50)
            print("The features that are most importent of the ",predicted_cls[0]," class:")
            print (tabulate(tabulte_list, headers=["Index",'Feature name', 'Present or Not']))
```

### Correctly Classified point

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.025  0.0068 0.0323 0.9053 0.0107 0.0035 0.0132 0.0013 0.0019]]
Actual Class : 4
--------------------------------------------------
102 Text feature [abnormalities] present in test data point [True]
175 Text feature [allowed] present in test data point [True]
274 Text feature [amino] present in test data point [True]
366 Text feature [activity] present in test data point [True]
462 Text feature [act] present in test data point [True]
491 Text feature [ala] present in test data point [True]
498 Text feature [along] present in test data point [True]
Out of the top  500  features  7 are present in query point
```

***Incorrectly Classified point***

```python
test_point_index = 100
no_feature = 500
stop=False
while stop==False:
    predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
    if int(predicted_cls[0])!=int(test_y[test_point_index]):
        print("Predicted Class :", predicted_cls[0])
        print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
        print("Actual Class :", test_y[test_point_index])
        indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
        print("-"*50)
        get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
        stop=True
    else:
        test_point_index+=2
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0305 0.0045 0.3231 0.3519 0.2493 0.0326 0.0041 0.0014 0.0025]]
Actual Class : 3
--------------------------------------------------
175 Text feature [allowed] present in test data point [True]
274 Text feature [amino] present in test data point [True]
326 Text feature [agency] present in test data point [True]
366 Text feature [activity] present in test data point [True]
383 Text feature [appears] present in test data point [True]
429 Text feature [affi] present in test data point [True]
442 Text feature [aberrant] present in test data point [True]
466 Text feature [advantage] present in test data point [True]
498 Text feature [along] present in test data point [True]
Out of the top  500  features  9 are present in query point
```
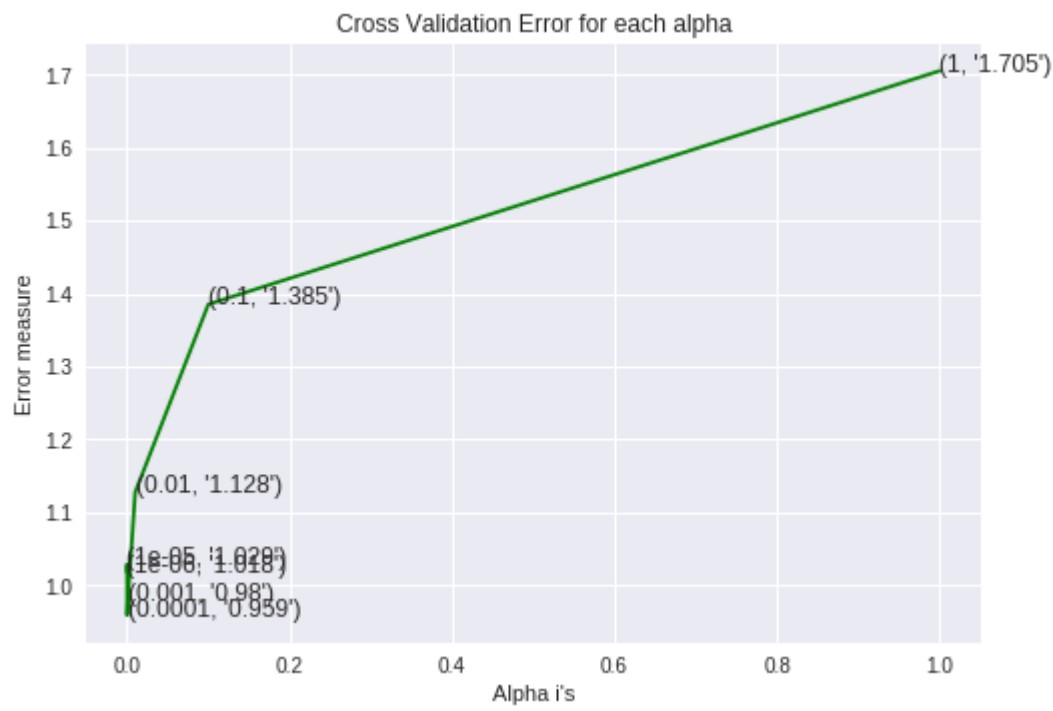
**Without Class balancing**

In [0]:
```python
alpha = [10 ** x for x in range(-6, 1)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


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

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

```
for alpha = 1e-06
Log Loss : 1.01800285112994
for alpha = 1e-05
Log Loss : 1.0287623635674068
for alpha = 0.0001
Log Loss : 0.9588971206703185
for alpha = 0.001
Log Loss : 0.9802808741194752
for alpha = 0.01
Log Loss : 1.1275786754451347
for alpha = 0.1
Log Loss : 1.3850841888718475
for alpha = 1
Log Loss : 1.7049526078637938
```
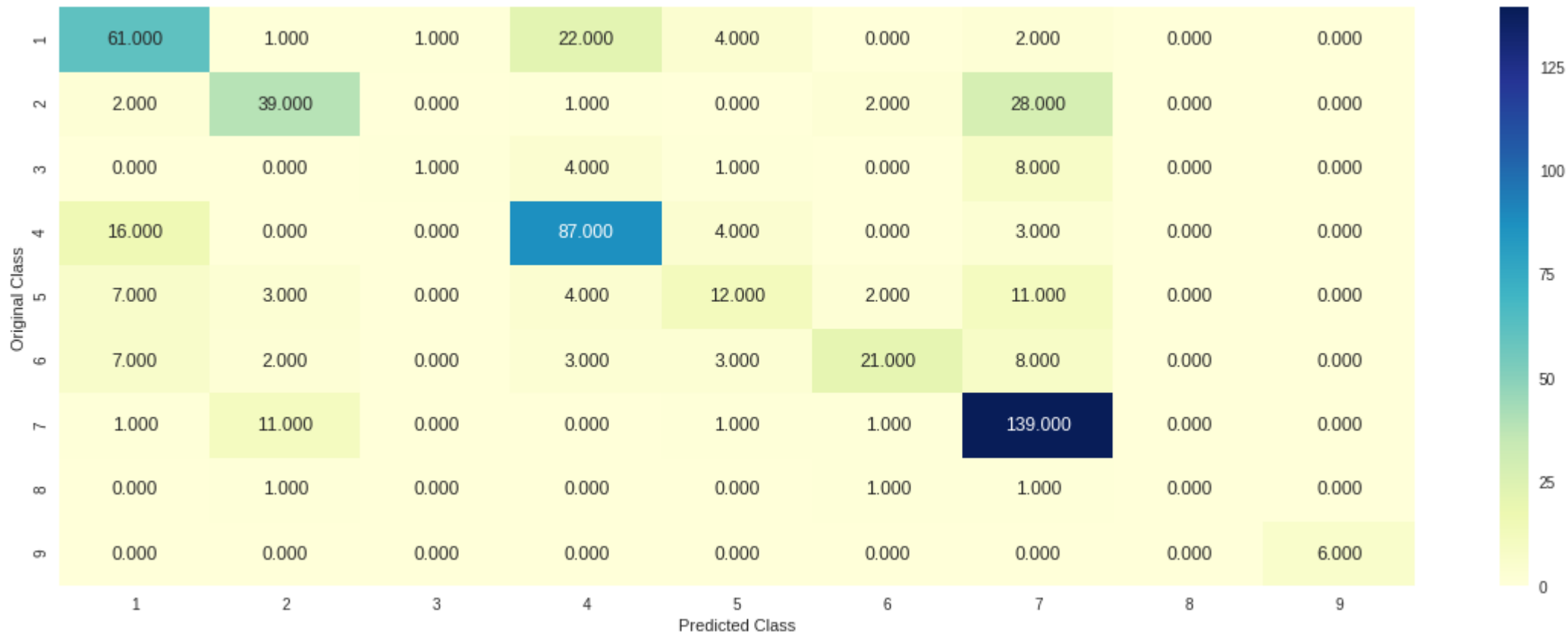


```
For values of best alpha =  0.0001 The train log loss is: 0.42182121033130965
For values of best alpha =  0.0001 The cross validation log loss is: 0.9588971206703185
For values of best alpha =  0.0001 The test log loss is: 0.9268173226303703
```
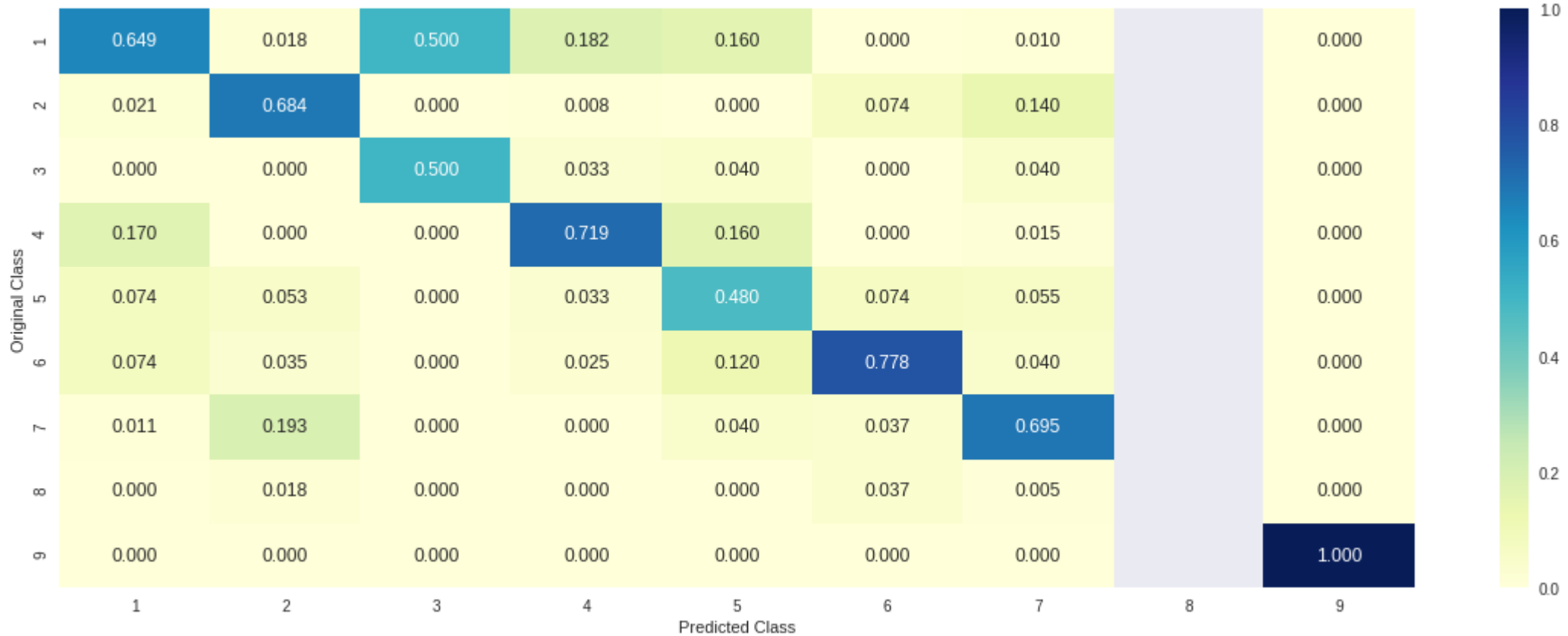
**Testing model with best hyper parameters**

```
In [0]:   clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
          predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```
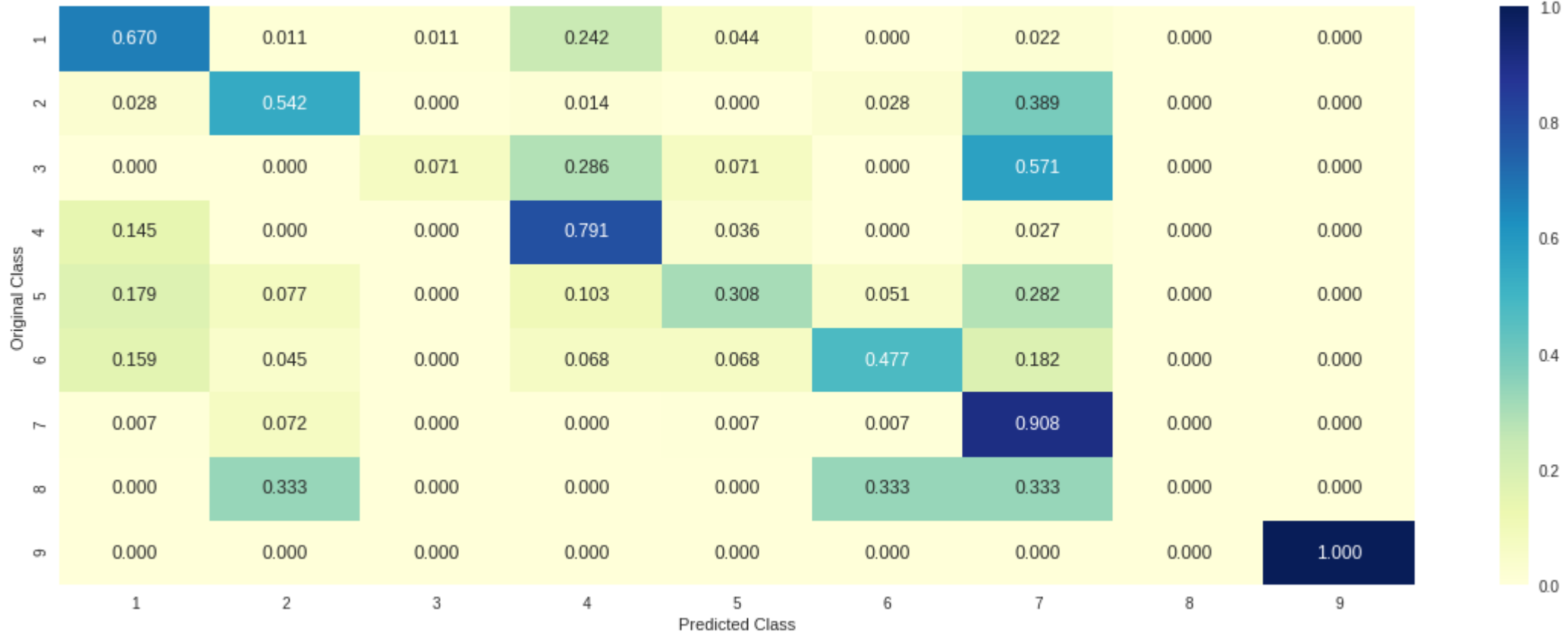
Log loss : 0.9588971206703185
Number of mis-classified points : 0.31203007518796994
-------------------- Confusion matrix --------------------



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



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



*Feature Importance, Correctly Classified point*

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

Predicted Class : 4
Predicted Class Probabilities: [[0.025  0.0068 0.0323 0.9053 0.0107 0.0035 0.0132 0.0013 0.0019]]
Actual Class : 4
------------------------------------------------
66 Text feature [abnormalities] present in test data point [True]
144 Text feature [allowed] present in test data point [True]
322 Text feature [amino] present in test data point [True]
378 Text feature [activity] present in test data point [True]
392 Text feature [act] present in test data point [True]
423 Text feature [along] present in test data point [True]
Out of the top  500  features  6 are present in query point

*Feature Importance, Inorrectly Classified point*

In [0]:
```python
test_point_index = 100
no_feature = 500
stop=False
while stop==False:
    predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
    if int(predicted_cls[0])!=int(test_y[test_point_index]):
        print("Predicted Class :", predicted_cls[0])
        print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
        print("Actual Class :", test_y[test_point_index])
        indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
        print("-"*50)
        get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
        stop=True
    else:
        test_point_index+=2
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0304 0.0052 0.2449 0.4089 0.2653 0.0363 0.0064 0.0008 0.0019]]
Actual Class : 3
--------------------------------------------------
144 Text feature [allowed] present in test data point [True]
322 Text feature [amino] present in test data point [True]
378 Text feature [activity] present in test data point [True]
423 Text feature [along] present in test data point [True]
440 Text feature [appears] present in test data point [True]
478 Text feature [affi] present in test data point [True]
492 Text feature [agency] present in test data point [True]
Out of the top  500  features   7 are present in query point
```

## Linear Support Vector Machines

In [0]:
```python
test_point_index = 100
no_feature = 500
stop=False
while stop==False:
    predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
    if int(predicted_cls[0])!=int(test_y[test_point_index]):
        print("Predicted Class :", predicted_cls[0])
        print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
        print("Actual Class :", test_y[test_point_index])
        indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
        print("-"*50)
        get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
        stop=True
    else:
        test_point_index+=2
```

```python
In [0]:  # read more about support vector machines with linear kernals here http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html

         # -------------------------------
         # default parameters
         # SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
         # cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_state=None)

         # Some of methods of SVM()
         # fit(X, y, [sample_weight])    Fit the SVM model according to the given training data.
         # predict(X)    Perform classification on samples in X.


         alpha = [10 ** x for x in range(-5, 3)]
         cv_log_error_array = []
         for i in alpha:
             print("for C =", i)
             #clf = SVC(C=i,kernel='linear',probability=True, class_weight='balanced')
             clf = SGDClassifier( class_weight='balanced', alpha=i, penalty='l2', loss='hinge', random_state=42)
             clf.fit(train_x_onehotCoding, train_y)
             sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
             sig_clf.fit(train_x_onehotCoding, train_y)
             sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
             cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
             print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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


         best_alpha = np.argmin(cv_log_error_array)
         # clf = SVC(C=i,kernel='linear',probability=True, class_weight='balanced')
         clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
         clf.fit(train_x_onehotCoding, train_y)
         sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
         sig_clf.fit(train_x_onehotCoding, train_y)

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

```
for C = 1e-05
Log Loss : 1.0577074289349597
for C = 0.0001
Log Loss : 1.027644661829584
for C = 0.001
Log Loss : 1.0228544865940221
for C = 0.01
Log Loss : 1.1203767905244797
for C = 0.1
Log Loss : 1.4660917415254184
for C = 1
Log Loss : 1.7462306563117043
for C = 10
Log Loss : 1.7462360195024083
for C = 100
Log Loss : 1.7462360232909615
```
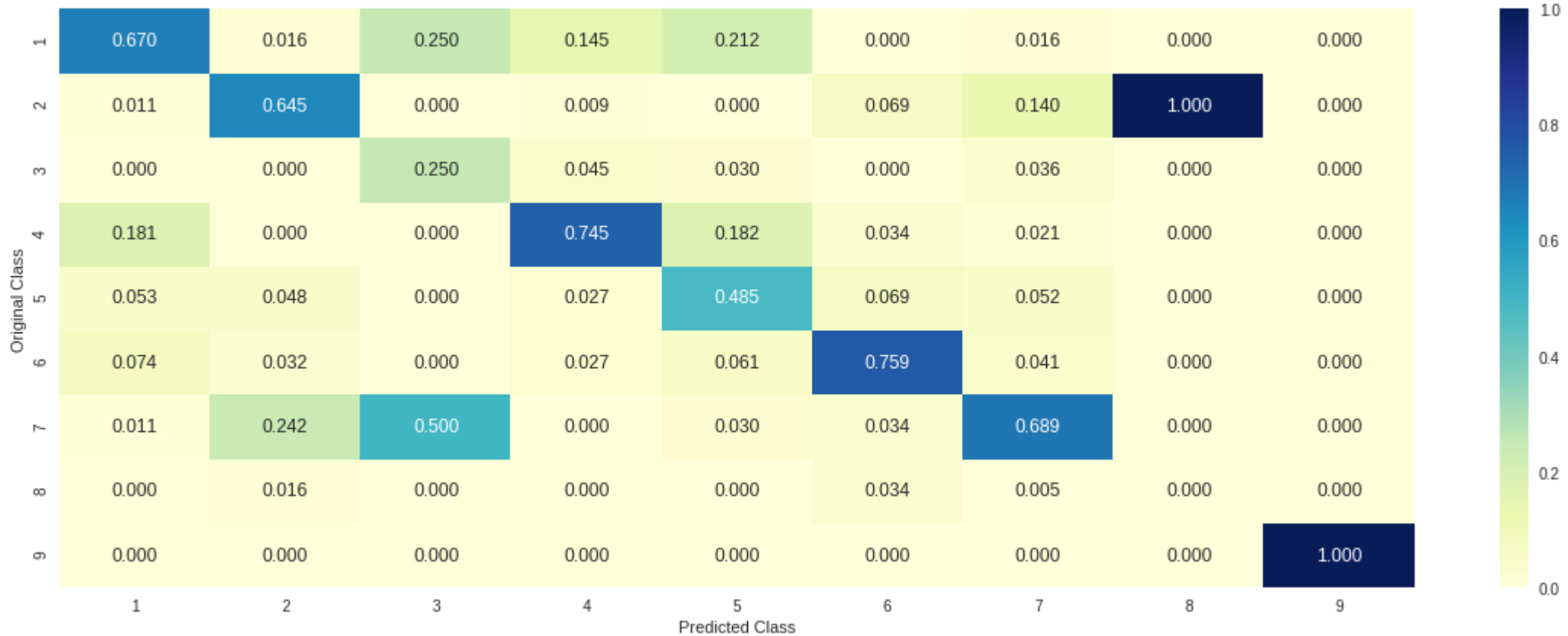


```
For values of best alpha =  0.001 The train log loss is: 0.5417333673200837
For values of best alpha =  0.001 The cross validation log loss is: 1.0228544865940221
For values of best alpha =  0.001 The test log loss is: 0.9932048482590475
```

```
In [0]: clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42,class_weight='balanced')
        predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf)
```

```
Log loss : 1.0228544865940221
Number of mis-classified points : 0.3176691729323308
------------------- Confusion matrix -------------------
```



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



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



**Feature importance for Correctly classified point**

```
In [0]: clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
        clf.fit(train_x_onehotCoding,train_y)
        test_point_index = 1
        # test_point_index = 100
        no_feature = 500
        predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
        print("Predicted Class :", predicted_cls[0])
        print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
        print("Actual Class :", test_y[test_point_index])
        indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
        print("-"*50)
        get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0376 0.0227 0.0243 0.8407 0.0166 0.0083 0.045  0.0022 0.0026]]
Actual Class : 4
------------------------------------------------
229 Text feature [allowed] present in test data point [True]
320 Text feature [amino] present in test data point [True]
364 Text feature [allows] present in test data point [True]
460 Text feature [agreement] present in test data point [True]
465 Text feature [abnormalities] present in test data point [True]
476 Text feature [ala] present in test data point [True]
492 Text feature [affect] present in test data point [True]
Out of the top  500  features  7 are present in query point
```

**Feature importance for incorrectly classified point**

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

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0768 0.1011 0.0208 0.0728 0.0257 0.1414 0.55   0.0046 0.0067]]
Actual Class : 7
--------------------------------------------------
78 Text feature [allow] present in test data point [True]
416 Text feature [amplify] present in test data point [True]
427 Text feature [aggregations] present in test data point [True]
442 Text feature [al] present in test data point [True]
456 Text feature [abbreviations] present in test data point [True]
Out of the top  500  features  5 are present in query point
```

## Random Forest Classifier

### One hot encoding

```
In [0]:  # --------------------------------
         # default parameters
         # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
         # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
         # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
         # class_weight=None)

         # Some of methods of RandomForestClassifier()
         # fit(X, y, [sample_weight])    Fit the SVM model according to the given training data.
         # predict(X)    Perform classification on samples in X.
         # predict_proba (X) Perform classification on samples in X.

         # some of attributes of  RandomForestClassifier()
         # feature_importances_ : array of shape = [n_features]
         # The feature importances (the higher, the more important the feature).


         alpha = [100,200,500,1000,2000]
         max_depth = [5, 10]
         cv_log_error_array = []
         for i in alpha:
             for j in max_depth:
                 print("for n_estimators =", i,"and max depth = ", j)
                 clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
                 clf.fit(train_x_onehotCoding, train_y)
                 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
                 sig_clf.fit(train_x_onehotCoding, train_y)
                 sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
                 cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
                 print("Log Loss :",log_loss(cv_y, sig_clf_probs))

         '''fig, ax = plt.subplots()
         features = np.dot(np.array(alpha)[:,None],np.array(max_depth)[None]).ravel()
         ax.plot(features, cv_log_error_array,c='g')
         for i, txt in enumerate(np.round(cv_log_error_array,3)):
             ax.annotate((alpha[int(i/2)],max_depth[int(i%2)],str(txt)), (features[i],cv_log_error_array[i]))
         plt.grid()
         plt.title("Cross Validation Error for each alpha")
         plt.xlabel("Alpha i's")
         plt.ylabel("Error measure")
         plt.show()
         '''

         best_alpha = np.argmin(cv_log_error_array)
         clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha%2)], random_state=42, n_jobs=-1)
         clf.fit(train_x_onehotCoding, train_y)
         sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
         sig_clf.fit(train_x_onehotCoding, train_y)

         predict_y = sig_clf.predict_proba(train_x_onehotCoding)
         print('For values of best estimator = ', alpha[int(best_alpha/2)], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
         print('For values of best estimator = ', alpha[int(best_alpha/2)], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(test_x_onehotCoding)
         print('For values of best estimator = ', alpha[int(best_alpha/2)], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
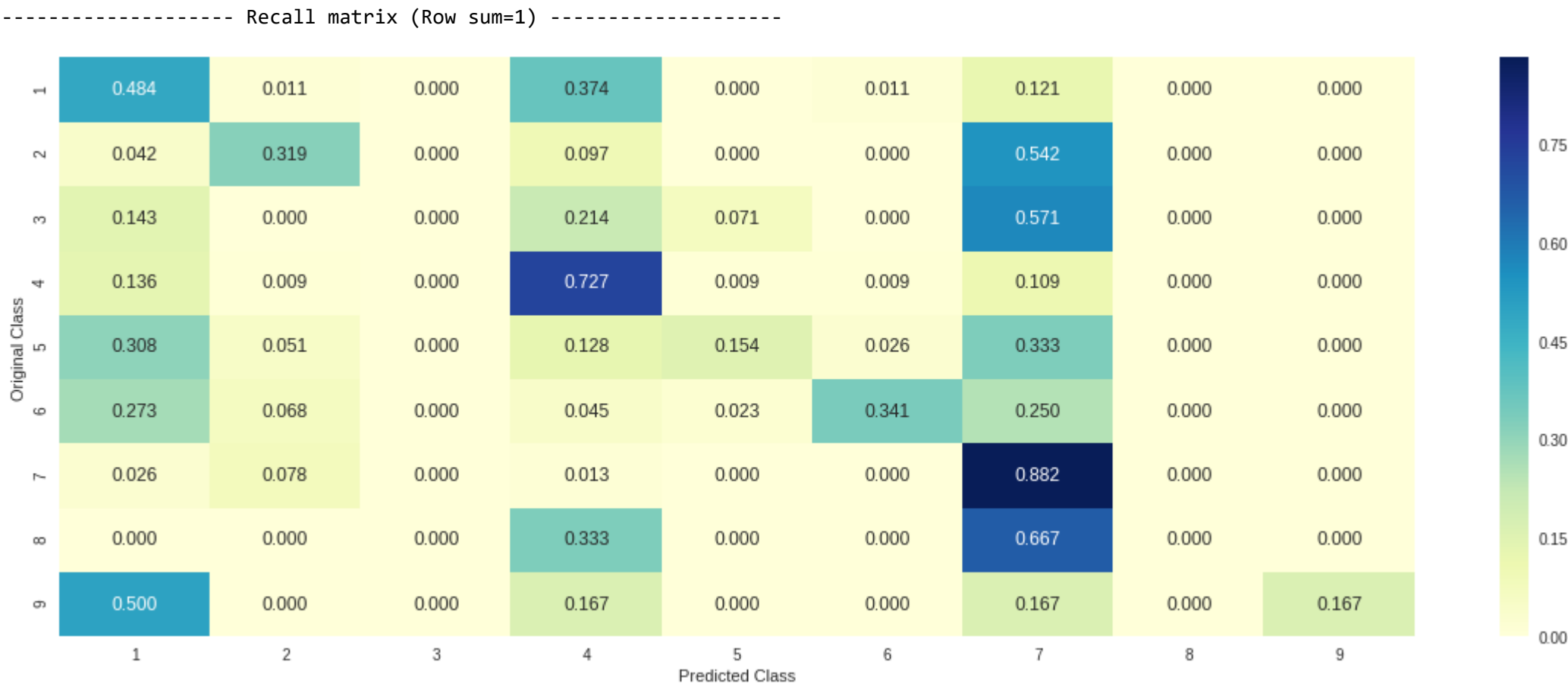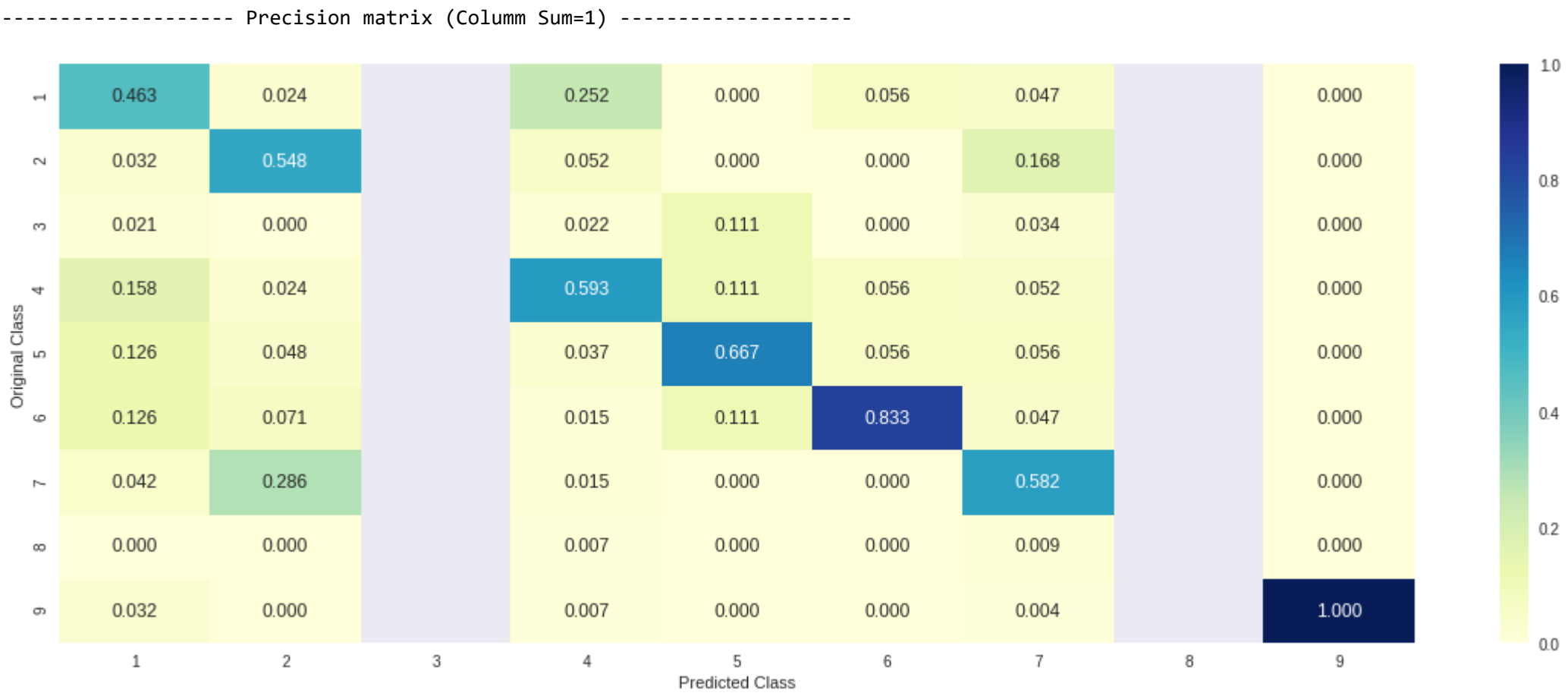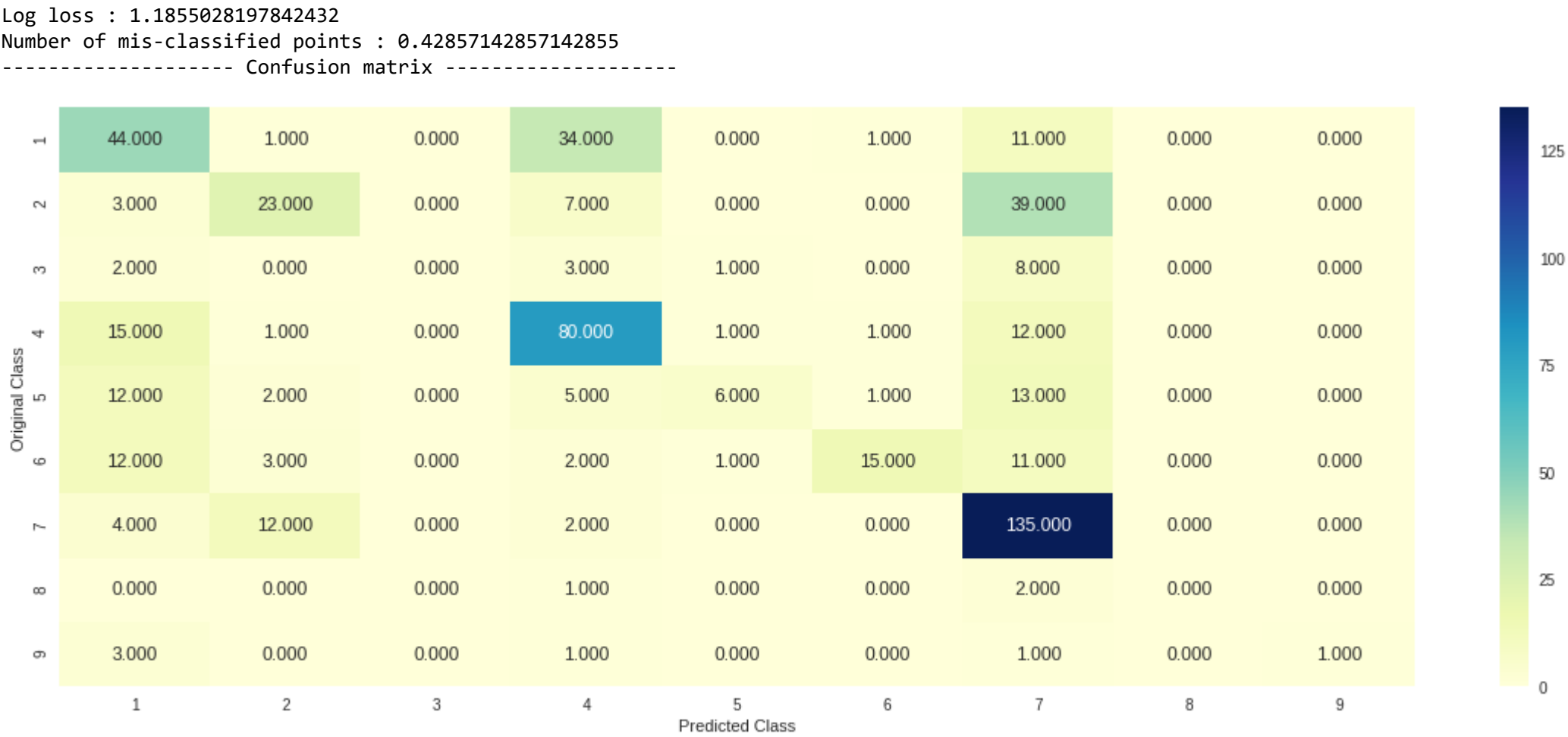
```
for n_estimators = 100 and max depth =  5
Log Loss : 1.2220094451266852
for n_estimators = 100 and max depth =  10
Log Loss : 1.233214354727003
for n_estimators = 200 and max depth =  5
Log Loss : 1.1958221839329501
for n_estimators = 200 and max depth =  10
Log Loss : 1.2152151149476638
for n_estimators = 500 and max depth =  5
Log Loss : 1.1908437380318881
for n_estimators = 500 and max depth =  10
Log Loss : 1.2091197041112527
for n_estimators = 1000 and max depth =  5
Log Loss : 1.1857287554700007
for n_estimators = 1000 and max depth =  10
Log Loss : 1.2053451548585152
for n_estimators = 2000 and max depth =  5
Log Loss : 1.1855028197842432
for n_estimators = 2000 and max depth =  10
Log Loss : 1.2022005734168828
For values of best estimator =  2000 The train log loss is: 0.8551956524228201
For values of best estimator =  2000 The cross validation log loss is: 1.1855028197842432
For values of best estimator =  2000 The test log loss is: 1.1145718647824898
```

In [0]:
```python
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha%2)], random_state=42, n_jobs=-1)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf)
```

Log loss : 1.1855028197842432
Number of mis-classified points : 0.42857142857142855
------------------- Confusion matrix -------------------



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



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



**Feature importance for correctly classified point**

In [0]:
```python
# test_point_index = 10
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha%2)], random_state=42, n_jobs=-1)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

Predicted Class : 4
Predicted Class Probabilities: [[0.0496 0.0097 0.0244 0.8283 0.0336 0.0262 0.0231 0.0029 0.0023]]
Actual Class : 4
-------------------------------------------------
59 Text feature [affecting] present in test data point [True]
Out of the top  100  features  1 are present in query point

**Feature importance for inorrectly classified point**

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0478 0.1158 0.0249 0.0478 0.0489 0.0493 0.6556 0.0048 0.0052]]
Actuall Class : 7
--------------------------------------------------
8 Text feature [alone] present in test data point [True]
41 Text feature [around] present in test data point [True]
51 Text feature [according] present in test data point [True]
59 Text feature [affecting] present in test data point [True]
76 Text feature [accessible] present in test data point [True]
99 Text feature [aliquot] present in test data point [True]
Out of the top  100  features  6 are present in query point
```

## Response Coding

```
In [0]:  alpha = [10,50,100,200,500,1000]
         max_depth = [2,3,5,10]
         cv_log_error_array = []
         for i in alpha:
             for j in max_depth:
                 print("for n_estimators =", i,"and max depth = ", j)
                 clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
                 clf.fit(train_x_responseCoding, train_y)
                 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
                 sig_clf.fit(train_x_responseCoding, train_y)
                 sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
                 cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
                 print("Log Loss :",log_loss(cv_y, sig_clf_probs))
         '''
         fig, ax = plt.subplots()
         features = np.dot(np.array(alpha)[:,None],np.array(max_depth)[None]).ravel()
         ax.plot(features, cv_log_error_array,c='g')
         for i, txt in enumerate(np.round(cv_log_error_array,3)):
             ax.annotate((alpha[int(i/4)],max_depth[int(i%4)],str(txt)), (features[i],cv_log_error_array[i]))
         plt.grid()
         plt.title("Cross Validation Error for each alpha")
         plt.xlabel("Alpha i's")
         plt.ylabel("Error measure")
         plt.show()
         '''

         best_alpha = np.argmin(cv_log_error_array)
         clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_depth=max_depth[int(best_alpha%4)], random_state=42, n_jobs=-1)
         clf.fit(train_x_responseCoding, train_y)
         sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
         sig_clf.fit(train_x_responseCoding, train_y)

         predict_y = sig_clf.predict_proba(train_x_responseCoding)
         print('For values of best alpha = ', alpha[int(best_alpha/4)], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(cv_x_responseCoding)
         print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
         predict_y = sig_clf.predict_proba(test_x_responseCoding)
         print('For values of best alpha = ', alpha[int(best_alpha/4)], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
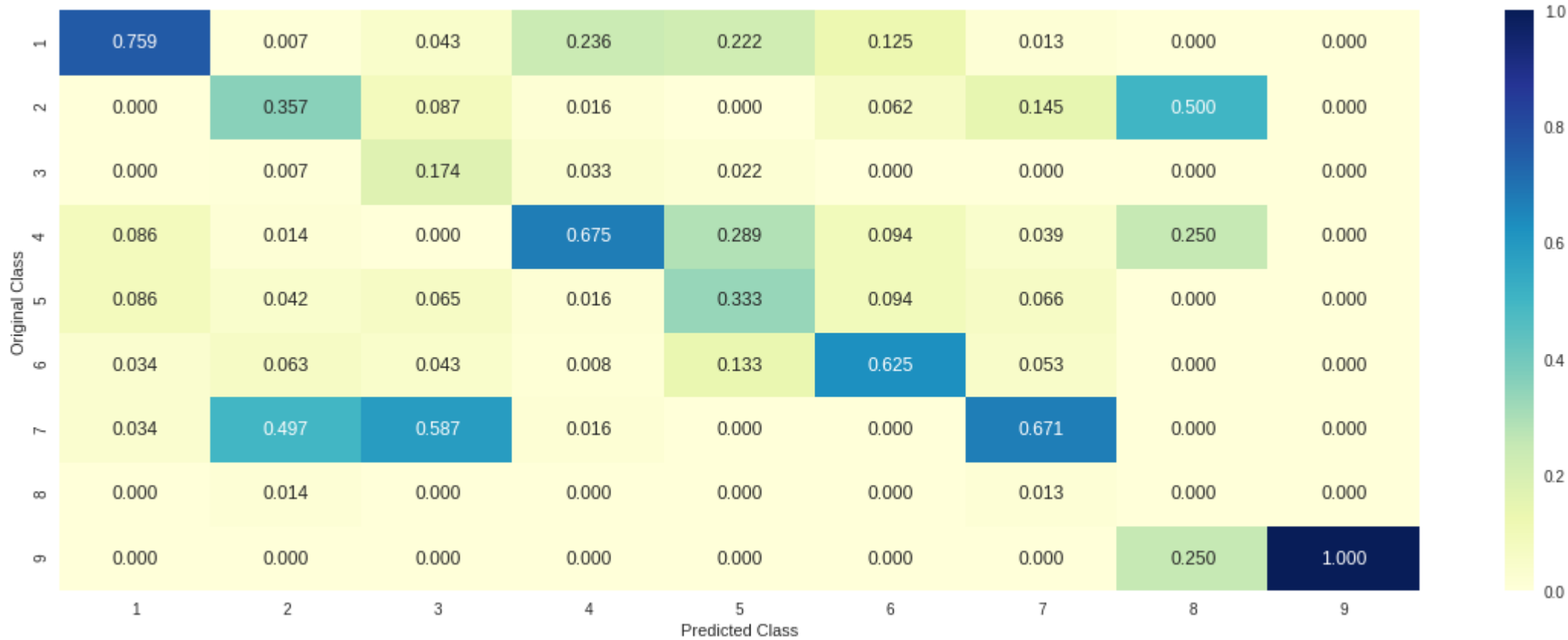
```
for n_estimators = 10 and max depth =  2
Log Loss : 2.252191750868235
for n_estimators = 10 and max depth =  3
Log Loss : 2.1659885231732456
for n_estimators = 10 and max depth =  5
Log Loss : 1.8097126023666052
for n_estimators = 10 and max depth =  10
Log Loss : 1.9843737201723362
for n_estimators = 50 and max depth =  2
Log Loss : 1.7333256486323743
for n_estimators = 50 and max depth =  3
Log Loss : 1.677006587662424
for n_estimators = 50 and max depth =  5
Log Loss : 1.6163467873791229
for n_estimators = 50 and max depth =  10
Log Loss : 1.5322612102969622
for n_estimators = 100 and max depth =  2
Log Loss : 1.5968578352368263
for n_estimators = 100 and max depth =  3
Log Loss : 1.4964023350682314
for n_estimators = 100 and max depth =  5
Log Loss : 1.3206570072592991
for n_estimators = 100 and max depth =  10
Log Loss : 1.450849327741403
for n_estimators = 200 and max depth =  2
Log Loss : 1.684096727622993
for n_estimators = 200 and max depth =  3
Log Loss : 1.4587635286995015
for n_estimators = 200 and max depth =  5
Log Loss : 1.332285915462042
for n_estimators = 200 and max depth =  10
Log Loss : 1.4375180789217783
for n_estimators = 500 and max depth =  2
Log Loss : 1.6052147328140849
for n_estimators = 500 and max depth =  3
Log Loss : 1.4642173256893227
for n_estimators = 500 and max depth =  5
Log Loss : 1.255197011758587
for n_estimators = 500 and max depth =  10
Log Loss : 1.4341334544390327
for n_estimators = 1000 and max depth =  2
Log Loss : 1.596243544653608
for n_estimators = 1000 and max depth =  3
Log Loss : 1.5057324822002502
for n_estimators = 1000 and max depth =  5
Log Loss : 1.2842127834975978
for n_estimators = 1000 and max depth =  10
Log Loss : 1.5032968957191686
For values of best alpha =  500 The train log loss is: 0.05857731418730803
For values of best alpha =  500 The cross validation log loss is: 1.255197011758587
For values of best alpha =  500 The test log loss is: 1.1849988487902927
```

```
In [0]: clf = RandomForestClassifier(max_depth=max_depth[int(best_alpha%4)], n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_features='auto',random_state=42)
        predict_and_plot_confusion_matrix(train_x_responseCoding, train_y,cv_x_responseCoding,cv_y, clf)
```

Log loss : 1.2551970117585867
Number of mis-classified points : 0.4793233082706767
-------------------- Confusion matrix --------------------



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



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



**Feature importance for correctly classified points**

```
In [0]: #clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_depth=max_depth[int(best_alpha%4)], random_state=42, n_jobs=-1)
        clf.fit(train_x_responseCoding, train_y)
        sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
        sig_clf.fit(train_x_responseCoding, train_y)


        test_point_index = 1
        no_feature = 27
        predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
        print("Predicted Class :", predicted_cls[0])
        print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_responseCoding[test_point_index].reshape(1,-1)),4))
        print("Actual Class :", test_y[test_point_index])
        indices = np.argsort(-clf.feature_importances_)
        print("-"*50)
        for i in indices:
            if i<9:
                print("Gene is important feature")
            elif i<18:
                print("Variation is important feature")
            else:
                print("Text is important feature")
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0294 0.0137 0.0837 0.7845 0.0147 0.0339 0.0067 0.0157 0.0178]]
Actual Class : 4
--------------------------------------------------
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Text is important feature
Text is important feature
Gene 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
Variation is important feature
Text is important feature
Variation is important feature
Gene is important feature
Text is important feature
Text is important feature
Gene is important feature
Text is important feature
Gene is important feature
Text is important feature
Text is important feature
Text is important feature
```

**Feature importance for incorrectly classified points**

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

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0198 0.2149 0.1666 0.0247 0.0249 0.2338 0.2778 0.0192 0.0182]]
Actual Class : 7
--------------------------------------------------
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Text is important feature
Text is important feature
Gene 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
Variation is important feature
Text is important feature
Variation is important feature
Gene is important feature
Text is important feature
Text is important feature
Gene is important feature
Text is important feature
Gene is important feature
Text is important feature
Text is important feature
Text is important feature
```

## Stacking Classifier

```
In [0]: from mlxtend.classifier import StackingClassifier

        clf1 = SGDClassifier(alpha=0.0001, penalty='l2', loss='log', class_weight='balanced', random_state=0)
        clf1.fit(train_x_onehotCoding, train_y)
        sig_clf1 = CalibratedClassifierCV(clf1, method="sigmoid")

        clf2 = SGDClassifier(alpha=0.001, penalty='l2', loss='hinge', class_weight='balanced', random_state=0)
        clf2.fit(train_x_onehotCoding, train_y)
        sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")


        clf3 = MultinomialNB(alpha=0.001)
        clf3.fit(train_x_onehotCoding, train_y)
        sig_clf3 = CalibratedClassifierCV(clf3, method="sigmoid")

        sig_clf1.fit(train_x_onehotCoding, train_y)
        print("Logistic Regression :  Log Loss: %0.2f" % (log_loss(cv_y, sig_clf1.predict_proba(cv_x_onehotCoding))))
        sig_clf2.fit(train_x_onehotCoding, train_y)
        print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf2.predict_proba(cv_x_onehotCoding))))
        sig_clf3.fit(train_x_onehotCoding, train_y)
        print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predict_proba(cv_x_onehotCoding))))
        print("-"*50)
        alpha = [0.0001,0.001,0.01,0.1,1,10]
        best_alpha = 999
        for i in alpha:
            lr = LogisticRegression(C=i)
            sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probas=True)
            sclf.fit(train_x_onehotCoding, train_y)
            print("Stacking Classifer : for the value of alpha: %f Log Loss: %0.3f" % (i, log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))))
            log_error =log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))
            if best_alpha > log_error:
                best_alpha = log_error
```

```
Logistic Regression :  Log Loss: 0.96
Support vector machines : Log Loss: 1.03
Naive Bayes : Log Loss: 1.20
--------------------------------------------------
Stacking Classifer : for the value of alpha: 0.000100 Log Loss: 2.172
Stacking Classifer : for the value of alpha: 0.001000 Log Loss: 1.986
Stacking Classifer : for the value of alpha: 0.010000 Log Loss: 1.376
Stacking Classifer : for the value of alpha: 0.100000 Log Loss: 1.073
Stacking Classifer : for the value of alpha: 1.000000 Log Loss: 1.268
Stacking Classifer : for the value of alpha: 10.000000 Log Loss: 1.653
```
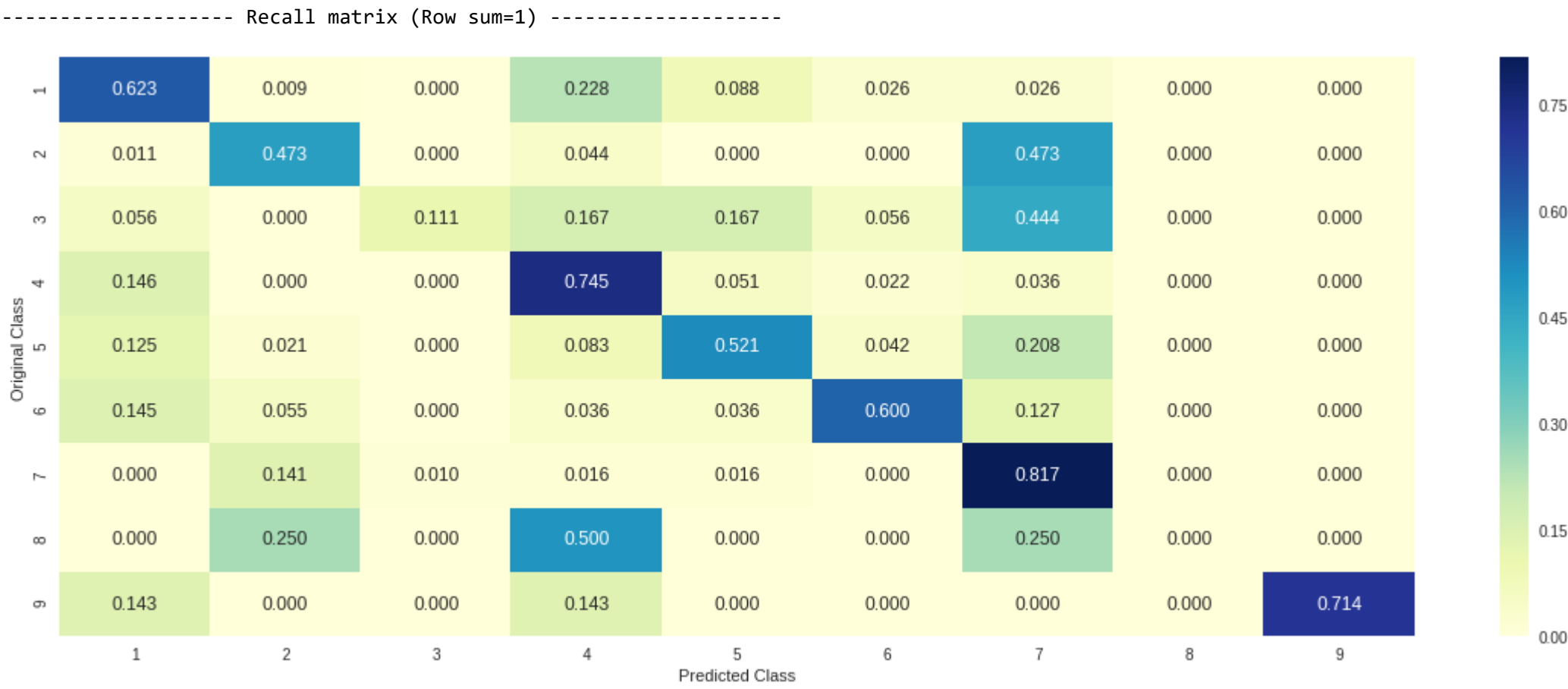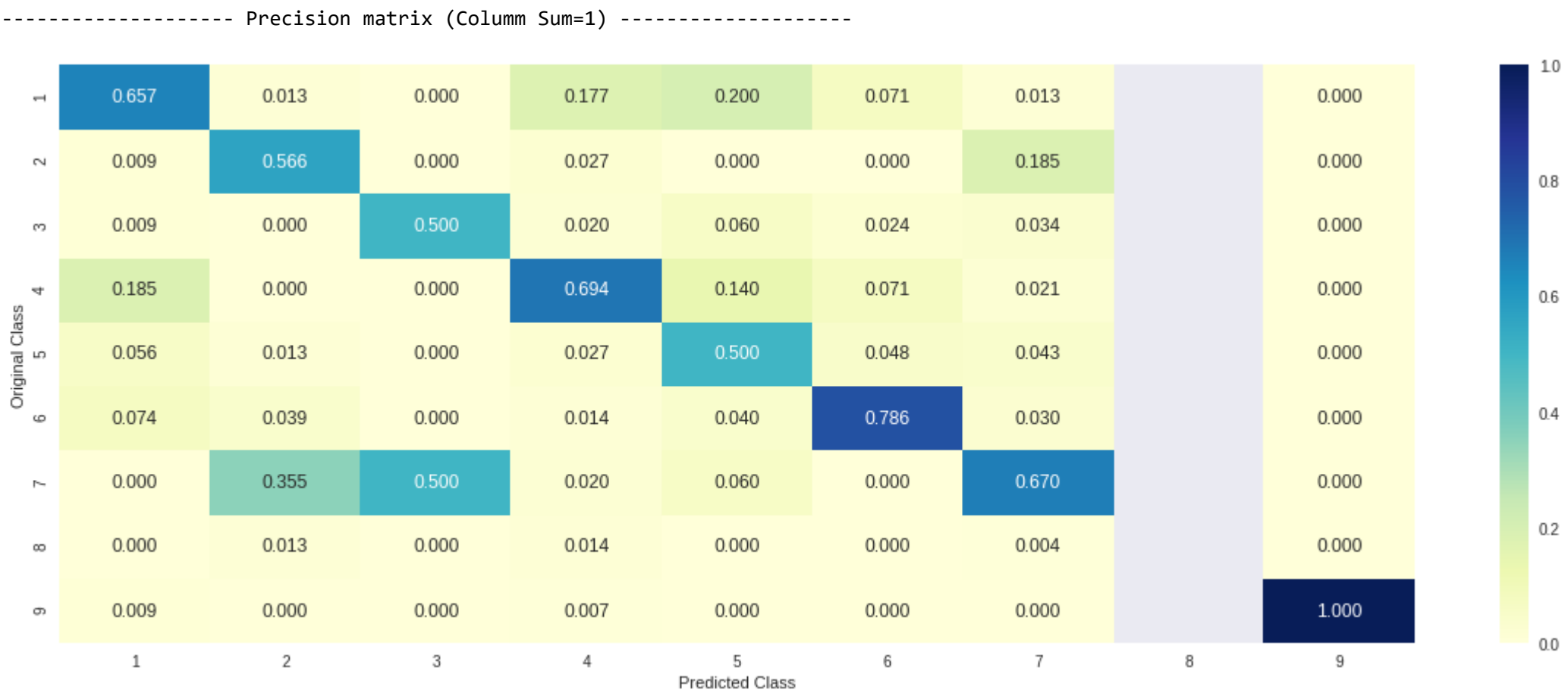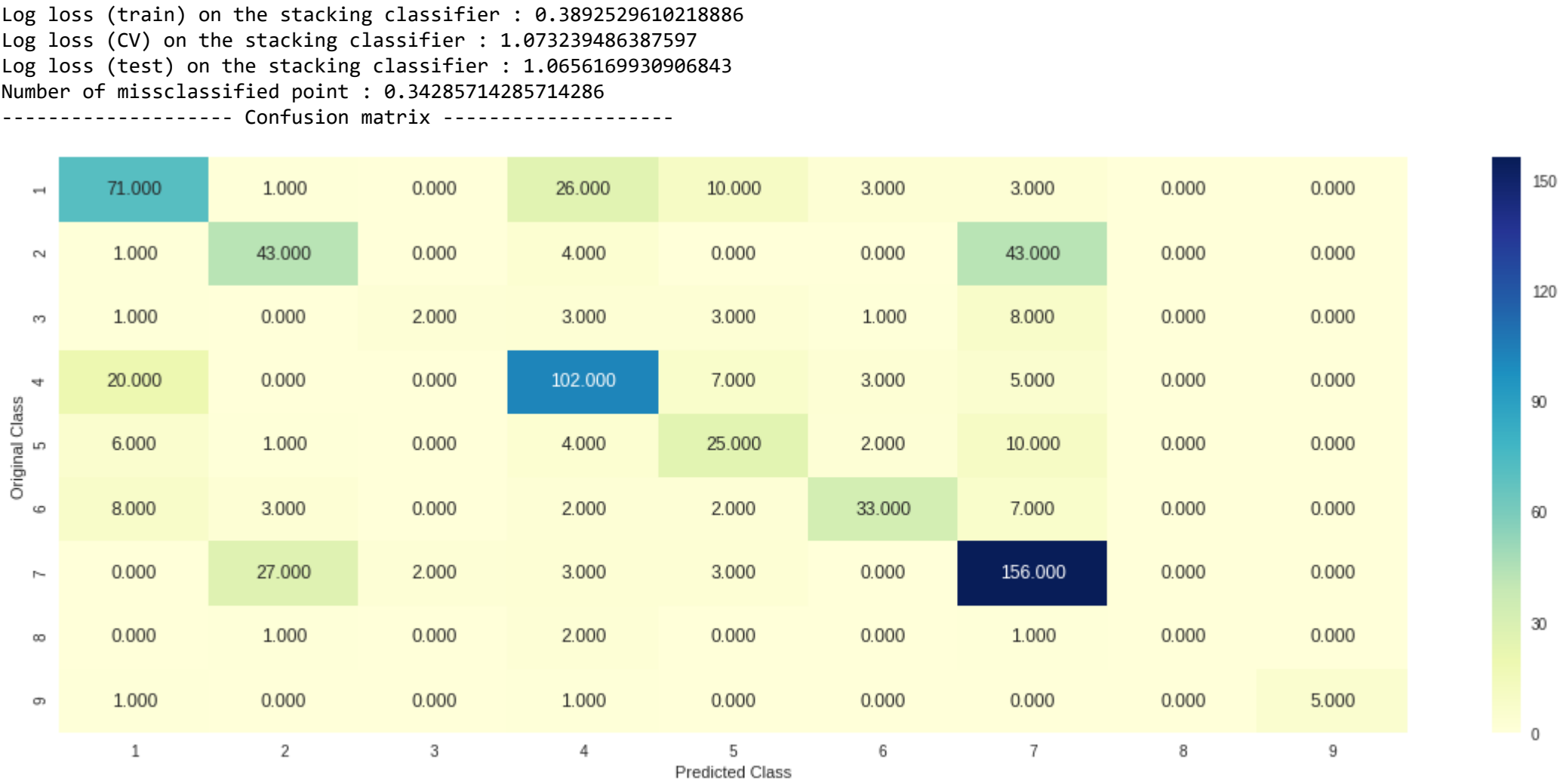
```
In [0]: lr = LogisticRegression(C=0.1)
        sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probas=True)
        sclf.fit(train_x_onehotCoding, train_y)

        log_error = log_loss(train_y, sclf.predict_proba(train_x_onehotCoding))
        print("Log loss (train) on the stacking classifier :",log_error)

        log_error = log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))
        print("Log loss (CV) on the stacking classifier :",log_error)

        log_error = log_loss(test_y, sclf.predict_proba(test_x_onehotCoding))
        print("Log loss (test) on the stacking classifier :",log_error)

        print("Number of missclassified point :", np.count_nonzero((sclf.predict(test_x_onehotCoding)- test_y))/test_y.shape[0])
        plot_confusion_matrix(test_y=test_y, predict_y=sclf.predict(test_x_onehotCoding))
```

```
Log loss (train) on the stacking classifier : 0.3892529610218886
Log loss (CV) on the stacking classifier : 1.073239486387597
Log loss (test) on the stacking classifier : 1.0656169930906843
Number of missclassified point : 0.34285714285714286
-------------------- Confusion matrix --------------------
```



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



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



## Maximum Voting classifier

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

```
Log loss (train) on the VotingClassifier : 0.5057273421750462
Log loss (CV) on the VotingClassifier : 1.0068529698160142
Log loss (test) on the VotingClassifier : 0.9773495713061825
Number of missclassified point : 0.34285714285714286
------------------- Confusion matrix -------------------
```



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



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



# Results

```python
from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["Model","Sampling","Hyper parameter","Train Log Loss","CV Log Loss", "Test Log Loss", "Percentage of misclassified points"]
x.add_row(["Naive Bayes","Without class balancing","alpha = 0.001","0.57","1.19","1.18","37.22%"])
x.add_row(["","","","","",""])
x.add_row(["K Nearest Neighbours","Without class balancing","k = 15","0.86","1.08","1.02","36.09%"])
x.add_row(["","","","","",""])
x.add_row(["Logistic Regression (TFIDF unigram)","SMOTE","alpha = 0.001","0.42","1.09","1.04","35.52%"])
x.add_row(["","With class balancing","alpha = 0.0001","0.42","0.96","0.92","31.76%"])
x.add_row(["","Without class balancing","alpha = 0.0001","0.42","0.96","0.92","31.20%"])
x.add_row(["","","","","",""])
x.add_row(["Logistic Regression (TFIDF unigrams and bigrams)","With class balancing","alpha = 0.0001","0.43","0.97","0.91","32.14%"])
x.add_row(["","Without class balancing","alpha = 0.0001","0.42","0.97","0.92","31.57%"])
x.add_row(["","","","","",""])
x.add_row(["Logistic Regression (BOW unigrams and bigrams)","With class balancing","alpha = 0.01","0.80","1.14","1.17","37.60%"])
x.add_row(["","Without class balancing","alpha = 0.01","0.78","1.14","1.18","36.28%"])
x.add_row(["","","","","",""])
x.add_row(["Support Vector Machines","With class balancing","alpha = 0.001","0.54","1.02","0.99","31.76%"])
x.add_row(["","","","","",""])
x.add_row(["Random Forests (Onehotencoding)","Without class balancing","n_estimators = 2000, max_depth = 5","0.85","1.18","1.11","42.85%"])
x.add_row(["Random Forests (Response coding)","Without class balancing","n_estimators = 500, max_depth = 5","0.06","1.25","1.18","47.93%"])
x.add_row(["","","","","",""])
x.add_row(["Stacking Classifier","","alpha = 0.1","0.39","1.07","1.06","34.28%"])
x.add_row(["","","","","",""])
x.add_row(["Max Voting Classifier","","alpha = 0.1","0.50","1.00","0.97","34.28%"])
print(x.get_string())
```

| Model | Sampling | Hyper parameter | Train Log Loss | CV Log Loss | Test Log Loss | Percentage of misclassified points |
|---|---|---|---|---|---|---|
| Naive Bayes | Without class balancing | alpha = 0.001 | 0.57 | 1.19 | 1.18 | 37.22% |
| | | | | | | |
| K Nearest Neighbours | Without class balancing | k = 15 | 0.86 | 1.08 | 1.02 | 36.09% |
| | | | | | | |
| Logistic Regression (TFIDF unigram) | SMOTE | alpha = 0.001 | 0.42 | 1.09 | 1.04 | 35.52% |
| | With class balancing | alpha = 0.0001 | 0.42 | 0.96 | 0.92 | 31.76% |
| | Without class balancing | alpha = 0.0001 | 0.42 | 0.96 | 0.92 | 31.20% |
| | | | | | | |
| Logistic Regression (TFIDF unigrams and bigrams) | With class balancing | alpha = 0.0001 | 0.43 | 0.97 | 0.91 | 32.14% |
| | Without class balancing | alpha = 0.0001 | 0.42 | 0.97 | 0.92 | 31.57% |
| | | | | | | |
| Logistic Regression (BOW unigrams and bigrams) | With class balancing | alpha = 0.01 | 0.80 | 1.14 | 1.17 | 37.60% |
| | Without class balancing | alpha = 0.01 | 0.78 | 1.14 | 1.18 | 36.28% |
| | | | | | | |
| Support Vector Machines | With class balancing | alpha = 0.001 | 0.54 | 1.02 | 0.99 | 31.76% |
| | | | | | | |
| Random Forests (Onehotencoding) | Without class balancing | n_estimators = 2000, max_depth = 5 | 0.85 | 1.18 | 1.11 | 42.85% |
| Random Forests (Response coding) | Without class balancing | n_estimators = 500, max_depth = 5 | 0.06 | 1.25 | 1.18 | 47.93% |
| | | | | | | |
| Stacking Classifier | | alpha = 0.1 | 0.39 | 1.07 | 1.06 | 34.28% |
| | | | | | | |
| Max Voting Classifier | | alpha = 0.1 | 0.50 | 1.00 | 0.97 | 34.28% |