

Personalized cancer diagnosis

Assignment 2-

Instead of using all the words in the dataset, use only the top 1000 words based of tf-idf values

1. Business Problem

1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

Context:

Source: <https://www.kaggle.com/c/msk-redefining-cancer-treatment/discussion/35336#198462>

Problem statement :

Classify the given genetic variations/mutations based on evidence from text-based clinical literature.

1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

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

1.3. Real-world/Business objectives and constraints.

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

2. Machine Learning Problem Formulation

2.1. Data

2.1.1. Data Overview

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

- 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>

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

2.3. Train, CV and Test Datasets

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

3. Exploratory Data Analysis

In [1]:

```
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
```

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

In [2]:

```
data = pd.read_csv('E:\\Machine Learning\\Assignments\\Personalized_Cancer_Diagnosis\\Data\\traini
ng_variants')
print('Number of data points : ', data.shape[0])
print('Number of features : ', data.shape[1])
print('Features : ', data.columns.values)
data.head()
```

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

Out[2]:

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

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

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

3.1.2. Reading Text Data

In [3]:

```
# note the separator in this file
data_text = pd.read_csv("E:\\Machine Learning\\Assignments\\Personalized_Cancer_Diagnosis\\Data\\training_text", sep="\\|\\|", engine="python", names=["ID", "TEXT"], skiprows=1)
print('Number of data points : ', data_text.shape[0])
print('Number of features : ', data_text.shape[1])
print('Features : ', data_text.columns.values)
data_text.head()
```

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

Out[3]:

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

3.1.3. Preprocessing of text

In [4]:

```
# loading stop words from nltk library
stop_words = set(stopwords.words('english'))

def nlp_preprocessing(total_text, index, column):
    if type(total_text) is not int:
        string = ""
```

```

# replace every special char with space
total_text = re.sub('[^a-zA-Z0-9\n]', ' ', total_text)
# replace multiple spaces with single space
total_text = re.sub('\s+', ' ', total_text)
# converting all the chars into lower-case.
total_text = total_text.lower()

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

data_text[column][index] = string

```

In [5]:

```

#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 : 178.901123113 seconds

```

In [6]:

```

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

```

Out[6]:

	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 [7]:

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

Out[7]:

	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 [8]:

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

In [9]:

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

Out [9]:

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

3.1.4. Test, Train and Cross Validation Split

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

In [10]:

```
from sklearn.model_selection import train_test_split
y_true = result['Class'].values
result.Gene = result.Gene.str.replace('\s+', '_')
result.Variation = result.Variation.str.replace('\s+', '_')

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

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

In [11]:

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

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

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

In [12]:

```
# 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 y_i 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), '%)')
```

```

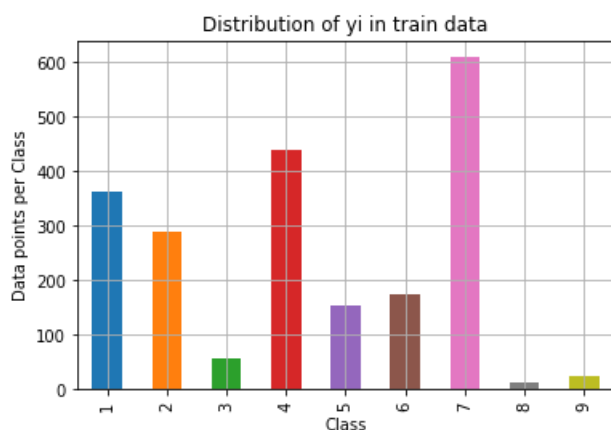
print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

# 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(
nd((test_class_distribution.values[i]/test_df.shape[0]*100), 3), '%)')

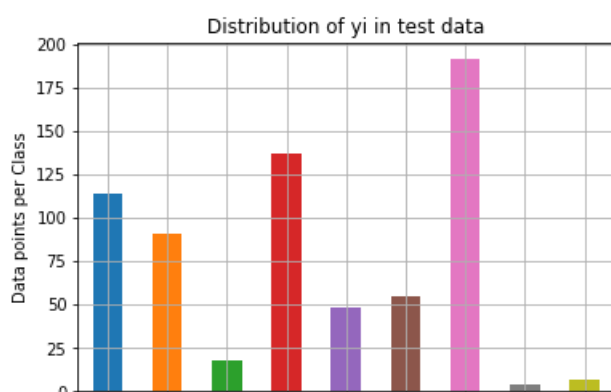
print('-'*80)
my_colors = 'rgbkymc'
cv_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in cross validation data')
plt.grid()
plt.show()

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

```



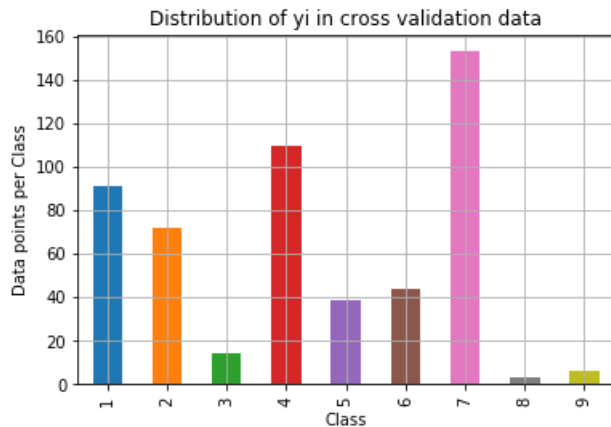
Number of data points in class 7 : 609 (28.672 %)
 Number of data points in class 4 : 439 (20.669 %)
 Number of data points in class 1 : 363 (17.09 %)
 Number of data points in class 2 : 289 (13.606 %)
 Number of data points in class 6 : 176 (8.286 %)
 Number of data points in class 5 : 155 (7.298 %)
 Number of data points in class 3 : 57 (2.684 %)
 Number of data points in class 9 : 24 (1.13 %)
 Number of data points in class 8 : 12 (0.565 %)



1. 2. 3. 4. 5. 6. 7. 8. 9.

Class

Number of data points in class 7 : 191 (28.722 %)
 Number of data points in class 4 : 137 (20.602 %)
 Number of data points in class 1 : 114 (17.143 %)
 Number of data points in class 2 : 91 (13.684 %)
 Number of data points in class 6 : 55 (8.271 %)
 Number of data points in class 5 : 48 (7.218 %)
 Number of data points in class 3 : 18 (2.707 %)
 Number of data points in class 9 : 7 (1.053 %)
 Number of data points in class 8 : 4 (0.602 %)



Number of data points in class 7 : 153 (28.759 %)
 Number of data points in class 4 : 110 (20.677 %)
 Number of data points in class 1 : 91 (17.105 %)
 Number of data points in class 2 : 72 (13.534 %)
 Number of data points in class 6 : 44 (8.271 %)
 Number of data points in class 5 : 39 (7.331 %)
 Number of data points in class 3 : 14 (2.632 %)
 Number of data points in class 9 : 6 (1.128 %)
 Number of data points in class 8 : 3 (0.564 %)

3.2 Prediction using a 'Random' Model

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

In [13]:

```
# This function plots the confusion matrices given  $y_i$ ,  $y_{i\_hat}$ .
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j

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

    # C = [[1, 2],
    #       [3, 4]]
    # C.T = [[1, 3],
    #         [2, 4]]
    # C.sum(axis = 1)  axis=0 corresponds to columns and axis=1 corresponds to rows in two
    # dimensional array
    # C.sum(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 corresponds to columns and axis=1 corresponds to rows in two
dimensional array
# C.sum(axis=0) = [[4, 6]]
# (C/C.sum(axis=0)) = [[1/4, 2/6],
# [3/4, 4/6]]

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

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

# representing B in heatmap format
print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
plt.figure(figsize=(20,7))
sns.heatmap(A, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()

```

In [14]:

```

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

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

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

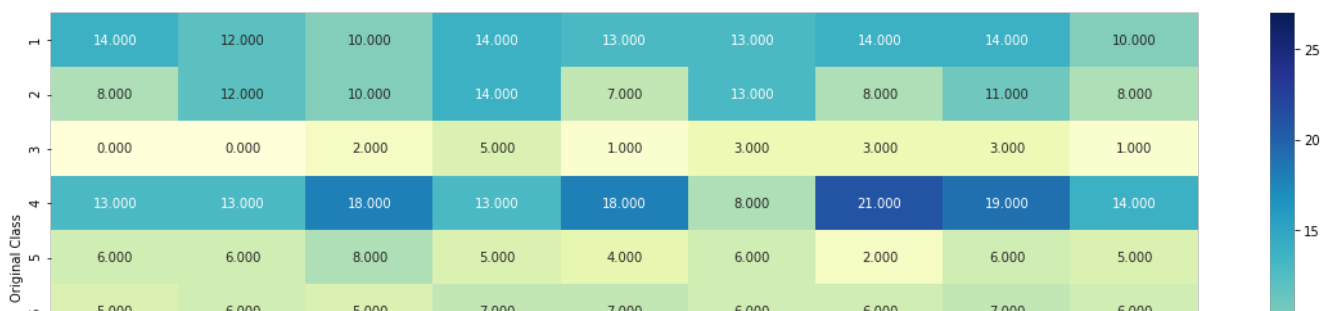
predicted_y = np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

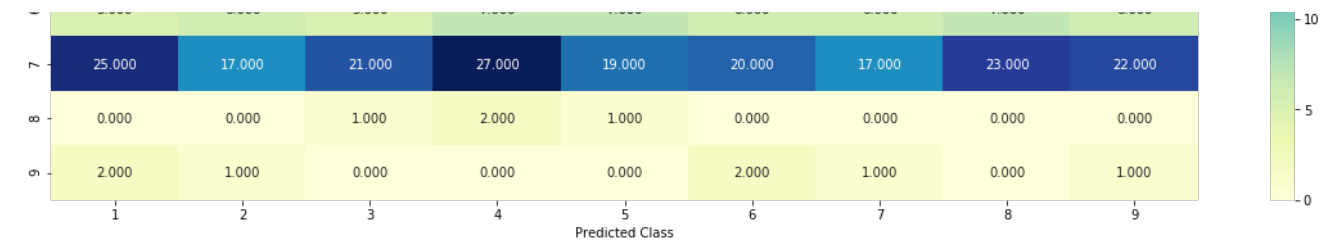
```

Log loss on Cross Validation Data using Random Model 2.443224340690172

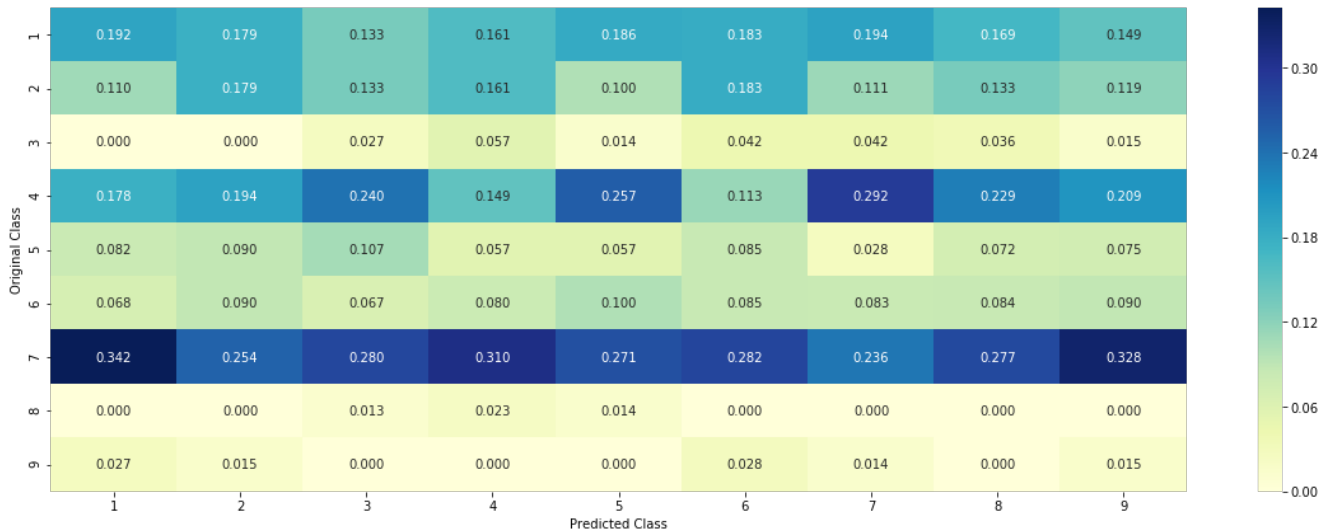
Log loss on Test Data using Random Model 2.4835341367475485

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

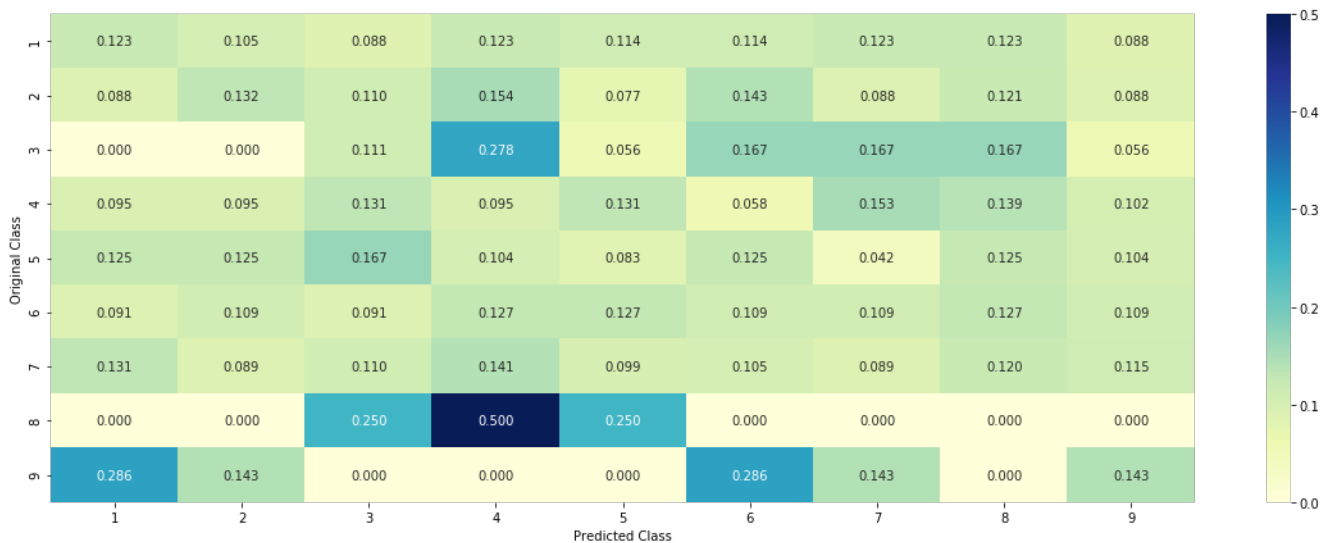




----- Precision matrix (Column Sum=1) -----



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



3.3 Univariate Analysis

In [15]:

```
# 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 occurred in class1 + 10*alpha / number of time it occurred in total data+90*alpha)
# gv_dict is like a look up table, for every gene it store a (1*9) representation of it
# for a value of feature in df:
# if it is in train data:
```

```

# we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
# if it is not there is train:
# we add [1/9, 1/9, 1/9, 1/9,1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
# return 'gv_fea'
# -----

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

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

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

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

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    #      {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.0681818181818177,
    0.136363636363635, 0.25, 0.193181818181818, 0.037878787878788, 0.037878787878788,
    0.037878787878788],
    #      'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366,
    0.27040816326530615, 0.061224489795918366, 0.066326530612244902, 0.051020408163265307, 0.051020408
    163265307, 0.056122448979591837],
    #      'EGFR': [0.0568181818181816, 0.21590909090909091, 0.0625, 0.0681818181818177,

```

```

0.0681818181818177, 0.0625, 0.34659090909090912, 0.0625, 0.0568181818181816],
# 'BRCA2': [0.1333333333333333, 0.0606060606060606, 0.0606060606060606,
0.0787878787878782, 0.1393939393939394, 0.3454545454545454, 0.0606060606060606,
0.0606060606060606, 0.0606060606060606],
# 'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917,
0.46540880503144655, 0.075471698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081
761006289, 0.062893081761006289],
# 'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295,
0.072847682119205295, 0.066225165562913912, 0.066225165562913912, 0.27152317880794702,
0.066225165562913912, 0.066225165562913912],
# 'BRAF': [0.06666666666666666, 0.17999999999999999, 0.07333333333333334,
0.07333333333333334, 0.09333333333333338, 0.08000000000000002, 0.29999999999999999,
0.06666666666666666, 0.06666666666666666],
# ...
# }
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 da
ta
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] 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])
# gv_fea.append([-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

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

3.2.1 Univariate Analysis on Gene Feature

Q1. Gene, What type of feature it is ?

Ans. Gene is a categorical variable

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

In [16]:

```

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

```

```

Number of Unique Genes : 236
BRCA1      158
TP53       115
EGFR       98
BRCA2      82
PTEN       71
KIT        66
BRAF       58
ERBB2      46
ALK        42
PDGFRA     40
Name: Gene, dtype: int64

```

In [17]:

```

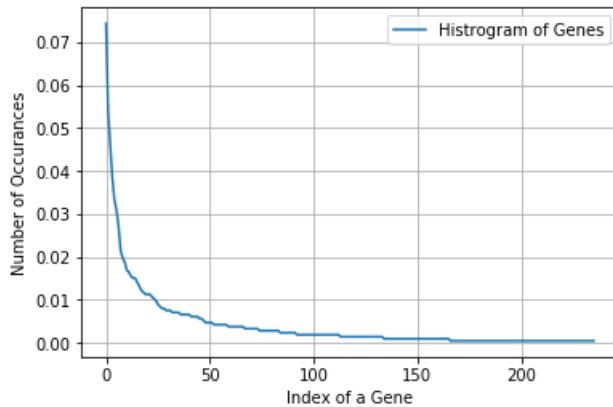
print("Ans: There are", unique_genes.shape[0] , "different categories of genes in the train data, and they are distributed as follows",)

```

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

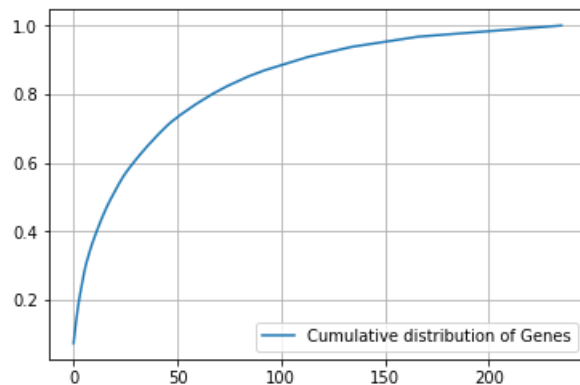
In [18]:

```
s = sum(unique_genes.values);
h = unique_genes.values/s;
plt.plot(h, label="Histogram of Genes")
plt.xlabel('Index of a Gene')
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```



In [19]:

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



Q3. How to featurize this Gene feature ?

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

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

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 [20]:

```
#response-coding of the Gene feature
# alpha is used for laplace smoothing
alpha = 1
# train gene feature
from sklearn.feature_extraction import CountVec... (code continues)
```

```
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 [21]:

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

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

In [22]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer()
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 [23]:

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

Out[23]:

```
894      PDGFRA
719      ERBB2
1745     MSH2
1326     MLH1
1992     MAP2K1
Name: Gene, dtype: object
```

In [24]:

```
gene_vectorizer.get_feature_names()
```

Out[24]:

```
['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ar',
 'araf',
 'arid1b',
 'arid2',
 'arid5b',
 'asx11',
 'asx12',
 'atm',
 'atr',
 'atrx',
 'aurkb',
 'axin1',
 'b2m',
 'bap1',
 'bard1',
 'bcl10',
 'bcl2',
 'bcl2l11',
 'bcor',
 'braf',
 'brca1',
 'brca2',
 'brd4',
 ...]
```

'brip1',
'btk',
'card11',
'carm1',
'casp8',
'cbl',
'ccnd1',
'ccnd2',
'ccnd3',
'ccne1',
'cdh1',
'cdk12',
'cdk4',
'cdk6',
'cdk8',
'cdkn1a',
'cdkn1b',
'cdkn2a',
'cdkn2b',
'cdkn2c',
'chek2',
'cic',
'crebbp',
'ctcf',
'ctla4',
'ctnnb1',
'ddr2',
'dicer1',
'dnmt3a',
'dnmt3b',
'egfr',
'eiflax',
'elf3',
'ep300',
'epas1',
'epcam',
'erbb2',
'erbb3',
'erbb4',
'ercc2',
'ercc3',
'ercc4',
'erg',
'errfi1',
'esr1',
'etv1',
'etv6',
'ewsr1',
'ezh2',
'fam58a',
'fanca',
'fancc',
'fat1',
'fbxw7',
'fgf19',
'fgf3',
'fgf4',
'fgfr1',
'fgfr2',
'fgfr3',
'fgfr4',
'flt3',
'foxa1',
'foxl2',
'foxo1',
'foxp1',
'fubp1',
'gata3',
'gna11',
'gnaq',
'gnas',
'h3f3a',
'hnfla',
'hras',
'idh1',
'idh2',
'igflr',

'ikbke',
'il7r',
'jak1',
'jak2',
'kdm5a',
'kdm5c',
'kdm6a',
'kdr',
'keap1',
'kit',
'klf4',
'kmt2a',
'kmt2b',
'kmt2c',
'kmt2d',
'knstrn',
'kras',
'lats1',
'map2k1',
'map2k2',
'map2k4',
'map3k1',
'mapk1',
'mdm4',
'med12',
'mef2b',
'men1',
'met',
'mga',
'mlh1',
'mpl',
'msh2',
'msh6',
'mtor',
'myc',
'mycn',
'myd88',
'myod1',
'ncor1',
'nfl',
'nf2',
'nfe2l2',
'nfkb1a',
'nkx2',
'notch1',
'notch2',
'nras',
'nsd1',
'ntrk1',
'ntrk2',
'ntrk3',
'nup93',
'pak1',
'pax8',
'pbrm1',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pik3r3',
'pim1',
'pms1',
'pms2',
'pole',
'ppp2r1a',
'ppp6c',
'prdm1',
'ptch1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rab35',
'rac1',


```
'rad21',
'rad50',
'rad51b',
'rad51c',
'rad51d',
'rad541',
'raf1',
'rb1',
'rbm10',
'ret',
'rheb',
'rhoa',
'rictor',
'rit1',
'ros1',
'rras2',
'runx1',
'rxra',
'rybp',
'sdhb',
'setd2',
'sf3b1',
'shq1',
'smad2',
'smad3',
'smad4',
'smarca4',
'smarcb1',
'smo',
'sos1',
'sox9',
'spop',
'src',
'stat3',
'stk11',
'tert',
'tet2',
'tgfbr1',
'tgfbr2',
'tmprss2',
'tp53',
'tp53bp1',
'tsc1',
'tsc2',
'u2af1',
'vhl',
'whsc1',
'whsc1l1',
'xpo1',
'xrcc2',
'yap1']
```

In [25]:

```
print("train_gene_feature_onehotCoding is converted feature using tfidf encoding method. The shape
of gene feature:", train_gene_feature_onehotCoding.shape)
```

```
train_gene_feature_onehotCoding is converted feature using tfidf encoding method. The shape of
gene feature: (2124, 236)
```

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 [26]:

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

# read more about SGDClassifier() at http://scikit-
# learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
```

```

# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

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

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

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

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

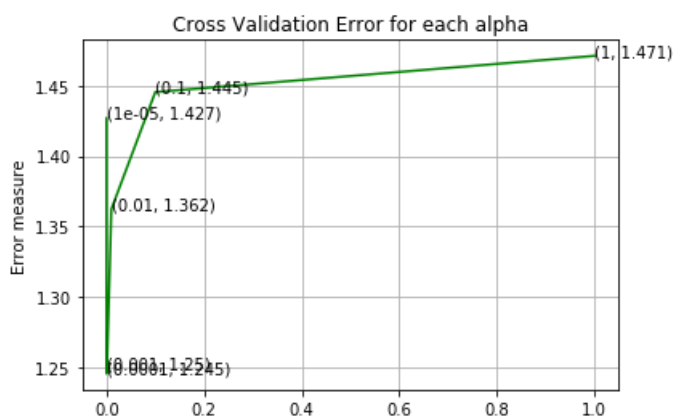
predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
print('For values of best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

```

```

For values of alpha = 1e-05 The log loss is: 1.426807306948854
For values of alpha = 0.0001 The log loss is: 1.2454475140592103
For values of alpha = 0.001 The log loss is: 1.2496019712668545
For values of alpha = 0.01 The log loss is: 1.3622292293706546
For values of alpha = 0.1 The log loss is: 1.4453946247918192
For values of alpha = 1 The log loss is: 1.4711378904748011

```



For values of best alpha = 0.0001 The train log loss is: 1.0539602595040236
 For values of best alpha = 0.0001 The cross validation log loss is: 1.2454475140592103
 For values of best alpha = 0.0001 The test log loss is: 1.2006828905934004

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 [27]:

```
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 236 genes in train dataset?
 Ans

1. In test data 643 out of 665 : 96.69172932330827
2. In cross validation data 517 out of 532 : 97.18045112781954

3.2.2 Univariate Analysis on Variation Feature

Q7. Variation, What type of feature is it ?

Ans. Variation is a categorical variable

Q8. How many categories are there?

In [28]:

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

```
Number of Unique Variations : 1922
Truncating_Mutations    55
Deletion                 53
Amplification            49
Fusions                  24
Overexpression           4
G12V                     3
Q61H                     3
V321M                    2
I31M                     2
E17K                     2
Name: Variation, dtype: int64
```

In [29]:

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

Ans: There are 1922 different categories of variations in the train data, and they are distributed as follows

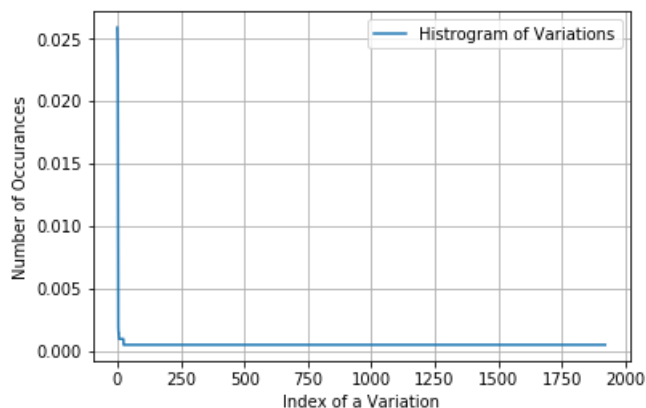
In [30]:

```
s = sum(unique_variations.values);
```

```

n = unique_variations.values/s;
plt.plot(h, label="Histogram of Variations")
plt.xlabel('Index of a Variation')
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()

```



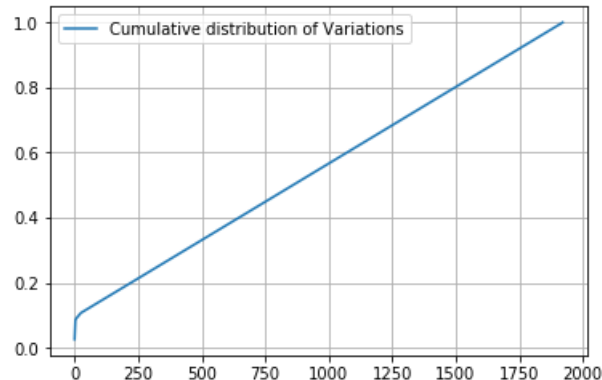
In [31]:

```

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

```

```
[0.02589454 0.05084746 0.07391714 ... 0.99905838 0.99952919 1.          ]
```



Q9. How to featurize this Variation feature ?

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

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

1. One hot Encoding
2. Response coding

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

In [32]:

```

# 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 [33]:

```
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 [34]:

```
# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer()
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 [35]:

```
print("train_variation_feature_onehotEncoded is converted feature using the Tfidf encoding method. The shape of Variation feature:", train_variation_feature_onehotCoding.shape)
```

train_variation_feature_onehotEncoded is converted feature using the Tfidf encoding method. The shape of Variation feature: (2124, 1953)

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

Let's build a model just like the earlier!

In [36]:

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

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

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

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

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

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

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

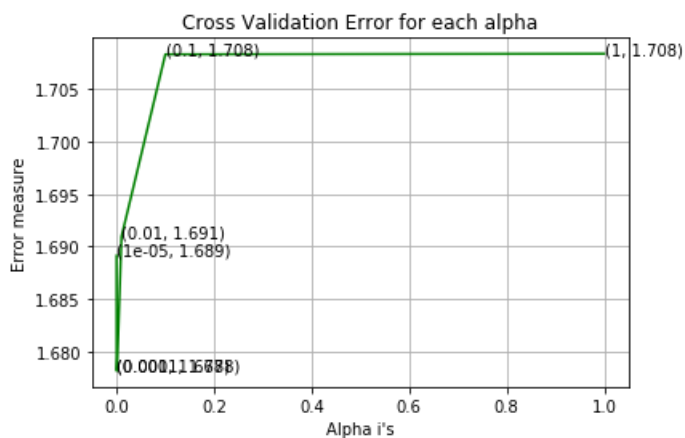
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], cv_log_error_array[i]))
plt.grid()
```

```
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

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

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

For values of alpha = 1e-05 The log loss is: 1.689123678946004
For values of alpha = 0.0001 The log loss is: 1.678151232848739
For values of alpha = 0.001 The log loss is: 1.6782257699906953
For values of alpha = 0.01 The log loss is: 1.6909346334903939
For values of alpha = 0.1 The log loss is: 1.7083199905684636
For values of alpha = 1 The log loss is: 1.7083749832875543



For values of best alpha = 0.0001 The train log loss is: 0.771899875633995
For values of best alpha = 0.0001 The cross validation log loss is: 1.678151232848739
For values of best alpha = 0.0001 The test log loss is: 1.7314055618658766

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 [37]:

```
print("Q12. How many data points are covered by total ", unique_variations.shape[0], " genes in te
st 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.s
hape[0])*100)
```

Q12. How many data points are covered by total 1922 genes in test and cross validation data sets?

Ans

1. In test data 53 out of 665 : 7.969924812030076
2. In cross validation data 59 out of 532 : 11.090225563909774

3.2.3 Univariate Analysis on Text Feature

1. How many unique words are present in train data?
2. How are word frequencies distributed?
3. How to featurize text field?
4. Is the text feature useful in predicting y_i ?
5. Is the text feature stable across train, test and CV datasets?

In [38]:

```
# cls_text is a data frame
# for every row in data 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 [39]:

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

In [53]:

```
#from sklearn.feature_extraction.text import TfidfVectorizer
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer(min_df=3,max_features=1000)
train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])
# getting all the feature names (words)
train_text_features= text_vectorizer.get_feature_names()

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

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

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

Total number of unique words in train data : 1000

In [54]:

```
print(train_text_feature_onehotCoding.shape)
```

(2124, 1000)

In []:

In [55]:

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

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

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

In [56]:

```
#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 [57]:

```
# 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 [58]:

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

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

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

In [59]:

```
#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 [60]:

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

```
Counter({255.4086836711719: 1, 180.3444202948155: 1, 134.09338862549183: 1, 131.14652965736693: 1,
```


127.73290602016066: 1, 117.05851225564032: 1, 116.8110465657733: 1, 116.20867543623926: 1, 109.12237725987195: 1, 105.15240856561334: 1, 104.55335474719448: 1, 93.28244234897271: 1, 93.01531862116894: 1, 88.84937468705874: 1, 83.83625223737529: 1, 81.50698951733614: 1, 78.70911645905615: 1, 78.14570097806798: 1, 77.41092906698577: 1, 77.06920445537187: 1, 73.72810300226574: 1, 72.81489249422077: 1, 72.31189473279127: 1, 70.77400309165854: 1, 69.39071375803019: 1, 68.89129506421139: 1, 67.82388014484098: 1, 64.87062579045296: 1, 64.34756440221705: 1, 64.09719196067614: 1, 64.0726770867231: 1, 63.07548702183802: 1, 61.408336206671194: 1, 60.72274438305604: 1, 59.03266499638227: 1, 57.16659739241288: 1, 56.8710975 90022494: 1, 56.66369568700459: 1, 52.018487689129984: 1, 50.11767851273153: 1, 49.82914559814799: 1, 49.11526161282805: 1, 48.42469931506257: 1, 47.9850263428371: 1, 47.61809991370876: 1, 46.4452752278991: 1, 46.39813049881669: 1, 45.25264436651568: 1, 45.19646289338978: 1, 44.16947403787889: 1, 44.167975644627475: 1, 43.93988542880423: 1, 43.31545366876956: 1, 42.83715765827708: 1, 42.38419011461101: 1, 42.314322651079415: 1, 42.300664390946615: 1, 42.12365790413625: 1, 42.07289273126195: 1, 41.923966582085974: 1, 41.638093523317444: 1, 41.523994051737354: 1, 41.51485256015422: 1, 41.170306470756806: 1, 41.154224830674764: 1, 41.0956318640042: 1, 41.065653152835175: 1, 40.70836129836931: 1, 40.63481359060724: 1, 39.83604552319816: 1, 39.42542849285207: 1, 39.15273690512158: 1, 38.466103004056436: 1, 38.05748927569554: 1, 37.61117150657249: 1, 37.43994020235219: 1, 36.495359834589145: 1, 36.44470501099222: 1, 36.123243870671: 1, 35.50315566436088: 1, 35.469257444352166: 1, 35.28229614627149: 1, 35.13261540398904: 1, 34.620921640879665: 1, 34.53155257434234: 1, 34.4528596681028: 1, 34.43354206509663: 1, 34.26304204491943: 1, 34.23604854401906: 1, 34.115751645743956: 1, 34.053418977455394: 1, 33.42156528394029: 1, 33.0693134679787: 1, 32.42531116877464: 1, 32.34748134285403: 1, 32.334672608291065: 1, 32.33037321320598: 1, 32.212133704142545: 1, 32.18812194024285: 1, 31.991836828571873: 1, 31.96530709219318: 1, 31.928516781343806: 1, 31.871355716674298: 1, 31.8693594044294: 1, 31.656465790495172: 1, 31.58762816384161: 1, 31.541595077240054: 1, 31.455450540157628: 1, 31.40548902590355: 1, 31.371382312456763: 1, 30.919791170783974: 1, 30.797153356663: 1, 30.71798911100032: 1, 30.675109638767893: 1, 30.617911179178964: 1, 30.41425204526495: 1, 30.33954497584197: 1, 30.117131929698793: 1, 29.94787888863044: 1, 29.702076525790584: 1, 29.56634783162448: 1, 29.50427799410116: 1, 29.14003908993198: 1, 29.042655510418953: 1, 28.698735469563637: 1, 28.689412018374444: 1, 28.64035480574123: 1, 28.52341257002209: 1, 28.351311331800403: 1, 28.31320455051252: 1, 28.30362205511252: 1, 28.29014896979023: 1, 28.196315214735698: 1, 28.000268224052924: 1, 27.696584184125772: 1, 27.54124837286089: 1, 27.50152183956183: 1, 27.272726558554503: 1, 27.16702142115833: 1, 27.0150703971168: 1, 26.92603814745057: 1, 26.90516342983227: 1, 26.80608474727773: 1, 26.794228118942396: 1, 26.579821603553363: 1, 26.49890056117726: 1, 26.420328121082438: 1, 26.276365374438875: 1, 26.24278895987232: 1, 26.161032726913252: 1, 26.117317148266142: 1, 25.809291818064672: 1, 25.42560934924831: 1, 25.29364560041617: 1, 25.199103271378128: 1, 25.16694997344531: 1, 25.129335242248104: 1, 24.833493856897686: 1, 24.822031312907434: 1, 24.721005048106154: 1, 24.696665492861886: 1, 24.66046248188366: 1, 24.608597758399508: 1, 24.58982621429249: 1, 24.513328419983523: 1, 24.4901109203332: 1, 24.462179421013214: 1, 24.44232760812821: 1, 24.339371186238036: 1, 24.264154774574923: 1, 24.222129642465088: 1, 24.218163684094883: 1, 24.11324234487665: 1, 23.95975745966559: 1, 23.830710024032886: 1, 23.818373864939925: 1, 23.738213699491116: 1, 23.68920770586673: 1, 23.620873337788115: 1, 23.60531130890419: 1, 23.41471696372977: 1, 23.252650421713064: 1, 23.13068382174456: 1, 23.06314027639983: 1, 23.06045019979732: 1, 23.0368726 2686757: 1, 23.004676685400543: 1, 22.94532722345419: 1, 22.883445595630597: 1, 22.80544496854292: 1, 22.75216451890397: 1, 22.716760172407948: 1, 22.703704713603585: 1, 22.702295234436452: 1, 22.668756322084118: 1, 22.624636206451036: 1, 22.623569924781037: 1, 22.54470824874731: 1, 22.537727470595886: 1, 22.522814114548613: 1, 22.51857103801733: 1, 22.50488826153655: 1, 22.450154482233497: 1, 22.417413199290422: 1, 22.33829908271732: 1, 22.305085791400693: 1, 22.245569594709895: 1, 22.0948882704275: 1, 22.085833676751946: 1, 22.083257577733523: 1, 22.08107431384057: 1, 22.071660668204654: 1, 22.00669191922763: 1, 21.934083627869303: 1, 21.857207146141015: 1, 21.83014517294909: 1, 21.798186652306192: 1, 21.739104003980398: 1, 21.692005421336457: 1, 21.675764596532073: 1, 21.54799609852769: 1, 21.524837050991884: 1, 21.41408132279104: 1, 21.41170004020658: 1, 21.19088460257289: 1, 21.15999991391758: 1, 21.12028644961535: 1, 21.093952479555217: 1, 20.946533296084763: 1, 20.934611502276606: 1, 20.91724650188194: 1, 20.876412775482486: 1, 20.8698484276622: 1, 20.822081042764026: 1, 20.587312746331815: 1, 20.459661196211542: 1, 20.442092004010842: 1, 20.33289638633411: 1, 20.3252318357821: 1, 20.256555224970135: 1, 20.125268519634737: 1, 20.111145019028566: 1, 20.008106767473823: 1, 19.992920740094565: 1, 19.986512429400257: 1, 19.975647616752276: 1, 19.961080334839515: 1, 19.924561918387944: 1, 19.820985756864932: 1, 19.820659529970886: 1, 19.78775470849255: 1, 19.751393223546973: 1, 19.74056074842421: 1, 19.716890418436805: 1, 19.700853546366258: 1, 19.66911934287414: 1, 19.63058204384549: 1, 19.50361857349681: 1, 19.4673627 50489023: 1, 19.397200271440358: 1, 19.366794219530917: 1, 19.348665702518723: 1, 19.29944668781282: 1, 19.278280968049994: 1, 19.244032924957704: 1, 19.237020933596543: 1, 19.205045309339024: 1, 19.204468006876525: 1, 19.201130189975697: 1, 19.188998695259418: 1, 19.156132057876384: 1, 19.1225226755975: 1, 19.101928973748063: 1, 19.043381428941203: 1, 19.030095052148248: 1, 19.01455686672078: 1, 19.00859256066439: 1, 18.964636712879894: 1, 18.96356440498823: 1, 18.94862829968351: 1, 18.885454371433784: 1, 18.86326830611303: 1, 18.85988577612903: 1, 18.81985463469781: 1, 18.6872036922994: 1, 18.67960325278124: 1, 18.66583897505086: 1, 18.64615590948028: 1, 18.56717177259816: 1, 18.517809597729872: 1, 18.506565830363606: 1, 18.36398313514326: 1, 18.346587739755073: 1, 18.28577593435863: 1, 18.231920391879648: 1, 18.19381174180719: 1, 18.188844173442504: 1, 18.180489733913113: 1, 18.14135622059712: 1, 18.138985627021395: 1, 18.122794594920798: 1, 18.064137592834996: 1, 18.028227732474125: 1, 18.003690017389296: 1, 17.986169403263016: 1, 17.861714976636716: 1, 17.848701192334616: 1, 17.779273123478813: 1, 17.767117005485687: 1, 17.7528256981185: 1, 17.752020314086: 1, 17.694106446889904: 1, 17.684014165150792: 1, 17.612492386227547: 1,

17.56432812924743: 1, 17.529104914244154: 1, 17.502750575379377: 1, 17.461677945189674: 1, 17.454449866425023: 1, 17.441457928914428: 1, 17.43355508138719: 1, 17.43291011791032: 1, 17.43250550561128: 1, 17.428919205466133: 1, 17.408767452675644: 1, 17.403104659444843: 1, 17.3803512496207: 1, 17.37716101980046: 1, 17.35016078015181: 1, 17.30435669904955: 1, 17.281101088014708: 1, 17.275221096734267: 1, 17.194836502713333: 1, 17.17217868792897: 1, 17.149537968664237: 1, 17.07356056348511: 1, 17.070820850458087: 1, 17.059756483738738: 1, 17.038524061940063: 1, 17.016737788014964: 1, 16.97957038491065: 1, 16.96493791018695: 1, 16.961600446776643: 1, 16.953054876106183: 1, 16.9514127108774: 1, 16.926148564523658: 1, 16.900279056860462: 1, 16.870590985396085: 1, 16.823826309802058: 1, 16.787785082409613: 1, 16.7635628737114: 1, 16.753251254712225: 1, 16.702469090005525: 1, 16.675300251720305: 1, 16.6698559841232: 1, 16.613033607531897: 1, 16.55640769184553: 1, 16.555236548188546: 1, 16.54726824920114: 1, 16.50449772165057: 1, 16.484299272862298: 1, 16.41253353632573: 1, 16.381520036436136: 1, 16.330308612666258: 1, 16.314899872754076: 1, 16.293921963027678: 1, 16.247551406120603: 1, 16.247173925867013: 1, 16.231215394037545: 1, 16.221402681405067: 1, 16.218272455479966: 1, 16.1045110717932: 1, 16.10060166007644: 1, 16.095529525679268: 1, 16.074085060970425: 1, 16.070015263192328: 1, 16.05704283544889: 1, 16.04246783552156: 1, 15.973338617538174: 1, 15.964122253157324: 1, 15.960402150868468: 1, 15.943671869507092: 1, 15.934955565223753: 1, 15.856500127665319: 1, 15.7640834174792: 1, 15.672554975720038: 1, 15.609982427679144: 1, 15.591066416607728: 1, 15.579428197717524: 1, 15.493062911451627: 1, 15.48880815638546: 1, 15.46853352988342: 1, 15.41282438095239: 1, 15.40238022482248: 1, 15.36329867305021: 1, 15.335316240724495: 1, 15.312801290883481: 1, 15.306857394888297: 1, 15.290009798664688: 1, 15.26548833206873: 1, 15.223833578944676: 1, 15.219776081352299: 1, 15.205960625597383: 1, 15.185776074469171: 1, 15.15926298949431: 1, 15.150310010937115: 1, 15.088556586667998: 1, 15.072414631479255: 1, 15.04694056525533: 1, 15.045470817629122: 1, 15.040699789838904: 1, 15.024547386100886: 1, 15.010564498501855: 1, 14.930087247181389: 1, 14.929567232042892: 1, 14.923351364285628: 1, 14.909244319152673: 1, 14.864122912663573: 1, 14.858668761505607: 1, 14.854641049672109: 1, 14.848863618678235: 1, 14.823537679113295: 1, 14.812735498924413: 1, 14.785127084725934: 1, 14.733772895219118: 1, 14.726933033996099: 1, 14.688111817263783: 1, 14.687165545128018: 1, 14.679894104486934: 1, 14.621982522384698: 1, 14.612848581371352: 1, 14.611998006134257: 1, 14.57228672776145: 1, 14.560755614698989: 1, 14.539882683919618: 1, 14.477390366256055: 1, 14.455098581756534: 1, 14.442394272860591: 1, 14.416114647127728: 1, 14.395990993535923: 1, 14.38333460497038: 1, 14.37840063931081: 1, 14.372110620691894: 1, 14.369122149527414: 1, 14.3685417772663: 1, 14.365345545112204: 1, 14.346562023871575: 1, 14.320780752708835: 1, 14.319629174115057: 1, 14.30283196561286: 1, 14.298863327339259: 1, 14.241462828509286: 1, 14.21902916324693: 1, 14.217035866236628: 1, 14.186346178733086: 1, 14.130146775753772: 1, 14.111093891732738: 1, 14.106729743253933: 1, 14.090976020774006: 1, 14.078067982043606: 1, 14.070218969387238: 1, 14.018341442547584: 1, 14.00231547625213: 1, 13.99541526171637: 1, 13.944745732454228: 1, 13.91426766867743: 1, 13.913435546741297: 1, 13.8991595432368: 1, 13.869879120009259: 1, 13.864388467774653: 1, 13.847227699023295: 1, 13.839235976317385: 1, 13.815024987587147: 1, 13.785957137137734: 1, 13.743722918540518: 1, 13.735515582727173: 1, 13.734703414005311: 1, 13.730319760908202: 1, 13.72493773726309: 1, 13.711469703995185: 1, 13.674664984899652: 1, 13.659825578343433: 1, 13.64720096920117: 1, 13.64076150153843: 1, 13.62958399211476: 1, 13.62438542824137: 1, 13.622358020975602: 1, 13.615085121588843: 1, 13.607281502832771: 1, 13.60568154692528: 1, 13.57911644067107: 1, 13.543778250259825: 1, 13.535976426464504: 1, 13.513789299514992: 1, 13.49362058284742: 1, 13.472487639378937: 1, 13.469839786287203: 1, 13.444308716309319: 1, 13.37279040527896: 1, 13.328593978697647: 1, 13.317023360340876: 1, 13.285934188144303: 1, 13.283791931264606: 1, 13.211968118388919: 1, 13.20443725080553: 1, 13.15212780802418: 1, 13.141696724523749: 1, 13.1353144616356: 1, 13.113633057801431: 1, 13.11304875530662: 1, 13.100364766927433: 1, 13.091503450384979: 1, 13.082723780155746: 1, 13.07540425885229: 1, 13.066493125655951: 1, 13.025836154000777: 1, 13.022855222672144: 1, 13.004911568201035: 1, 12.98799839193278: 1, 12.969787842165124: 1, 12.957224769869915: 1, 12.928573557558913: 1, 12.923600067741335: 1, 12.88648554052569: 1, 12.880276068147875: 1, 12.879477209919765: 1, 12.873124958752681: 1, 12.843101216148453: 1, 12.818854427480714: 1, 12.795092143895435: 1, 12.744075073213285: 1, 12.680495388077498: 1, 12.659048658842504: 1, 12.650733991676619: 1, 12.63725555470167: 1, 12.624583369203016: 1, 12.58928405792037: 1, 12.56386227114931: 1, 12.55899631438682: 1, 12.555155426068874: 1, 12.551170336407049: 1, 12.520755965462545: 1, 12.501987368392777: 1, 12.489471869393471: 1, 12.46071784708502: 1, 12.44741937131256: 1, 12.424945576172211: 1, 12.422001914096962: 1, 12.390721978544706: 1, 12.386722139777543: 1, 12.38404504901519: 1, 12.382933124403424: 1, 12.381081550138505: 1, 12.324745521883692: 1, 12.311370879170191: 1, 12.244056558093261: 1, 12.196939255109573: 1, 12.189468086490379: 1, 12.18726263352704: 1, 12.131240485837997: 1, 12.13015061191147: 1, 12.119121860442396: 1, 12.10358417926974: 1, 12.077005183594983: 1, 12.068679290025425: 1, 11.968109424253099: 1, 11.926633133109338: 1, 11.919384040496775: 1, 11.914222587462538: 1, 11.893320594610358: 1, 11.882437194137175: 1, 11.879771412773465: 1, 11.864025456542318: 1, 11.861856164903514: 1, 11.856200612030667: 1, 11.836958006051152: 1, 11.822567937735373: 1, 11.82218315872492: 1, 11.816526390700703: 1, 11.81474498355431: 1, 11.799744274542912: 1, 11.742007355114595: 1, 11.741946631598635: 1, 11.732634714191741: 1, 11.730461015661355: 1, 11.719272304354778: 1, 11.70872684843321: 1, 11.696554772919605: 1, 11.680081938843461: 1, 11.66410454755786: 1, 11.644882503885524: 1, 11.636397303913105: 1, 11.612789784536341: 1, 11.598326325004315: 1, 11.591505147682257: 1, 11.567085101919643: 1, 11.565904460533872: 1, 11.564856004960479: 1, 11.536667012596826: 1, 11.529764788832882: 1, 11.492357225835036: 1, 11.477942859697658: 1, 11.45124370491169: 1, 11.443070937154333: 1, 11.437563081190921: 1, 11.430574807449375: 1, 11.42926902568096: 1, 11.311214561826223: 1, 11.305258449611085: 1, 11.292632228202242: 1, 11.288776204931635: 1, 11.262537323953987: 1, 11.246461392941061: 1, 11.231001800562746: 1, 11.228207127286328: 1, 11.217443490463188: 1, 11.202016872640884: 1, 11.194321025189499: 1, 11.174568520399117: 1, 11.163207172503185: 1, 11.162753906705037: 1, 11.156081132402967: 1,

11.154649877712602: 1, 11.146964885898596: 1, 11.139281534418952: 1, 11.13611998266071: 1, 11.127518793567322: 1, 11.11498627623893: 1, 11.114361762257985: 1, 11.113825426648521: 1, 11.113366574814641: 1, 11.10531461555147: 1, 11.082573249425668: 1, 11.076103617596738: 1, 11.070844368011837: 1, 11.061291070242477: 1, 11.037791558371406: 1, 11.012168189865202: 1, 11.011515191471293: 1, 10.991570627307414: 1, 10.986234423793327: 1, 10.978259751080788: 1, 10.977426796418309: 1, 10.965289394584557: 1, 10.960669424859693: 1, 10.932202785093468: 1, 10.92902684198268: 1, 10.920815874787925: 1, 10.917755578295: 1, 10.915440781820552: 1, 10.910826121165282: 1, 10.896111852842996: 1, 10.893968146368982: 1, 10.873184860412815: 1, 10.859548331176967: 1, 10.85110499443303: 1, 10.80380527073: 1, 10.77946961169958: 1, 10.748687118669588: 1, 10.741807894621152: 1, 10.7367752729409: 1, 10.72206578887616: 1, 10.718909239314035: 1, 10.715116055917399: 1, 10.712341260782708: 1, 10.712300782199483: 1, 10.69393590544361: 1, 10.689539063734019: 1, 10.669487396129346: 1, 10.6591671306079: 1, 10.654349302044798: 1, 10.641831282379938: 1, 10.635523013769772: 1, 10.630218509980601: 1, 10.612704813490208: 1, 10.612564966063887: 1, 10.575626355735857: 1, 10.572736772059542: 1, 10.556631301795761: 1, 10.544088653468195: 1, 10.54121088496635: 1, 10.521308494144613: 1, 10.51732941364631: 1, 10.509731384639629: 1, 10.503758992334767: 1, 10.481753069585864: 1, 10.47753205028488: 1, 10.461079667361282: 1, 10.447037894447492: 1, 10.426417643429822: 1, 10.417658790306405: 1, 10.412419315298106: 1, 10.39763684363565: 1, 10.351306634118394: 1, 10.350979133218873: 1, 10.34706491490113: 1, 10.320321527963197: 1, 10.3089469588102: 1, 10.295284433924168: 1, 10.288480389195811: 1, 10.279610161913368: 1, 10.267211234490164: 1, 10.266839450742582: 1, 10.264413563047997: 1, 10.240595946644488: 1, 10.236599833416756: 1, 10.205541146267073: 1, 10.196902076015265: 1, 10.182580471395301: 1, 10.177369657235808: 1, 10.118491868940506: 1, 10.103026549674759: 1, 10.084808923731389: 1, 10.083905375095872: 1, 10.066479117522343: 1, 10.060588010423068: 1, 10.054700791845027: 1, 10.051334116396262: 1, 10.042608528465854: 1, 10.03555151719483: 1, 10.022713126579752: 1, 10.006886024792747: 1, 10.000801648590702: 1, 9.992826908593106: 1, 9.97883970962311: 1, 9.966081194106573: 1, 9.9641477548875: 1, 9.960519723651172: 1, 9.947201796957964: 1, 9.9309928539397: 1, 9.921419542353055: 1, 9.912953656200003: 1, 9.911639223688043: 1, 9.910857476562098: 1, 9.909949491395746: 1, 9.881534255601021: 1, 9.880179017438603: 1, 9.870622764938348: 1, 9.866983622707354: 1, 9.842840314095108: 1, 9.839599169626057: 1, 9.818286999018985: 1, 9.81169672030304: 1, 9.807206892975472: 1, 9.80504332739838: 1, 9.764823771418135: 1, 9.759423824798754: 1, 9.749351610027329: 1, 9.741940845924807: 1, 9.738301660905053: 1, 9.734753506979285: 1, 9.72829443206632: 1, 9.721862480937698: 1, 9.721480280149837: 1, 9.720772912583879: 1, 9.711054873469433: 1, 9.707183613988406: 1, 9.70339478467253: 1, 9.702688993783898: 1, 9.67945855665573: 1, 9.677701368280628: 1, 9.644264122307803: 1, 9.643828516385222: 1, 9.64070626280922: 1, 9.629424234908203: 1, 9.604800688670583: 1, 9.604487478901317: 1, 9.601780761102551: 1, 9.592694365394483: 1, 9.591162871573143: 1, 9.586227423561722: 1, 9.581491352321436: 1, 9.578914113798614: 1, 9.547153479886848: 1, 9.538062127568407: 1, 9.512095992400603: 1, 9.50230767351044: 1, 9.464256009058735: 1, 9.463945296276586: 1, 9.456746528068289: 1, 9.454694457615581: 1, 9.409953857378687: 1, 9.408901613994228: 1, 9.393296627602102: 1, 9.38772956693125: 1, 9.384339668960356: 1, 9.383356613332849: 1, 9.371495824981215: 1, 9.363179752695276: 1, 9.335370842415548: 1, 9.335176354060206: 1, 9.325512524732996: 1, 9.29488073180698: 1, 9.292397278041982: 1, 9.290700439310944: 1, 9.288649658400802: 1, 9.276690047871742: 1, 9.273265157719091: 1, 9.26424353103488: 1, 9.24066858763401: 1, 9.236351671324307: 1, 9.229155241791995: 1, 9.210844271797386: 1, 9.201795006757555: 1, 9.181594392485794: 1, 9.171359897189605: 1, 9.163190434045532: 1, 9.15698128738799: 1, 9.154802300581878: 1, 9.154595606868067: 1, 9.14379536944781: 1, 9.139013298745105: 1, 9.138092804139777: 1, 9.13379375570124: 1, 9.122194902930525: 1, 9.120854738578759: 1, 9.109831868760384: 1, 9.104683830997251: 1, 9.089595398360471: 1, 9.079402777405654: 1, 9.079250653351766: 1, 9.070234614943823: 1, 9.06442989758896: 1, 9.052324419189379: 1, 9.02339280074982: 1, 9.01749774750834: 1, 9.015444094008089: 1, 9.007741267540823: 1, 9.003992256516844: 1, 8.983250178548621: 1, 8.968398759286705: 1, 8.95681521440377: 1, 8.937311716912966: 1, 8.93700343417962: 1, 8.93070613774506: 1, 8.9229819588102786: 1, 8.92240336968132: 1, 8.907256317773685: 1, 8.904093059036509: 1, 8.891703879838813: 1, 8.888071890793134: 1, 8.86944659219443: 1, 8.864283543207138: 1, 8.84882332945963: 1, 8.83955176328561: 1, 8.82757790205663: 1, 8.799827824019868: 1, 8.797113790788053: 1, 8.786948947045504: 1, 8.778088123485023: 1, 8.776916861947655: 1, 8.757384865787335: 1, 8.753516325448691: 1, 8.74151048640639: 1, 8.725610829310956: 1, 8.72199835500032: 1, 8.71804302058611: 1, 8.71457402660962: 1, 8.714007676130196: 1, 8.701060322728283: 1, 8.699566210225786: 1, 8.69895936314557: 1, 8.695576496990906: 1, 8.692637849981702: 1, 8.645648432601362: 1, 8.643380856256893: 1, 8.63993855883173: 1, 8.606906246533772: 1, 8.603728306221871: 1, 8.589975455548151: 1, 8.582529677001379: 1, 8.559533618369533: 1, 8.531566024421718: 1, 8.52542605310792: 1, 8.523391700020087: 1, 8.505570290319058: 1, 8.498481640778367: 1, 8.473972225783083: 1, 8.473196736651337: 1, 8.460560020334707: 1, 8.44572152480745: 1, 8.443899239886331: 1, 8.43738520116142: 1, 8.414754787647508: 1, 8.40022782616007: 1, 8.365682123067868: 1, 8.363009554711603: 1, 8.360360022508328: 1, 8.335686907110887: 1, 8.319282984209863: 1, 8.311554199890262: 1, 8.28162020964199: 1, 8.279677377755883: 1, 8.26261904998381: 1, 8.261947910926017: 1, 8.25838334763426: 1, 8.257989727683157: 1, 8.250505872534873: 1, 8.244228370211095: 1, 8.229114129232551: 1, 8.219473249934255: 1, 8.215682971655546: 1, 8.209156921145247: 1, 8.201215177526175: 1, 8.186062212615937: 1, 8.18445093595589: 1, 8.178691762876038: 1, 8.168195126101224: 1, 8.14118663012003: 1, 8.133356207678274: 1, 8.12264810782736: 1, 8.113667579192622: 1, 8.100068525467792: 1, 8.098130680627959: 1, 8.083841619991455: 1, 8.06628355790877: 1, 8.057519477963078: 1, 8.03418945807185: 1, 8.020872005168085: 1, 8.011724327579648: 1, 7.987549902224742: 1, 7.987086532812995: 1, 7.976672324610291: 1, 7.959738331044812: 1, 7.9176629739348305: 1, 7.905352101370597: 1, 7.8899540714968035: 1, 7.875558636055444: 1, 7.858401512248233: 1, 7.847817431444963: 1, 7.831649968465452: 1,

```

7.81546761196545: 1, 7.800859062869251: 1, 7.784413424106252: 1, 7.7828745718286: 1,
7.773290927997185: 1, 7.771442192985451: 1, 7.7653153372596675: 1, 7.745793640183624: 1,
7.737456397414969: 1, 7.72496346993475: 1, 7.707371385023874: 1, 7.705679417358357: 1,
7.703083242655954: 1, 7.690844679400546: 1, 7.67670626131244: 1, 7.667242413835371: 1,
7.6381708229154865: 1, 7.626811176253351: 1, 7.619179483759866: 1, 7.617002841742056: 1, 7.61627828
9009852: 1, 7.615966427331192: 1, 7.581513507093012: 1, 7.570199103938951: 1, 7.553553649744871: 1,
7.546893558721115: 1, 7.530773569704344: 1, 7.514656032865905: 1, 7.510088035361333: 1,
7.5088417163870425: 1, 7.506469279453428: 1, 7.496269820733249: 1, 7.461058811307257: 1, 7.44457872
9324948: 1, 7.411258044725352: 1, 7.396786966175162: 1, 7.378977660842302: 1, 7.371968032235573: 1,
7.3714009144596595: 1, 7.352707802785961: 1, 7.344903530636471: 1, 7.316212130973699: 1, 7.31464273
76271905: 1, 7.3014006642003535: 1, 7.28877629059352: 1, 7.284052637790303: 1, 7.280033742176097:
1, 7.256982502755523: 1, 7.206407128559852: 1, 7.141103517404148: 1, 7.139962801825403: 1,
7.137483027825605: 1, 7.136015258646175: 1, 7.098748269565287: 1, 7.069348527903912: 1,
6.972057143108663: 1, 6.963402112108019: 1, 6.940431219808227: 1, 6.891975616817441: 1,
6.888509833990778: 1, 6.8883596391138955: 1, 6.8325742649672865: 1, 6.831826197262968: 1,
6.803221626799594: 1, 6.794956629775335: 1, 6.7553998737760494: 1, 6.571062485971009: 1,
6.470473003425169: 1, 6.431501569197116: 1})

```

In [61]:

```

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

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

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

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

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

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

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

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

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

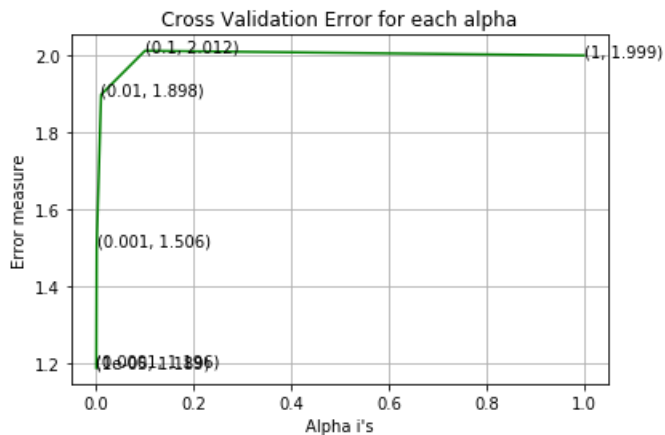
```

```

# For values of alpha = 1e-05, alpha[0], the cross validation log loss is: log_loss(
ss(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, p
redict_y, labels=clf.classes_, eps=1e-15))

```

For values of alpha = 1e-05 The log loss is: 1.1892492459756117
 For values of alpha = 0.0001 The log loss is: 1.195695724387353
 For values of alpha = 0.001 The log loss is: 1.5055121231895068
 For values of alpha = 0.01 The log loss is: 1.8978294826628537
 For values of alpha = 0.1 The log loss is: 2.0116339136561194
 For values of alpha = 1 The log loss is: 1.9993900184427222



For values of best alpha = 1e-05 The train log loss is: 0.7663741383515251
 For values of best alpha = 1e-05 The cross validation log loss is: 1.1892492459756117
 For values of best alpha = 1e-05 The test log loss is: 1.1146751501392769

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

Ans. Yes, it seems like!

In [62]:

```

def get_intersec_text(df):
    df_text_vec = TfidfVectorizer(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 [63]:

```

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

```

3.457 % of word of test data appeared in train data
 3.951 % of word of Cross Validation appeared in train data

4. Machine Learning Models

In [64]:

```

#Data preparation for ML models.

#Misc. functionns 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 [65]:

```
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 [66]:

```
# 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")
```

Stacking the three types of features

In [67]:

```
# merging gene, variance and text features

# building train, test and cross validation data sets
# a = [[1, 2],
```

```
#      [3, 4]]
# b = [[4, 5],
#      [6, 7]]
# hstack(a, b) = [[1, 2, 4, 5],
#                [ 3, 4, 6, 7]]

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

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

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

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

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

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

In [70]:

```
print("TFIDF 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)
```

```
TFIDF encoding features :
(number of data points * number of features) in train data = (2124, 3189)
(number of data points * number of features) in test data = (665, 3189)
(number of data points * number of features) in cross validation data = (532, 3189)
```

In [69]:

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

```
Response encoding features :
(number of data points * number of features) in train data = (2124, 27)
(number of data points * number of features) in test data = (665, 27)
(number of data points * number of features) in cross validation data = (532, 27)
```

4.1. Base Line Model

4.1.1. Naive Bayes

4.1.1.1. Hyper parameter tuning

.....Hyper parameter tuning

In [71]:

```
# 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.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

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 probabilitites 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()
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)

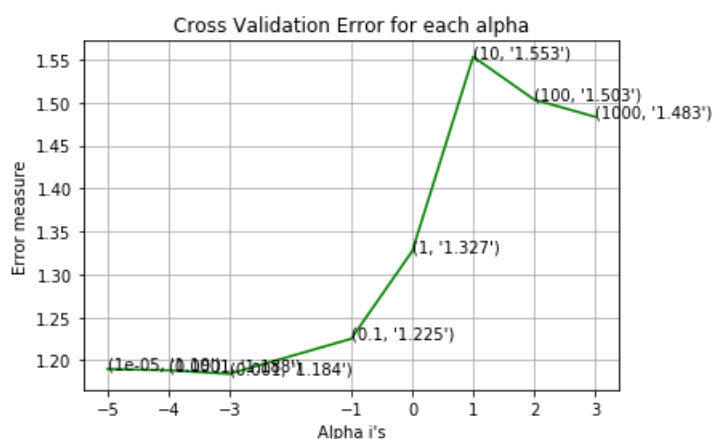
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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



```

for alpha = 1e-05
Log Loss : 1.1898024396943387
for alpha = 0.0001
Log Loss : 1.188439999818832
for alpha = 0.001
Log Loss : 1.1842459931589249
for alpha = 0.1
Log Loss : 1.225093845236475
for alpha = 1
Log Loss : 1.327002537636555
for alpha = 10
Log Loss : 1.553248515320758
for alpha = 100
Log Loss : 1.5033815964164934
for alpha = 1000
Log Loss : 1.4833322184388356

```



For values of best alpha = 0.001 The train log loss is: 0.5224203125858754
 For values of best alpha = 0.001 The cross validation log loss is: 1.1842459931589249
 For values of best alpha = 0.001 The test log loss is: 1.153402498691639

4.1.1.2. Testing the model with best hyper paramters

In [72]:

```

# 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.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

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

```

```

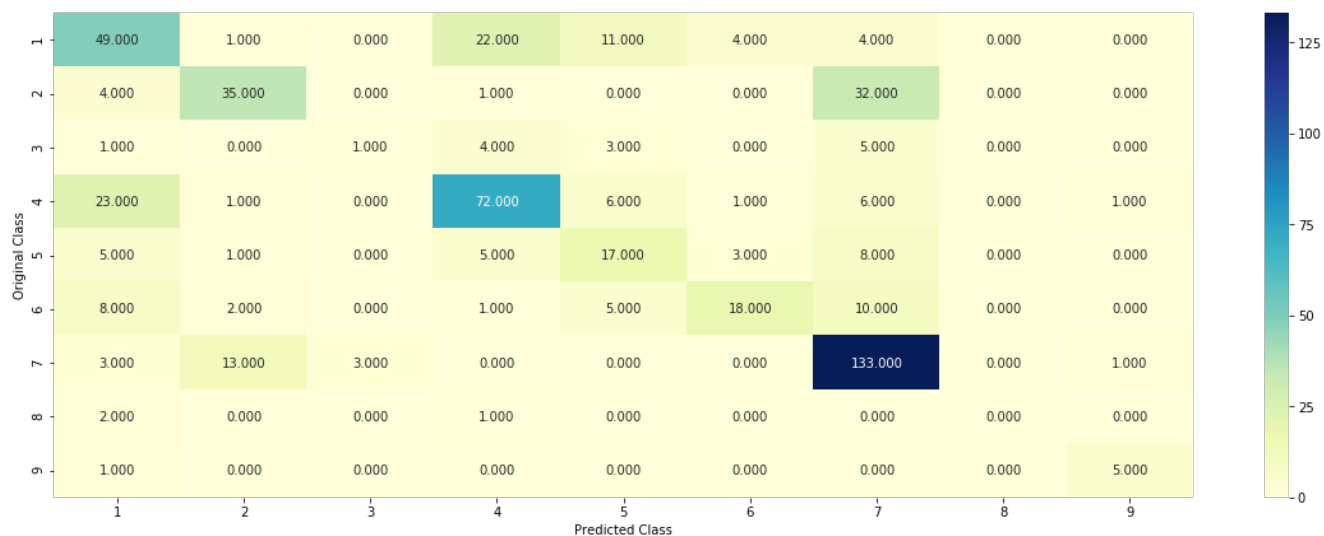
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
# to avoid rounding error while multiplying probabilities we 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()))

```

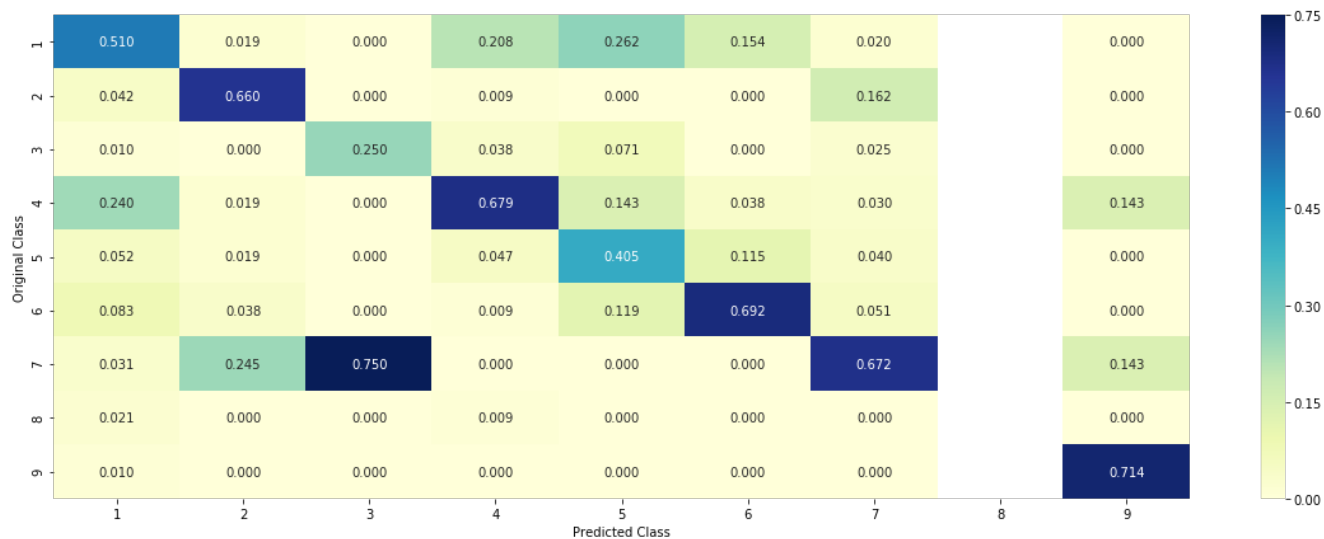
Log Loss : 1.1842459931589249

Number of missclassified point : 0.37969924812030076

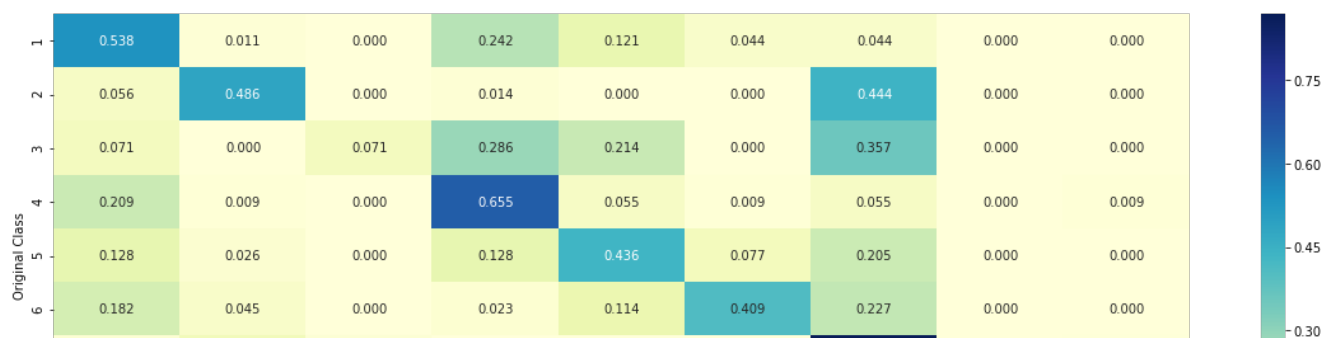
----- Confusion matrix -----

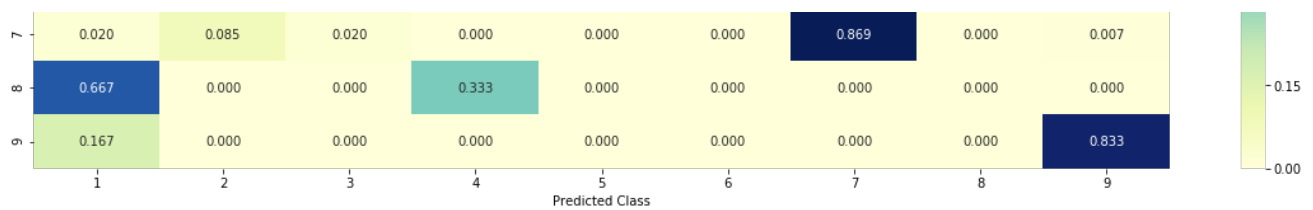


----- Precision matrix (Column Sum=1) -----



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





4.1.1.3. Feature Importance, Correctly classified point

In [73]:

```
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 : 7
Predicted Class Probabilities: [[0.0623 0.1902 0.0109 0.0708 0.033 0.0339 0.5898 0.0051 0.004 ]]
Actual Class : 7
-----
24 Text feature [101] present in test data point [True]
33 Text feature [114] present in test data point [True]
54 Text feature [02] present in test data point [True]
74 Text feature [12] present in test data point [True]
77 Text feature [030] present in test data point [True]
Out of the top 100 features 5 are present in query point
```

4.1.1.4. Feature Importance, Incorrectly classified point

In [74]:

```
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: [[0.0568 0.0469 0.0101 0.0646 0.0299 0.0429 0.7405 0.0046 0.0037]]
Actual Class : 7
-----
74 Text feature [12] present in test data point [True]
Out of the top 100 features 1 are present in query point
```

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

In [75]:

```
# 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.
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-geometric-intuition-with-a-toy-example-1/
#-----

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

alpha = [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 probabilities we use log-probability estimates
    print("Log Loss :", log_loss(cv_y, sig_clf.probs))

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

best_alpha = np.argmin(cv_log_error_array)
clf = 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_lo
ss(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, p
redict_y, labels=clf.classes_, eps=1e-15))

```

```

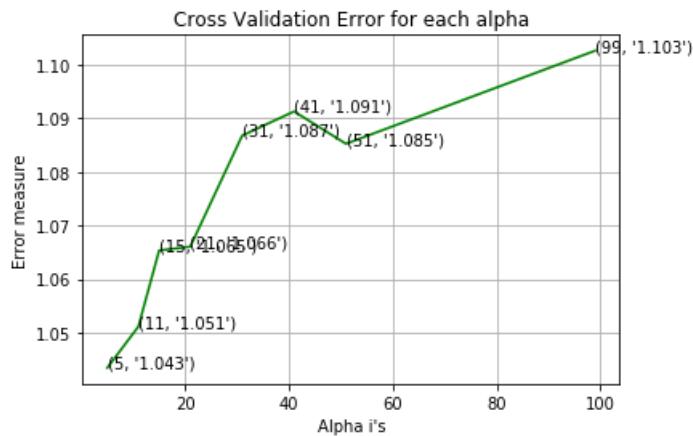
for alpha = 5
Log Loss : 1.0434285303443476
for alpha = 11
Log Loss : 1.0511245004990535

```

```

for alpha = 15
Log Loss : 1.065338805834702
for alpha = 21
Log Loss : 1.0660451815247383
for alpha = 31
Log Loss : 1.0867686210331065
for alpha = 41
Log Loss : 1.0912797574812965
for alpha = 51
Log Loss : 1.0852408575101333
for alpha = 99
Log Loss : 1.1026063754098931

```



For values of best alpha = 5 The train log loss is: 0.49982161166797434
 For values of best alpha = 5 The cross validation log loss is: 1.0434285303443476
 For values of best alpha = 5 The test log loss is: 1.046678499933953

4.2.2. Testing the model with best hyper paramters

In [76]:

```

# 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)

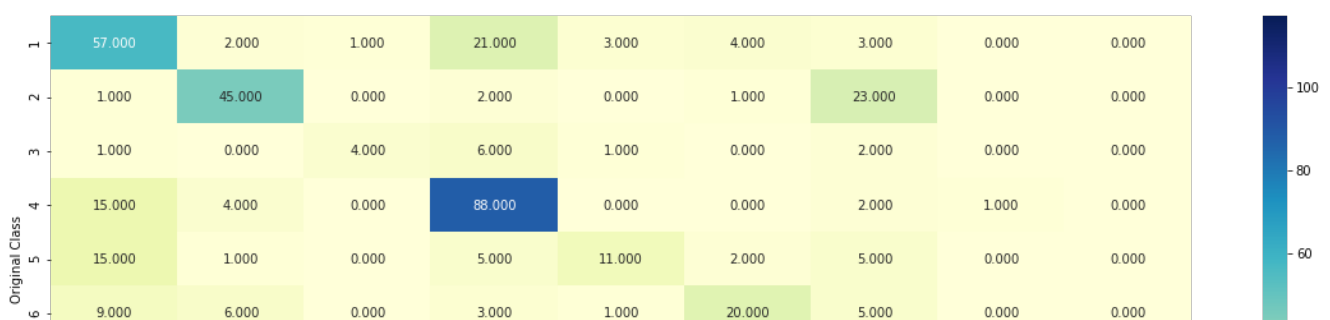
# methods of
# fit(X, y) : Fit the model using X as training data and y as target values
# predict(X):Predict the class labels for the provided data
# predict_proba(X):Return probability estimates for the test data X.
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-geometric-intuition-with-a-toy-example-1/
#-----
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseCoding, cv_y, clf)

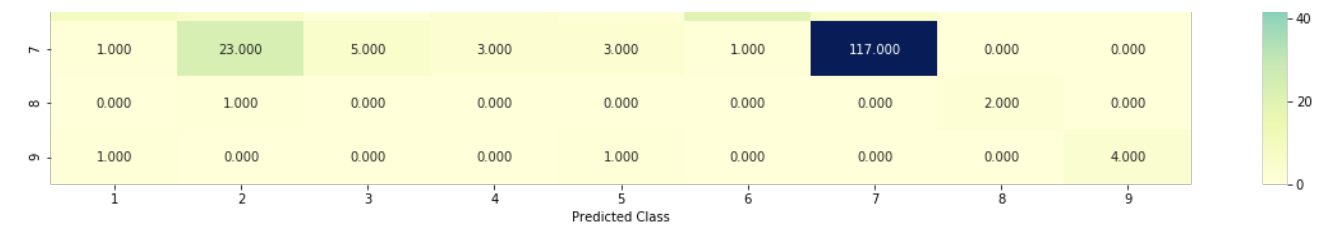
```

Log loss : 1.0434285303443476

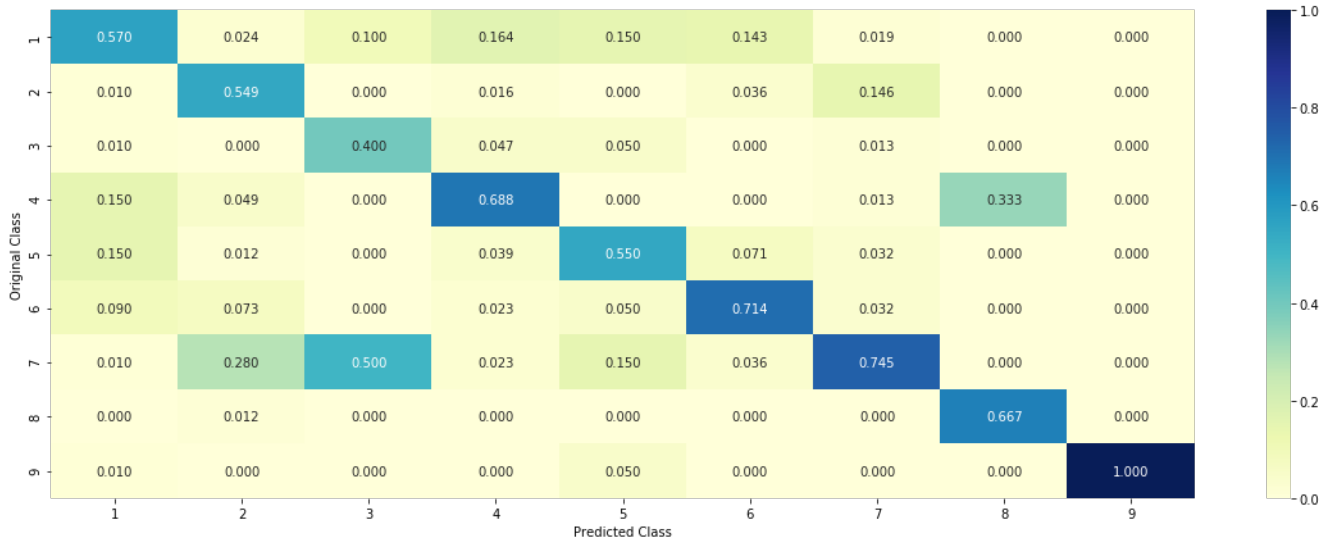
Number of mis-classified points : 0.3458646616541353

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

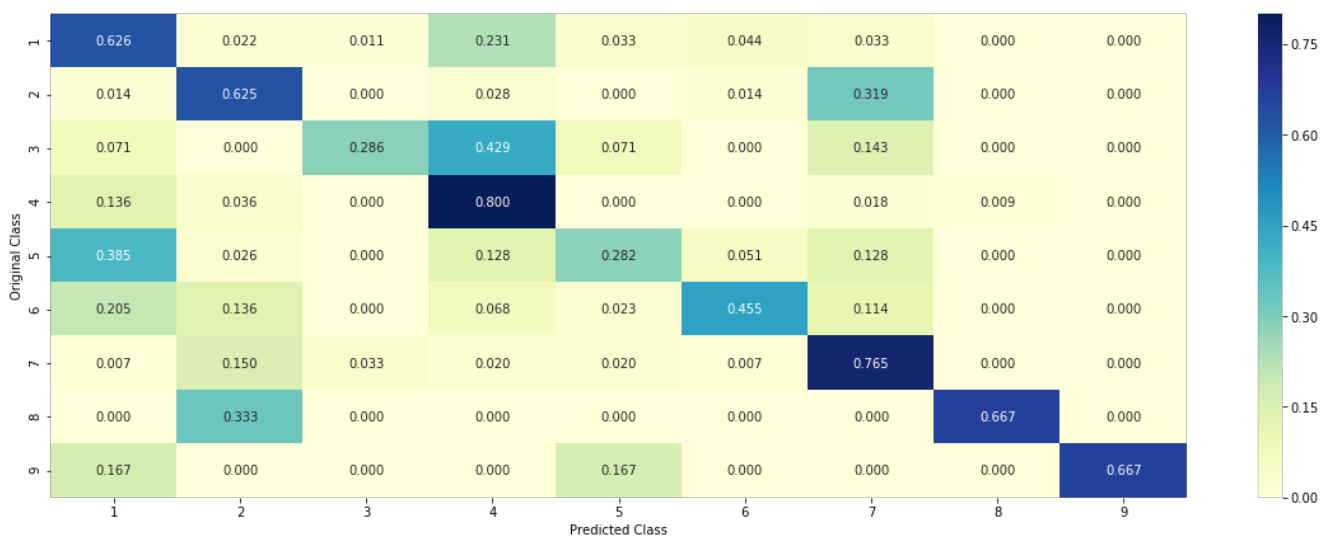




----- Precision matrix (Column Sum=1) -----



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



4.2.3.Sample Query point -1

In [77]:

```
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 : 6

Actual Class : 7

The 5 nearest neighbours of the test points belongs to classes [6 7 7 7 2]

Fequency of nearest points : Counter({7: 3, 6: 1, 2: 1})

4.2.4. Sample Query Point-2

In [78]:

```
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 be
longs 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 5 and the nearest neighbours of the test points belongs to classes [6 7 6 2 7]

Fequency of nearest points : Counter({6: 2, 7: 2, 2: 1})

4.3. Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper paramter tuning

In [79]:

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

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

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

# find more about CalibratedClassifierCV here at http://scikit-
learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()
```

```

# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
#-----
# video link:
#-----

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

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

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

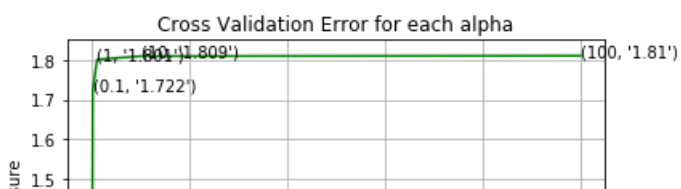
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

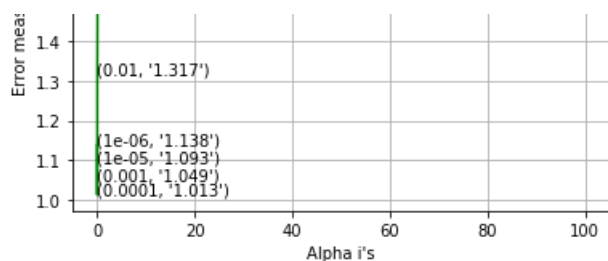
```

```

for alpha = 1e-06
Log Loss : 1.1381131331794327
for alpha = 1e-05
Log Loss : 1.0929509399434394
for alpha = 0.0001
Log Loss : 1.0128966329739684
for alpha = 0.001
Log Loss : 1.0487522351560081
for alpha = 0.01
Log Loss : 1.3166921136013736
for alpha = 0.1
Log Loss : 1.7224148215378734
for alpha = 1
Log Loss : 1.800967533828814
for alpha = 10
Log Loss : 1.8094254142192723
for alpha = 100
Log Loss : 1.8104034038093741

```





For values of best alpha = 0.0001 The train log loss is: 0.45981489120467645
 For values of best alpha = 0.0001 The cross validation log loss is: 1.0128966329739684
 For values of best alpha = 0.0001 The test log loss is: 1.0211072322802177

4.3.1.2. Testing the model with best hyper paramters

In [80]:

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

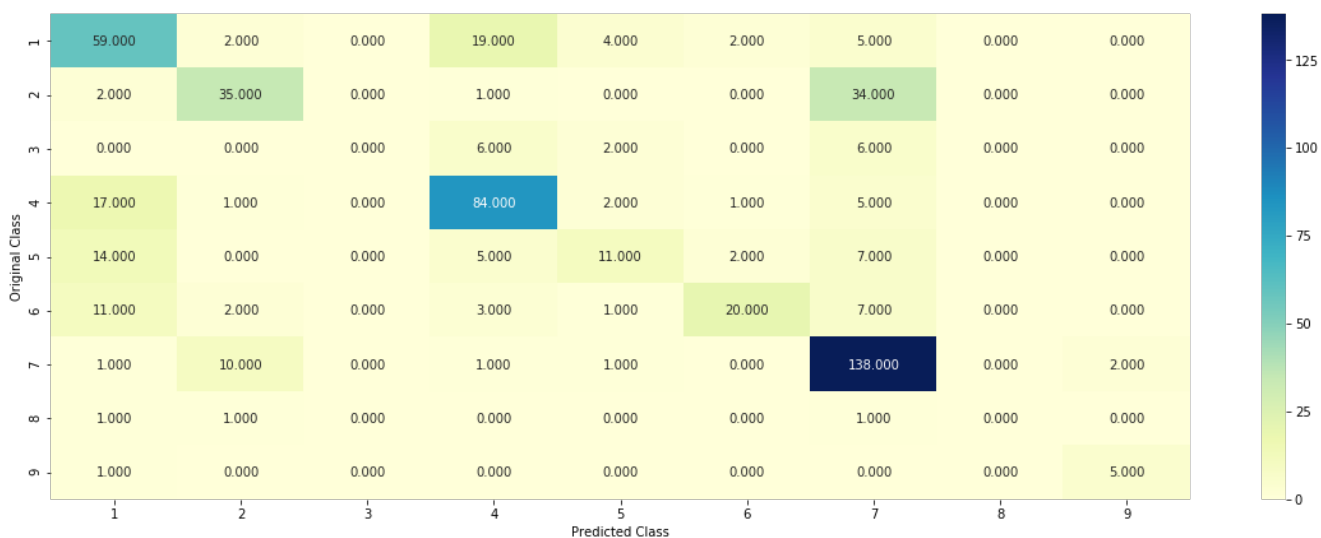
# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```

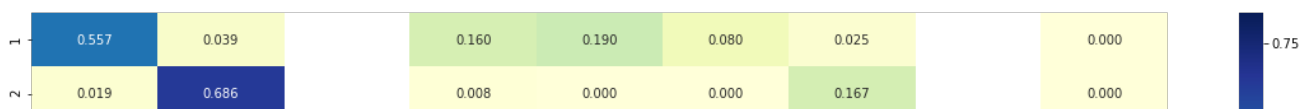
Log loss : 1.0128966329739684

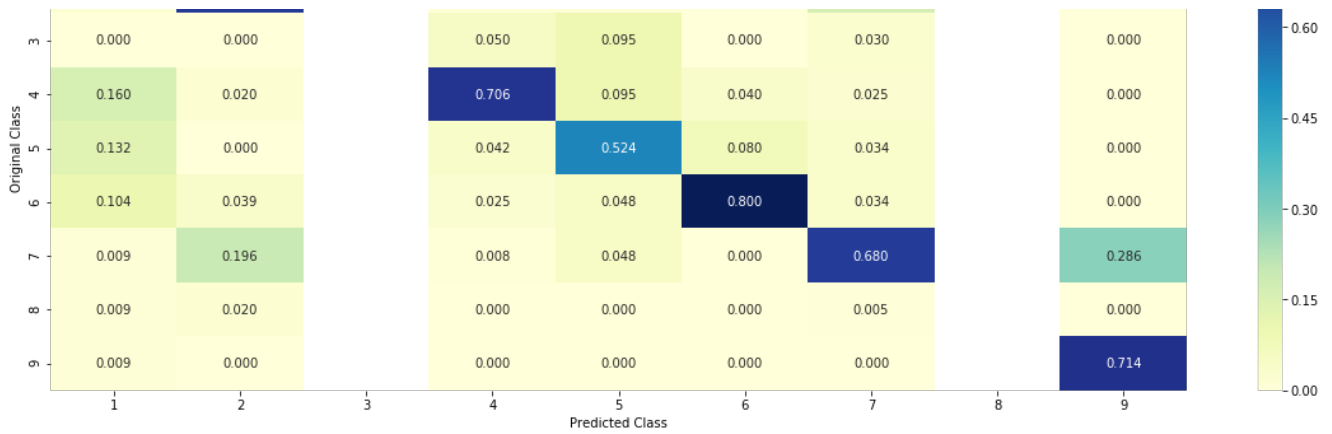
Number of mis-classified points : 0.3383458646616541

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

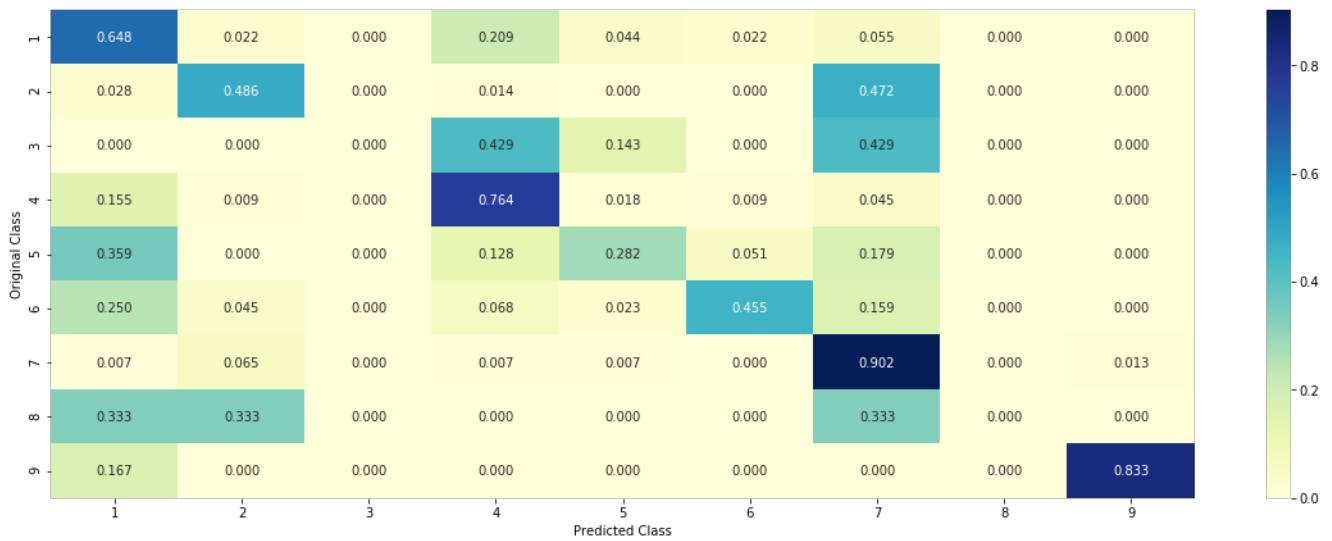


----- Precision matrix (Column Sum=1) -----





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



4.3.1.3. Feature Importance

In [81]:

```
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 important of the ", predicted_cls[0], " class:")
    print(tabulate(tabulte_list, headers=["Index", "Feature name", "Present or Not"]))
```

4.3.1.3.1. Correctly Classified point

In [82]:

```
# from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
```

```

coef_state=12,
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 : 7
Predicted Class Probabilities: [[1.200e-03 2.445e-01 1.000e-04 2.200e-03 9.000e-04 6.300e-03 7.300
e-01
1.470e-02 1.000e-04]]
Actual Class : 7
-----
7 Text feature [129] present in test data point [True]
61 Text feature [12q13] present in test data point [True]
74 Text feature [03] present in test data point [True]
131 Text feature [02] present in test data point [True]
211 Text feature [00] present in test data point [True]
250 Text feature [018] present in test data point [True]
Out of the top 500 features 6 are present in query point

```

4.3.1.3.2. Incorrectly Classified point

In [83]:

```

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: [[0.0295 0.0102 0.0024 0.0595 0.0047 0.2219 0.6608 0.0091 0.0019]]
Actual Class : 7
-----
Out of the top 500 features 0 are present in query point

```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

In [84]:

```

# 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_i
ter=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

```

```

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

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

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

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

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

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

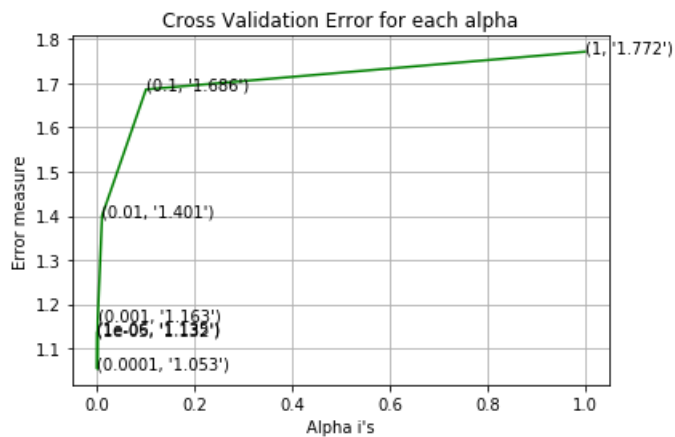
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 1e-06
Log Loss : 1.1354381283989008
for alpha = 1e-05
Log Loss : 1.1321472782538784
for alpha = 0.0001
Log Loss : 1.0533680070261722
for alpha = 0.001
Log Loss : 1.1628069895895135
for alpha = 0.01
Log Loss : 1.4005410489445727
for alpha = 0.1
Log Loss : 1.6862801728509882
for alpha = 1
Log Loss : 1.7716253046241515

```



For values of best alpha = 0.0001 The train log loss is: 0.4498471771179781
 For values of best alpha = 0.0001 The cross validation log loss is: 1.0533680070261722
 For values of best alpha = 0.0001 The test log loss is: 1.0417844976004644

4.3.2.2. Testing model with best hyper parameters

In [85]:

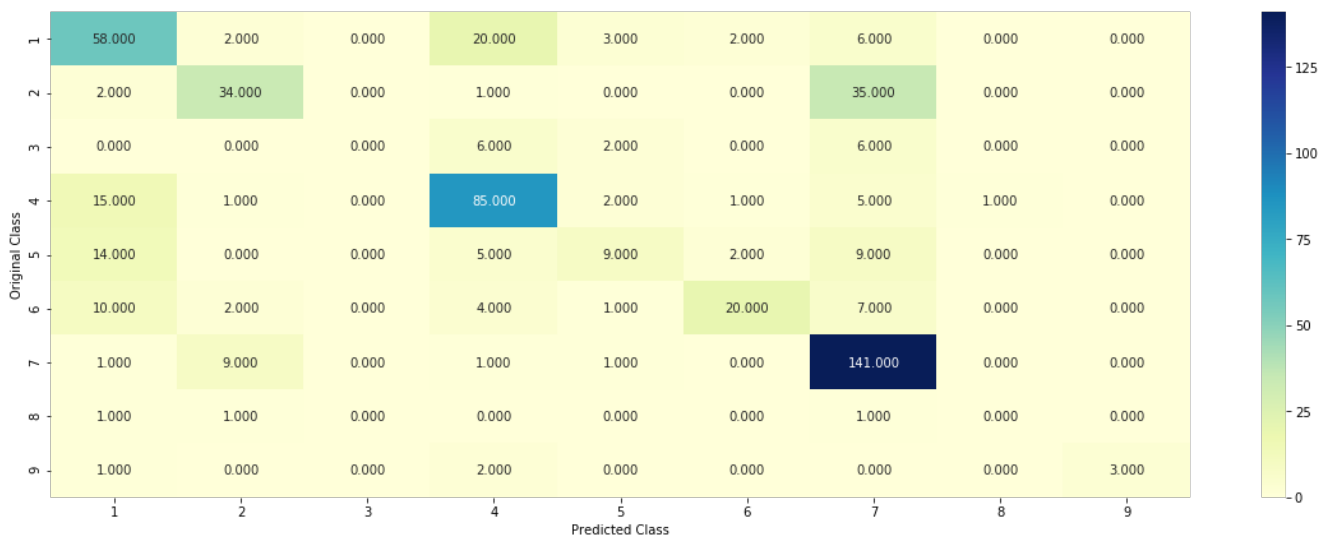
```
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0,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

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

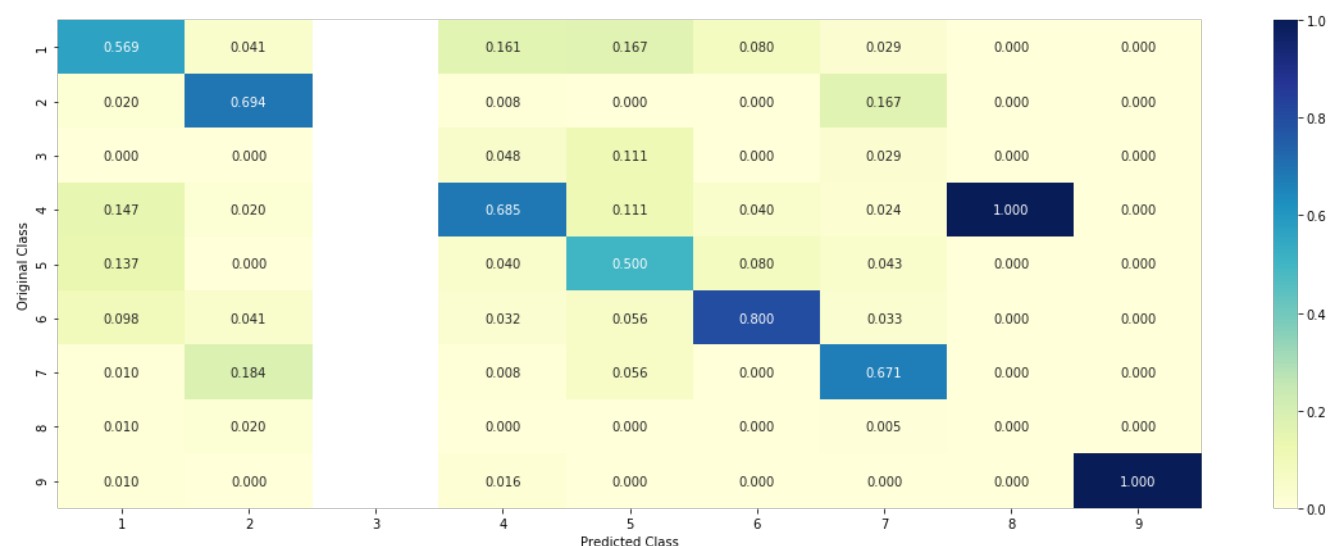
#-----
# video link:
#-----

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

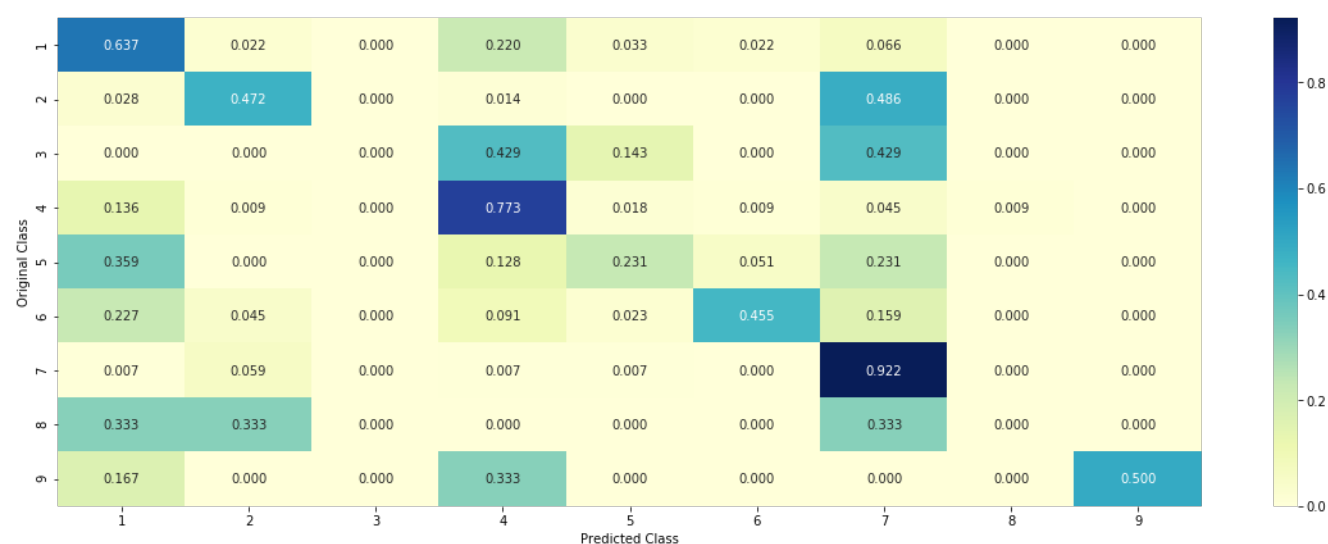
Log loss : 1.0533680070261722
 Number of mis-classified points : 0.34210526315789475
 ----- Confusion matrix -----



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



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



4.3.2.3. Feature Importance, Correctly Classified point

In [86]:

```
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 : 7

Predicted Class Probabilities: [[1.100e-03 2.167e-01 1.000e-04 2.100e-03 6.000e-04 4.800e-03 7.532e-01
2.130e-02 0.000e+00]]

Actual Class : 7

```

11 Text feature [129] present in test data point [True]
77 Text feature [12q13] present in test data point [True]
85 Text feature [03] present in test data point [True]
204 Text feature [00] present in test data point [True]
251 Text feature [02] present in test data point [True]
298 Text feature [018] present in test data point [True]
Out of the top 500 features 6 are present in query point

```

4.3.2.4. Feature Importance, Inorrectly Classified point

In [87]:

```

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: [[2.850e-02 9.400e-03 1.900e-03 6.150e-02 3.900e-03 1.860e-01 6.942
e-01
1.400e-02 6.000e-04]]
Actual Class : 7
-----
Out of the top 500 features 0 are present in query point

```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

In [88]:

```

# 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.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

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

```

```

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

alpha = [10 ** x for x in range(-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()
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)

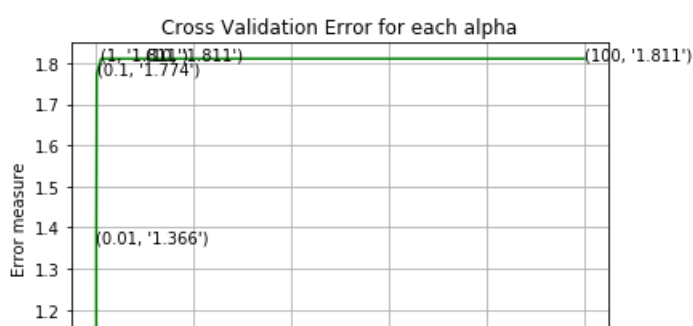
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

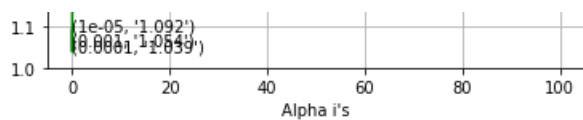
```

```

for C = 1e-05
Log Loss : 1.0921117673770042
for C = 0.0001
Log Loss : 1.0390641156737952
for C = 0.001
Log Loss : 1.0541275352330894
for C = 0.01
Log Loss : 1.365600976203778
for C = 0.1
Log Loss : 1.774224512949043
for C = 1
Log Loss : 1.8106217692080988
for C = 10
Log Loss : 1.8106216956546028
for C = 100
Log Loss : 1.8106219061360804

```





For values of best alpha = 0.0001 The train log loss is: 0.49326857709182725
 For values of best alpha = 0.0001 The cross validation log loss is: 1.0390641156737952
 For values of best alpha = 0.0001 The test log loss is: 1.0908378493314566

4.4.2. Testing model with best hyper parameters

In [89]:

```
# 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)

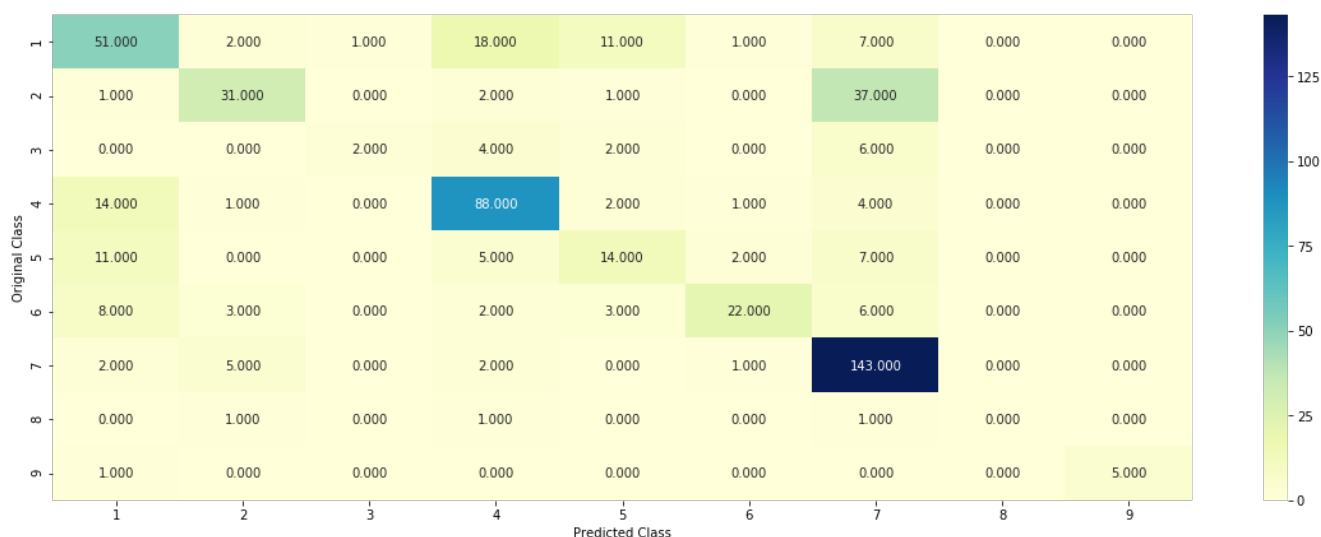
# Some of methods of SVM()
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
# predict(X) Perform classification on samples in X.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

# clf = SVC(C=alpha[best_alpha],kernel='linear',probability=True, class_weight='balanced')
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)
```

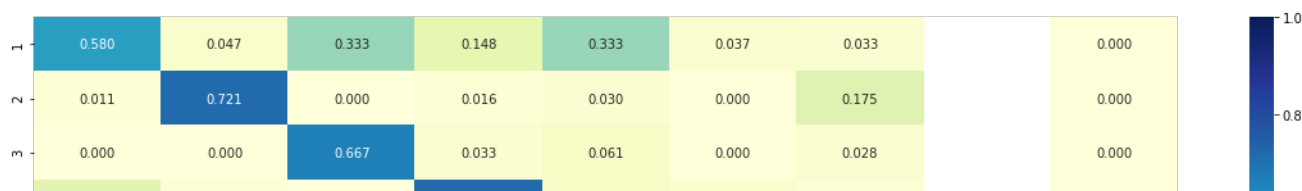
Log loss : 1.0390641156737952

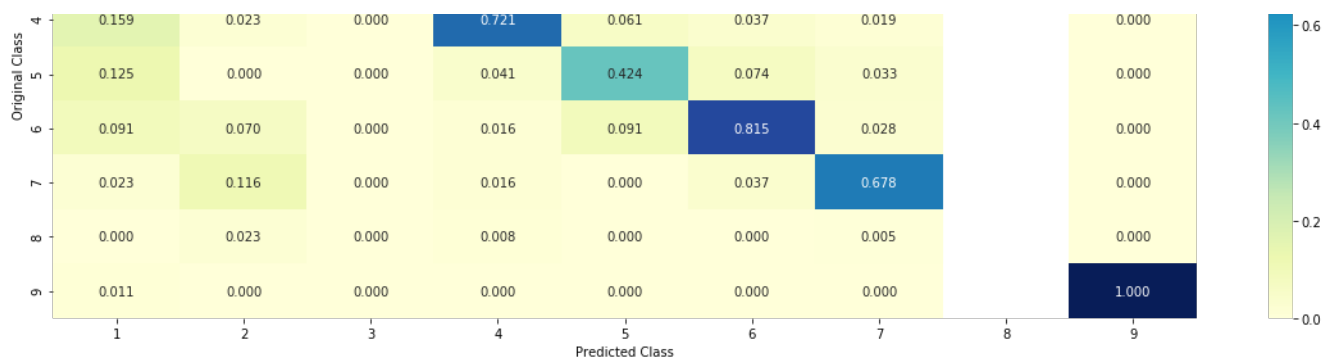
Number of mis-classified points : 0.3308270676691729

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

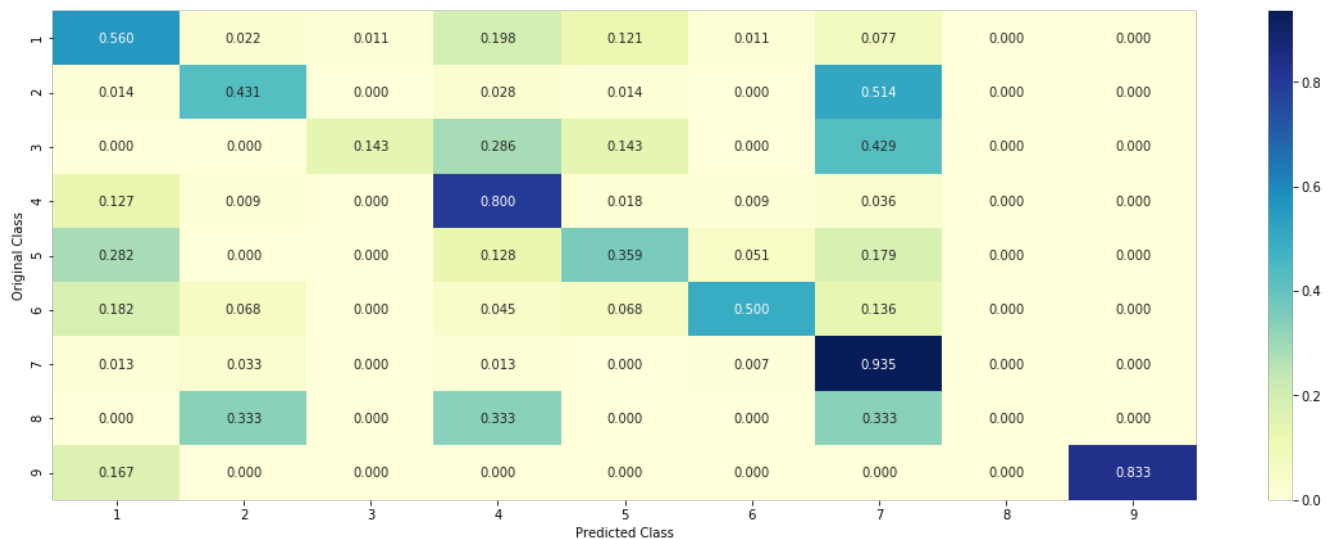


----- Precision matrix (Column Sum=1) -----





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



4.3.3. Feature Importance

4.3.3.1. For Correctly classified point

In [90]:

```
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 : 7

Predicted Class Probabilities: [[7.400e-03 2.823e-01 3.000e-04 8.100e-03 3.600e-03 9.500e-03 6.718e-01

1.660e-02 5.000e-04]]

Actual Class : 7

31 Text feature [129] present in test data point [True]
221 Text feature [03] present in test data point [True]
237 Text feature [00] present in test data point [True]
253 Text feature [12q13] present in test data point [True]
Out of the top 500 features 4 are present in query point

4.3.3.2. For Incorrectly classified point

In [91]:

```
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: [[0.0324 0.0083 0.0039 0.1491 0.011 0.1664 0.6195 0.0066 0.0027]]

Actual Class : 7

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

4.5 Random Forest Classifier

4.5.1. Hyper paramter tuning (With Tfidf Encoding)

In [92]:

```
# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_s
amples_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).

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

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

alpha = [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 : 1.2091101186993678
for n_estimators = 100 and max depth = 10
Log Loss : 1.2365656830349305
for n_estimators = 200 and max depth = 5
Log Loss : 1.2013494840592758
for n_estimators = 200 and max depth = 10
Log Loss : 1.2291548159334766
for n_estimators = 500 and max depth = 5
Log Loss : 1.1945011693162964
for n_estimators = 500 and max depth = 10
Log Loss : 1.221654131823169
for n_estimators = 1000 and max depth = 5
Log Loss : 1.1932758707533517
for n_estimators = 1000 and max depth = 10
Log Loss : 1.2162431181874434
for n_estimators = 2000 and max depth = 5
Log Loss : 1.1905755152282855
for n_estimators = 2000 and max depth = 10
Log Loss : 1.2145997702195375
For values of best estimator = 2000 The train log loss is: 0.8550990174597015
For values of best estimator = 2000 The cross validation log loss is: 1.1905755152282853
For values of best estimator = 2000 The test log loss is: 1.164721249859326

```

4.5.2. Testing model with best hyper parameters (Tfidf Encoding)

In [93]:

```

# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_s

```

```

# 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).

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

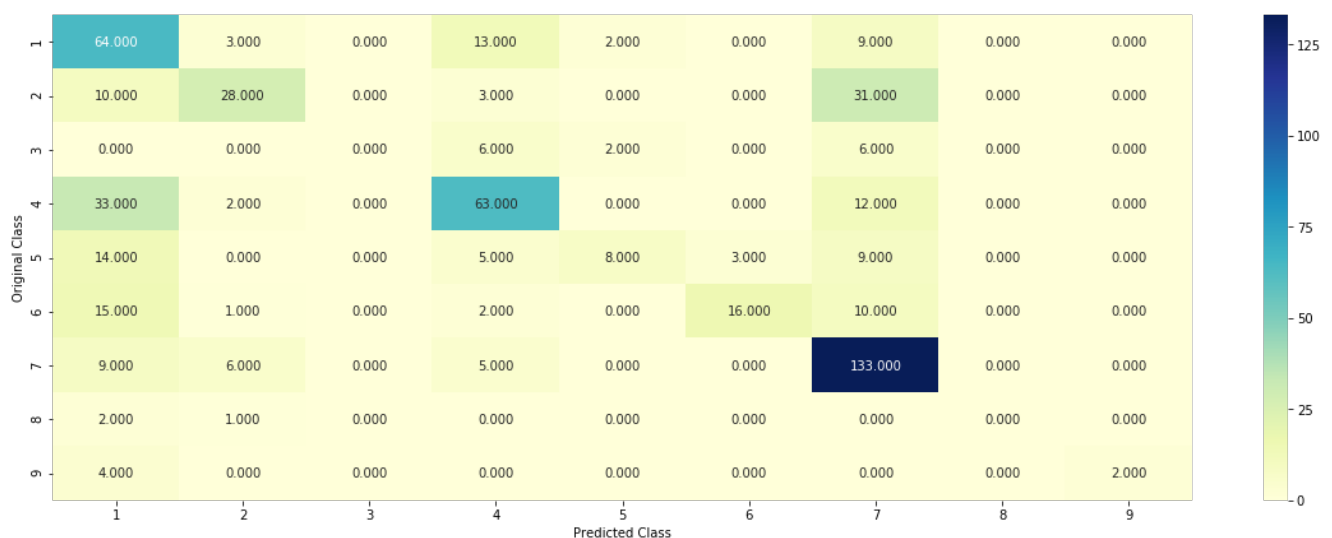
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_
depth[int(best_alpha*2)], random_state=42, n_jobs=-1)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding,cv_y, clf)

```

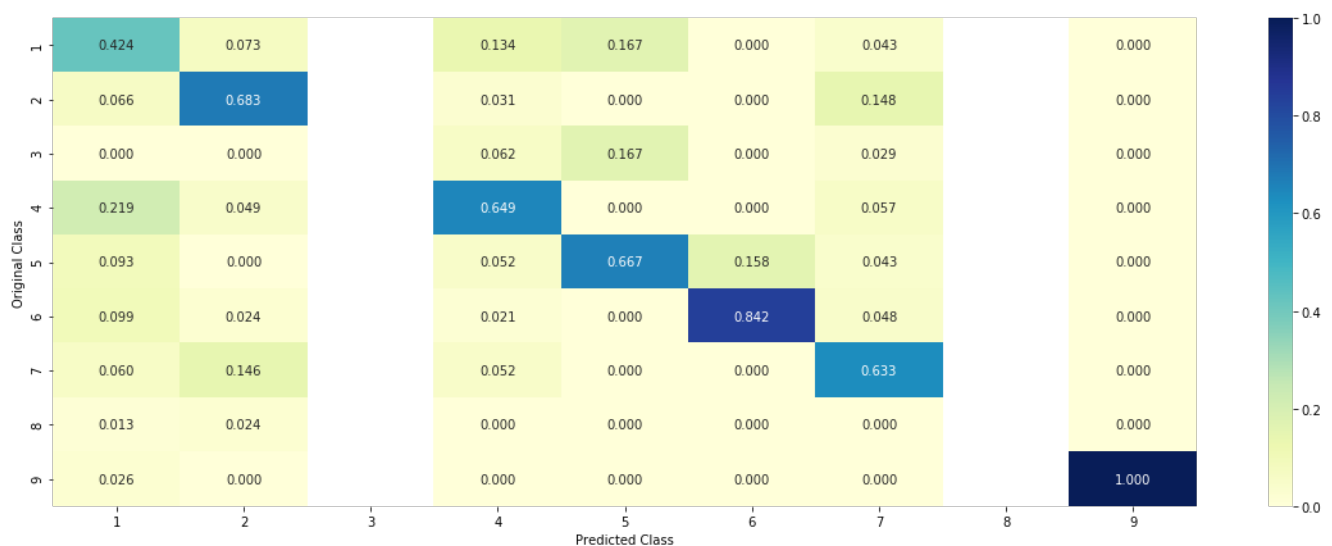
Log loss : 1.1905755152282853

Number of mis-classified points : 0.40977443609022557

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



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





4.5.3. Feature Importance

4.5.3.1. Correctly Classified point

In [94]:

```
# 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)
clf.fit(train_x_onehotCoding, train_y)
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 : 7

Predicted Class Probabilities: [[0.1452 0.2981 0.0135 0.0622 0.0485 0.0483 0.3476 0.0158 0.0208]]

Actual Class : 7

```
-----
11 Text feature [101] present in test data point [True]
18 Text feature [02] present in test data point [True]
33 Text feature [027] present in test data point [True]
37 Text feature [042] present in test data point [True]
61 Text feature [1000] present in test data point [True]
65 Text feature [12p13] present in test data point [True]
69 Text feature [009] present in test data point [True]
76 Text feature [10p11] present in test data point [True]
Out of the top 100 features 8 are present in query point
```

4.5.3.2. Inorrectly Classified point

In [95]:

```
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:",
```

```

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 : 7

Predicted Class Probabilities: [[0.0981 0.1202 0.0264 0.1247 0.0591 0.0891 0.4658 0.0102 0.0063]]

Actual Class : 7

61 Text feature [1000] present in test data point [True]

Out of the top 100 features 1 are present in query point

4.5.3. Hyper paramter tuning (With Response Coding)

In [96]:

```

# -----
# 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).

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

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

alpha = [10,50,100,200,500,1000]
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))
'''
# fit on the training data

```

```
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 : 2.060498452695094
for n_estimators = 10 and max depth = 3
Log Loss : 1.606713365024743
for n_estimators = 10 and max depth = 5
Log Loss : 1.4192388665172522
for n_estimators = 10 and max depth = 10
Log Loss : 2.269187615849894
for n_estimators = 50 and max depth = 2
Log Loss : 1.738662865487548
for n_estimators = 50 and max depth = 3
Log Loss : 1.4488530823726253
for n_estimators = 50 and max depth = 5
Log Loss : 1.4300935302380497
for n_estimators = 50 and max depth = 10
Log Loss : 1.7953191181370098
for n_estimators = 100 and max depth = 2
Log Loss : 1.642546723777754
for n_estimators = 100 and max depth = 3
Log Loss : 1.492862864438009
for n_estimators = 100 and max depth = 5
Log Loss : 1.334934089310618
for n_estimators = 100 and max depth = 10
Log Loss : 1.6969616634661162
for n_estimators = 200 and max depth = 2
Log Loss : 1.6778248867638328
for n_estimators = 200 and max depth = 3
Log Loss : 1.5422340427551107
for n_estimators = 200 and max depth = 5
Log Loss : 1.4052518393442586
for n_estimators = 200 and max depth = 10
Log Loss : 1.7204677071131564
for n_estimators = 500 and max depth = 2
Log Loss : 1.7313171678470458
for n_estimators = 500 and max depth = 3
Log Loss : 1.5765255159743978
for n_estimators = 500 and max depth = 5
Log Loss : 1.4347988455897978
for n_estimators = 500 and max depth = 10
Log Loss : 1.7608676086666228
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6775257280830034
for n_estimators = 1000 and max depth = 3
Log Loss : 1.5679259778770482
for n_estimators = 1000 and max depth = 5
Log Loss : 1.399787740698381
for n_estimators = 1000 and max depth = 10
```



```

Log Loss : 1.7333956686612713
For values of best alpha = 100 The train log loss is: 0.05930504051171449
For values of best alpha = 100 The cross validation log loss is: 1.334934089310618
For values of best alpha = 100 The test log loss is: 1.2885892520021862

```

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

In [97]:

```

# -----
# 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).

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-fores
t-and-their-construction-2/
# -----

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)

```

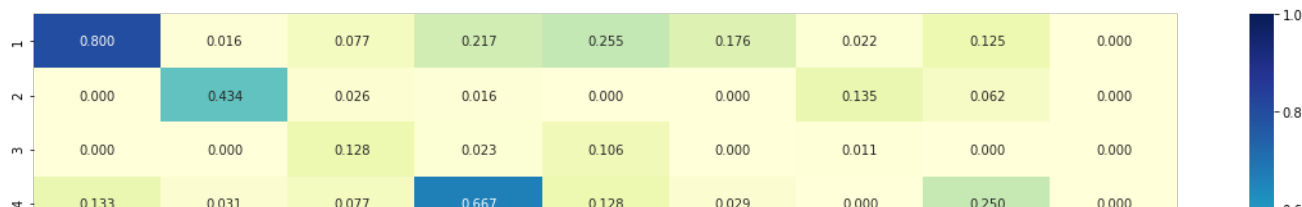
```

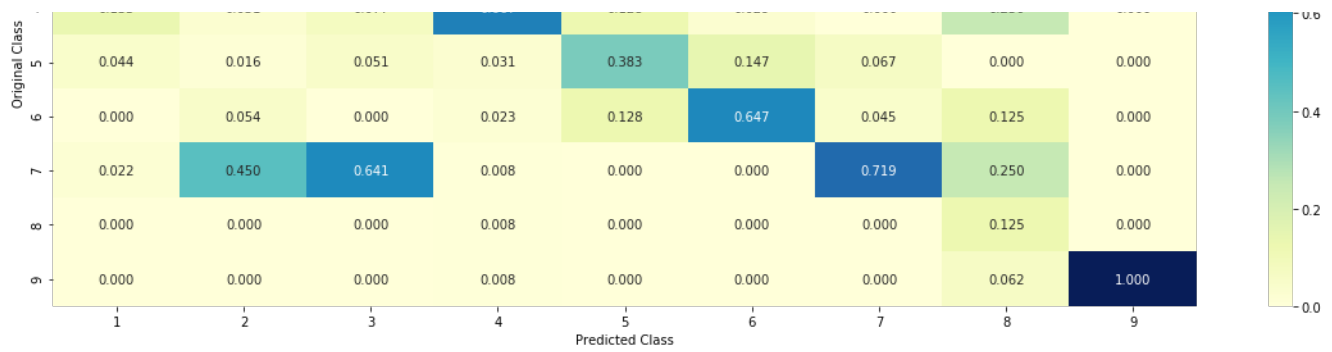
Log loss : 1.334934089310618
Number of mis-classified points : 0.4492481203007519
----- Confusion matrix -----

```

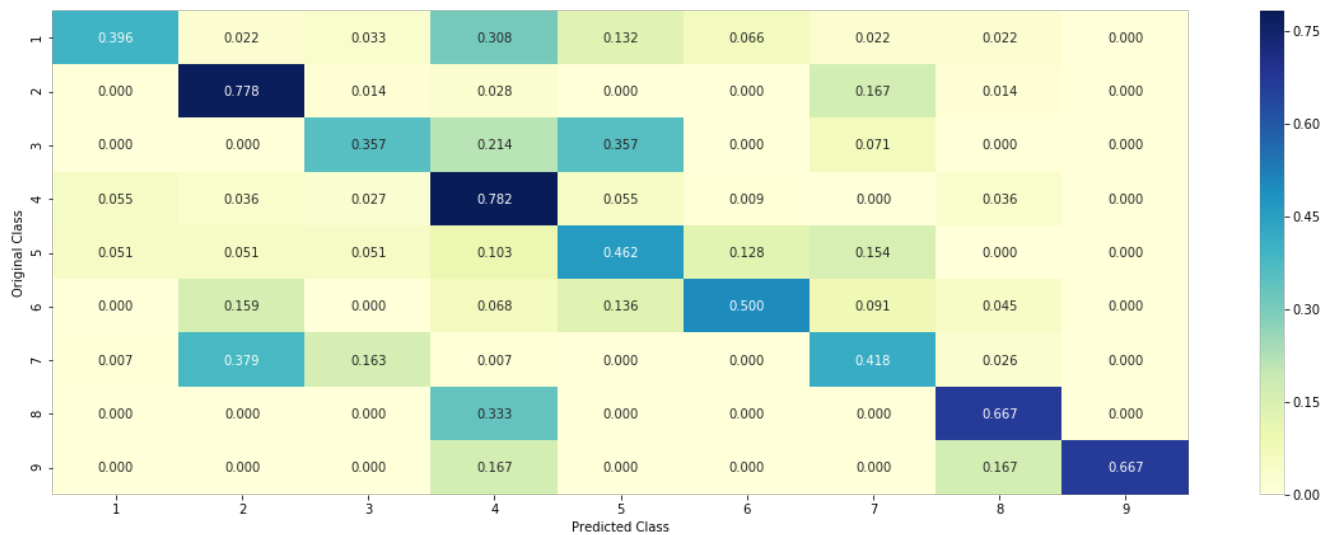


----- Precision matrix (Column Sum=1) -----





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



4.5.5. Feature Importance

4.5.5.1. Correctly Classified point

In [98]:

```
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 : 7

Predicted Class Probabilities: [[0.0017 0.0652 0.0015 0.0029 0.0012 0.0174 0.9067 0.0015 0.0018]]

Actual Class : 7

Variation is important feature

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

4.5.5.2. Incorrectly Classified point

In [99]:

```
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 : 6

Predicted Class Probabilities: [[0.0272 0.1406 0.1562 0.051 0.037 0.3378 0.1737 0.0574 0.0191]]

Actual Class : 7

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

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

4.7 Stack the models

4.7.1 testing with hyper parameter tuning

In [100]:

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

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

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

# read more about support vector machines with linear kernels 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.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

# read more about support vector machines with linear kernels here http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html
# -----
# 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).

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
```

```
# -----

clf1 = SGDClassifier(alpha=0.001, 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=1, 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_proba=True)
    sclf.fit(train_x_onehotCoding, train_y)
    print("Stacking Classifier : 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: 1.05
Support vector machines : Log Loss: 1.81
Naive Bayes : Log Loss: 1.18
-----
```

```
Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.177
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.031
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.495
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.157
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.362
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.787
```

4.7.2 testing the model with the best hyper parameters

In [101]:

```
lr = LogisticRegression(C=0.1)
sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_proba=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)

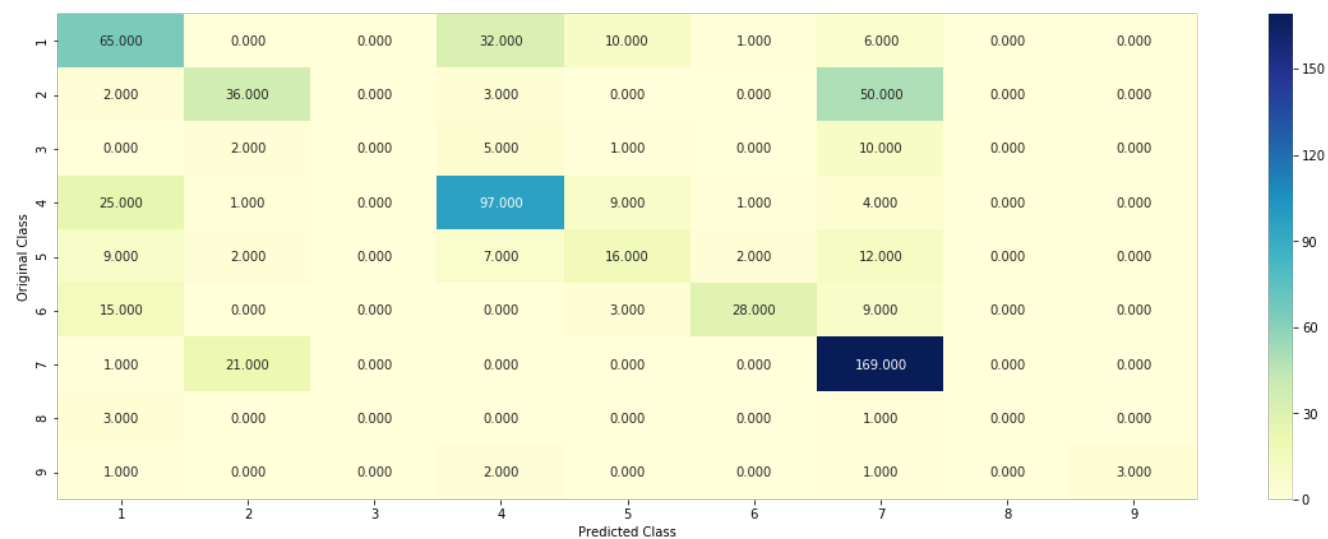
log_error = log_loss(test_y, sclf.predict_proba(test_x_onehotCoding))
print("Log loss (test) on the stacking classifier :", log_error)

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))
```

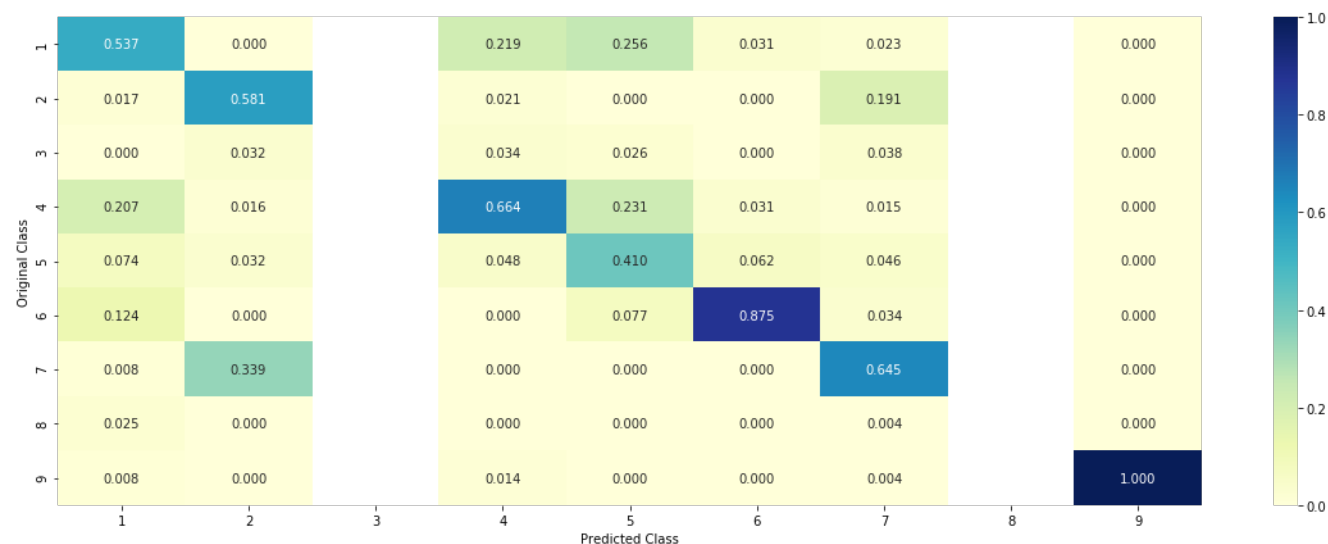
```
Log loss (train) on the stacking classifier : 0.5358049354834701
```

Log loss (train) on the stacking classifier : 0.3338049334834701
Log loss (CV) on the stacking classifier : 1.1570599523748386
Log loss (test) on the stacking classifier : 1.1577868259031538
Number of missclassified point : 0.3774436090225564

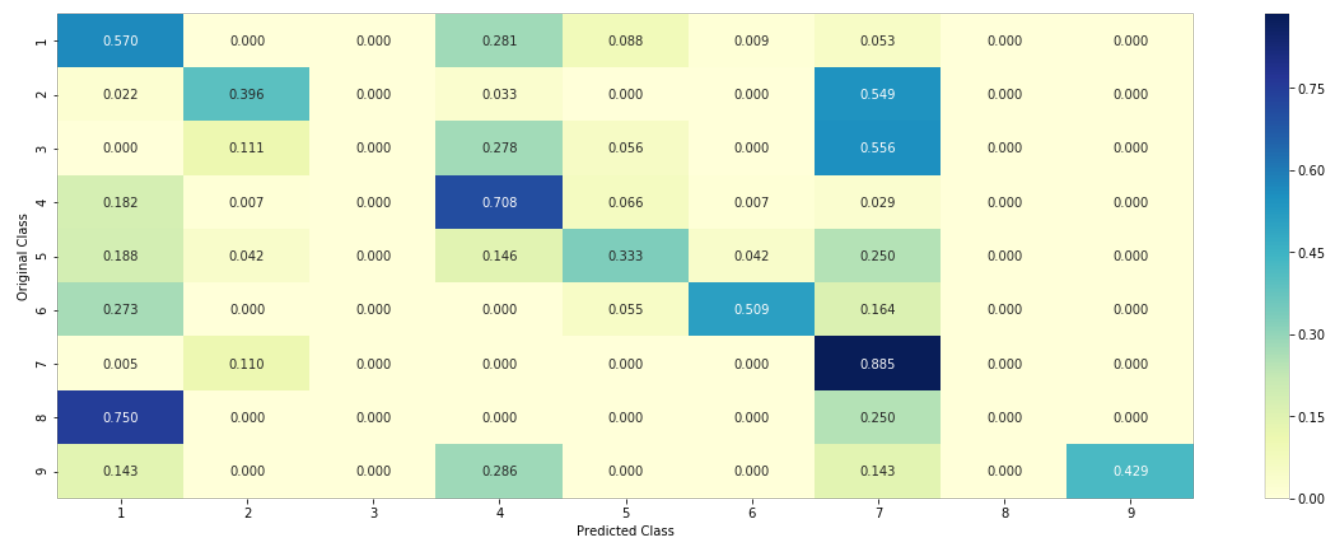
----- Confusion matrix -----



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



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



4.7.3 Maximum Voting classifier

In [102]:

```
#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))
```

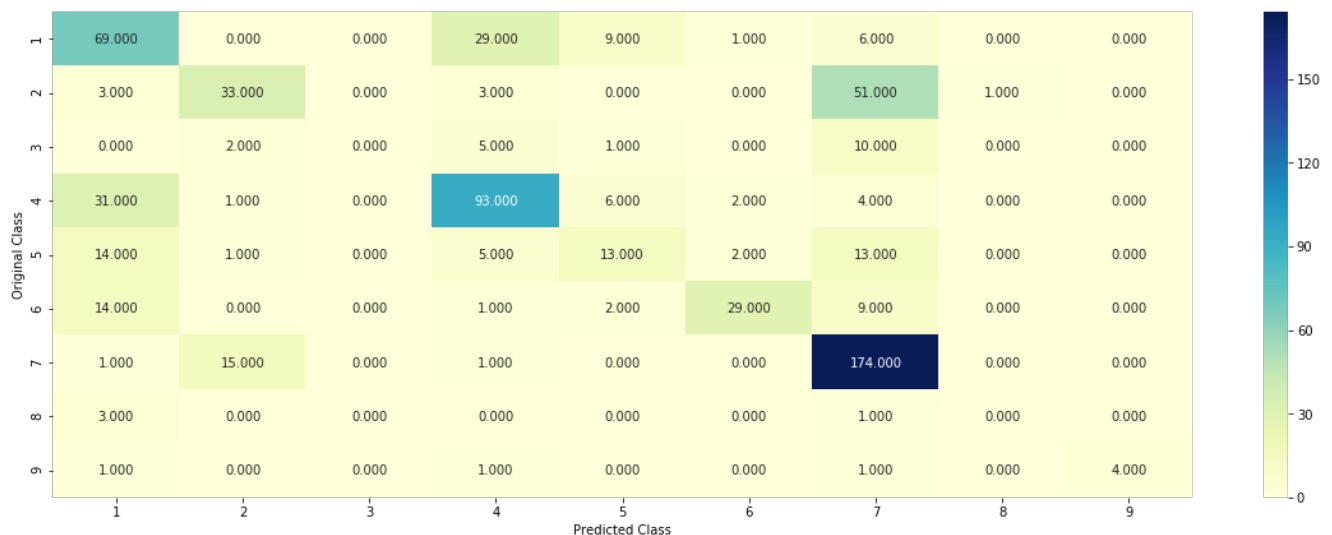
Log loss (train) on the VotingClassifier : 0.8350258772850658

Log loss (CV) on the VotingClassifier : 1.1969669416516497

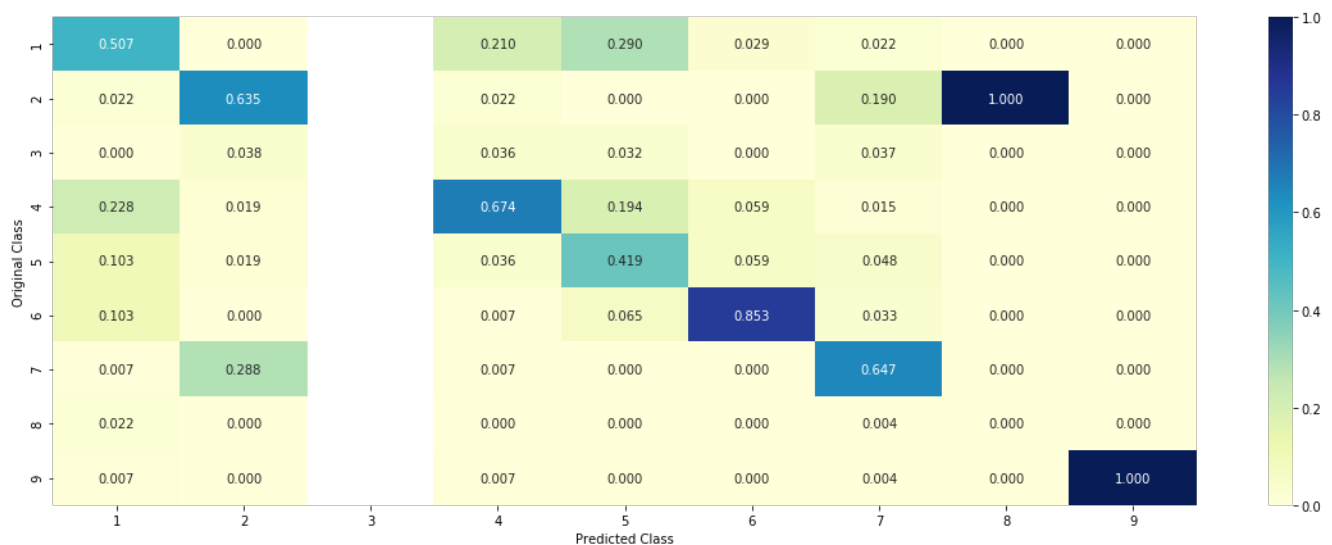
Log loss (test) on the VotingClassifier : 1.1862461699273428

Number of missclassified point : 0.37593984962406013

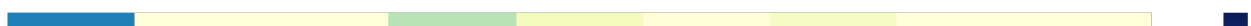
----- Confusion matrix -----

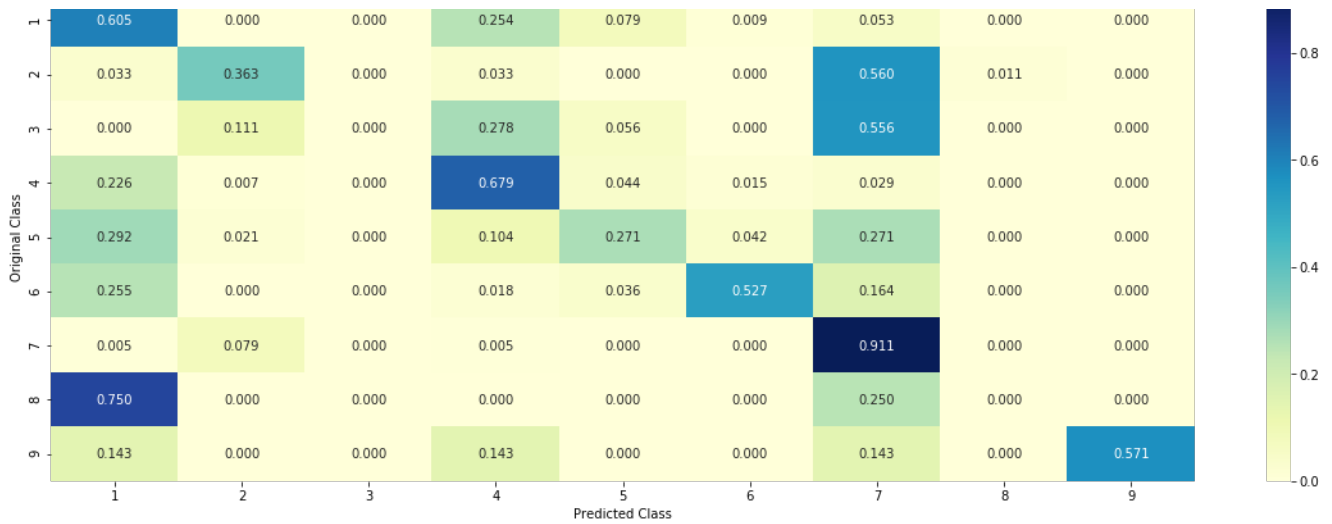


----- Precision matrix (Column Sum=1) -----



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





5. Conclusion

5.1 Steps taken

Instead of using all the words in the dataset, use only the top 1000 words based of tf-idf values

5.2 Model Comparison

In [3]:

```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
import pandas as pd
models = pd.DataFrame({'Model': ['Naive Bayes', "K-NN", "Logistic Regression with Class Balancing", 'Logistic Regression without Class Balancing', 'Linear SVM', 'Random forest', 'Random forest(response encoding)', 'Stacking classifier', 'Maximum voting classifier'], 'Log-Loss': [1.18, 1.04, 1.012, 1.05, 1.03, 1.19, 1.33, 1.15, 1.19], 'Percent of Misclassified points': [.379, .345, 0.338, .342, .330, .409, .449, .377, .375], 'Hyperparameter': [.001, 5, .0001, .0001, .0001, 2000, 100, .1, 0]}, columns = ["Model", "Log-Loss", "Percent of Misclassified points", "Hyperparameter"])
models
```

Out[3]:

	Model	Log-Loss	Percent of Misclassified points	Hyperparameter
0	Naive Bayes	1.180	0.379	0.0010
1	K-NN	1.040	0.345	5.0000
2	Logistic Regression with Class Balancing	1.012	0.338	0.0001
3	Logistic Regression without Class Balancing	1.050	0.342	0.0001
4	Linear SVM	1.030	0.330	0.0001
5	Random forest	1.190	0.409	2000.0000
6	Random forest(response encoding)	1.330	0.449	100.0000
7	Stacking classifier	1.150	0.377	0.1000
8	Maximum voting classifier	1.190	0.375	0.0000