

Personalized cancer diagnosis

1. Business Problem

1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

Context:

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

Problem statement :

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

1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

1. <https://www.forbes.com/sites/matthewherper/2017/06/03/a-new-cancer-drug-helped-almost-everyone-who-took-it-almost-heres-what-it-teaches-us/#2a44ee2f6b25>
2. <https://www.youtube.com/watch?v=UwbuW7oK8rk>
3. <https://www.youtube.com/watch?v=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'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
import seaborn as sns
from collections import Counter, defaultdict
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn.manifold import TSNE
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix, normalized_mutual_info_score
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.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB, GaussianNB
from sklearn.model_selection import train_test_split, GridSearchCV, StratifiedKFold

import math
from sklearn.ensemble import RandomForestClassifier
warnings.filterwarnings("ignore")

from mlxtend.classifier import StackingClassifier

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
import os
os.chdir('C:/Users/kingsubham27091995/Desktop/AppliedAiCouse/CASE STUDIES/Personalised Cancer Diag
nosis')
```

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

In [2]:

```
data_variants = pd.read_csv('training/training_variants')
print('Number of data points : ', data_variants.shape[0])
print('Number of features : ', data_variants.shape[1])
print('Features : ', data_variants.columns.values)
data_variants.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("training/training_text", sep="\\|\\|", engine="python", names=["ID", "TEXT"], skip
rows=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 : 218.35165418842536 seconds

```

In [6]:

```

# Merging both gene_variations and text data based on ID
result = pd.merge(data_variants, 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

If anyText feature has NAN values, replace it with 'Gene Variation' and treat it as Text

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]:

```
result.Gene = result.Gene.str.replace('\s+', '_')
result.Variation = result.Variation.str.replace('\s+', '_')
y_true = result[['Class']]
x_true = result.drop(['Class'], axis=1)

print("Feature columns in dataset: ")
print(x_true.head())
print()
print("Target columns in dataset: ")
print(y_true.head())
```

Feature columns in dataset:

```
   ID  Gene  Variation \
0   0  FAM58A  Truncating_Mutations
1   1   CBL          W802*
2   2   CBL          Q249E
3   3   CBL          N454D
4   4   CBL          L399V
```

TEXT

```
0  cyclin dependent kinases cdks regulate variety...
1  abstract background non small cell lung cancer...
2  abstract background non small cell lung cancer...
3  recent evidence demonstrated acquired uniparen...
4  oncogenic mutations monomeric casitas b lineag...
```

Target columns in dataset:

```
   Class
0      1
1      2
2      2
3      3
4      4
```

In [11]:

```
# Split the data into test and train by maintaining same distribution of output variable 'y_true'
[stratify=y_true]
x_train, x_test, y_train, y_test = train_test_split(x_true, 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]
x_train, x_cv, 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 [12]:

```
print('Number of data points in train data:', x_train.shape[0])
```

```
print('Number of data points in test data:', x_test.shape[0])
print('Number of data points in cross validation data:', x_cv.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 [13]:

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

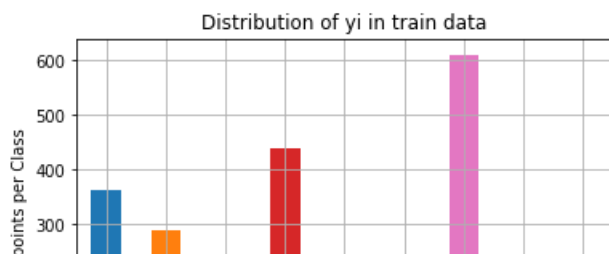
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

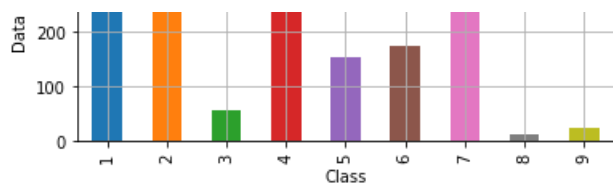
print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of y_i in test data')
plt.grid()
plt.show()

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

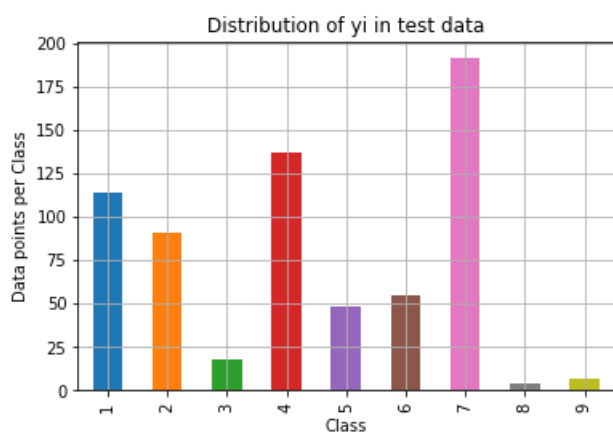
print('-'*80)
my_colors = 'rgbkymc'
cv_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of y_i in cross validation data')
plt.grid()
plt.show()

# 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]/y_cv.shape[0]*100), 3), '%)')
```

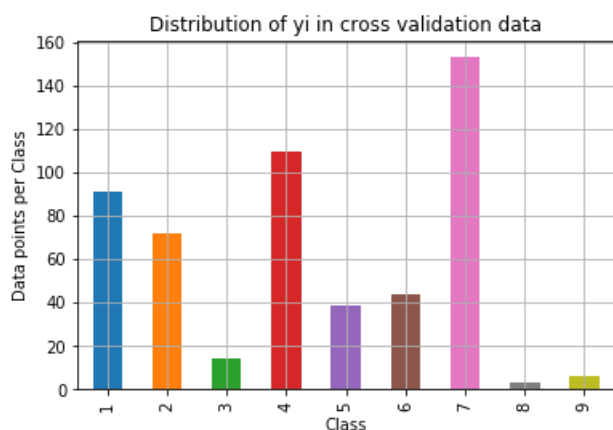




Number of data points in class 7 : 609 (28.672 %)
 Number of data points in class 4 : 439 (20.669 %)
 Number of data points in class 1 : 363 (17.09 %)
 Number of data points in class 2 : 289 (13.606 %)
 Number of data points in class 6 : 176 (8.286 %)
 Number of data points in class 5 : 155 (7.298 %)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



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



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

Summary from Histogram:

1. Imbalanced Data
2. Classes 1,2,4,7 are dominant classes .
3. Distribution of y_i's in Train,Test,CV Data are approximately same.

3.2 Prediction using a 'Random' Model

In a 'Random' Model, we generate the '9' class probabilities randomly such that they sum to 1.

In [14]:

```
def plot_matrix(matrix, labels):
    plt.figure(figsize=(20,7))
    sns.heatmap(matrix, annot=True, cmap="YlGnBu", fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()

# This function plots the confusion matrices given y_i, y_i_hat.
def plot_confusion_matrix(test_y, predict_y):
    cm = 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

    recall_table = ((cm.T)/(cm.sum(axis=1))).T
    # How did we calculate recall_table :
    # divide each element of the confusion matrix with the sum of elements in that column
    # C = [[1, 2],
    #      [3, 4]]
    # C.T = [[1, 3],
    #        [2, 4]]
    # C.sum(axis = 1) axis=0 corresponds to columns and axis=1 corresponds to rows in two dimensional array
    # C.sum(axis = 1) = [[3, 7]]
    # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7],
    #                             [2/3, 4/7]]
    # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3],
    #                               [3/7, 4/7]]
    # sum of row elements = 1

    precision_table = (cm/cm.sum(axis=0))
    # How did we calculate precision table :
    # divide 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]
    print()
    print("-"*20, "Confusion matrix", "-"*20)
    plot_matrix(cm, labels)

    print("-"*20, "Precision matrix (Column Sum=1)", "-"*20)
    plot_matrix(precision_table, labels)

    print("-"*20, "Recall matrix (Row sum=1)", "-"*20)
    plot_matrix(recall_table, labels)
```

In [15]:

```
# 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 = x_test.shape[0]
cv_data_len = x_cv.shape[0]
```


'kmt2a',
'kmt2b',
'kmt2c',
'kmt2d',
'knstrn',
'kras',
'lats1',
'map2k1',
'map2k2',
'map2k4',
'map3k1',
'mapk1',
'mdm2',
'mdm4',
'med12',
'mef2b',
'met',
'mga',
'mlh1',
'mpl',
'msh2',
'msh6',
'mtor',
'myc',
'mycn',
'myd88',
'myod1',
'nfl',
'nf2',
'nfe2l2',
'nfkb1a',
'nkx2',
'notch1',
'notch2',
'npm1',
'nras',
'nsd1',
'ntrk1',
'ntrk2',
'ntrk3',
'nup93',
'pak1',
'pbrm1',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pik3r3',
'pim1',
'pms2',
'pole',
'ppm1d',
'ppp2r1a',
'ppp6c',
'prdm1',
'ptch1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rab35',
'rac1',
'rad21',
'rad50',
'rad51b',
'rad51c',
'rad51d',
'rad54l',
'raf1',
'rara',
'rasa1',
'rb1',
'rbm10',
'ret'.

```

'rcu',
'rheb',
'rhoa',
'rit1',
'rnf43',
'ros1',
'rras2',
'runx1',
'rxra',
'rybp',
'sdhb',
'sdhc',
'setd2',
'sf3b1',
'shoc2',
'shq1',
'smad2',
'smad3',
'smad4',
'smarca4',
'smarcb1',
'smo',
'sos1',
'sox9',
'spop',
'src',
'stag2',
'stat3',
'stk11',
'tcf3',
'tcf7l2',
'tert',
'tet1',
'tet2',
'tgfbr1',
'tgfbr2',
'tmprss2',
'tp53',
'tp53bp1',
'tsc1',
'tsc2',
'u2af1',
'vegfa',
'vhl',
'whsc1',
'xpo1',
'xrcc2',
'yap1']

```

In [26]:

```

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

```

train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 240)

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 [27]:

```

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=0.1)

```

```

# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate=optimal, eta
=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.clas
ses_, 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))

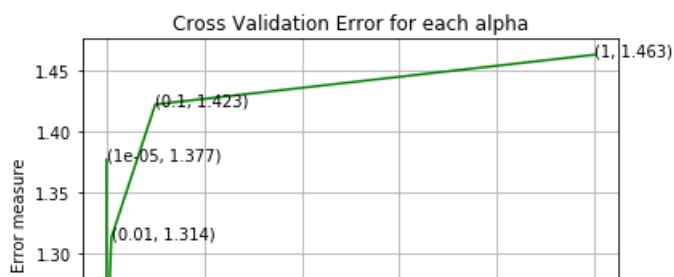
predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha],
      "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

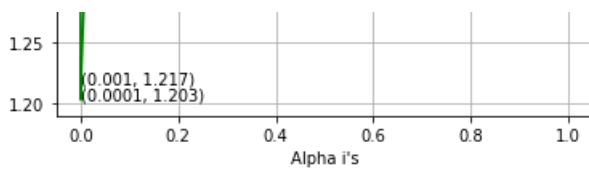
```

```

For values of alpha = 1e-05 The log loss is: 1.3774056725736534
For values of alpha = 0.0001 The log loss is: 1.203391638750355
For values of alpha = 0.001 The log loss is: 1.216736551413754
For values of alpha = 0.01 The log loss is: 1.313567738348165
For values of alpha = 0.1 The log loss is: 1.4227072523958515
For values of alpha = 1 The log loss is: 1.4634435086693665

```





For values of best alpha = 0.0001 The train log loss is: 1.0486598182030518
 For values of best alpha = 0.0001 The cross validation log loss is: 1.203391638750355
 For values of best alpha = 0.0001 The test log loss is: 1.1834720946379564

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. Here the Train log-loss, Test log-loss and CV log-loss are not very different

If the gap between Train,Test and CV log-loss would be bigger then it would mean that you are Overfitting

This checks Out of Total Data, how much part is Overlapping

In [28]:

```
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=x_test[x_test['Gene'].isin(list(set(x_train['Gene'])))].shape[0]
cv_coverage=x_cv[x_cv['Gene'].isin(list(set(x_train['Gene'])))].shape[0]

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

Q6. How many data points in Test and CV datasets are covered by the 240 genes in train dataset?
 Ans
 1. In test data 649 out of 665 : 97.59398496240601
 2. In cross validation data 520 out of 532 : 97.74436090225564

Observation:

Out of 665 data points in Test Data, 649 are present in Train Data

Out of 532 data points in CV Data, 520 are present in Train Data

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 [29]:

```
unique_variations = x_train['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 : 1919
Truncating_Mutations      68
Amplification              46
Deletion                  45
Fusions                   20
Overexpression             5
G12V                      3
```

```
Q61H          3
G12S          2
Q61R          2
G13D          2
Name: Variation, dtype: int64
```

In [30]:

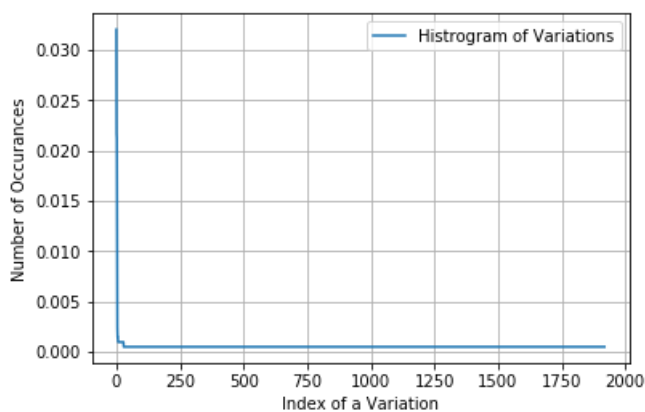
```
print("Ans: There are", unique_variations.shape[0] ,
      "different categories of variations in the train data, and they are distributed as follows",)
```

Ans: There are 1919 different categories of variations in the train data, and they are distributed as follows

Plotting the Distribution

In [31]:

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

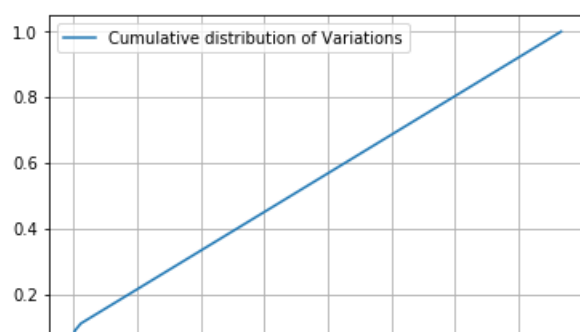


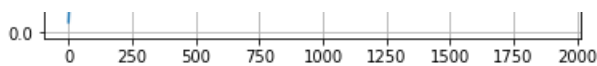
Plotting CDF

In [32]:

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

```
[0.03201507 0.05367232 0.07485876 ... 0.99905838 0.99952919 1.          ]
```





Q9. How to featurize this Variation feature ?

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

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

1. One hot Encoding
2. Response coding

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

Response Coding

In [33]:

```
# alpha is used for laplace smoothing
alpha = 1

# train gene feature
train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_train))

# test gene feature
test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_test))

# cross validation gene feature
cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_cv))
```

In [34]:

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

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

One Hot Encoding

In [35]:

```
# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer()
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(x_train['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(x_test['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(x_cv['Variation'])
```

In [36]:

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

train_variation_feature_onehotEncoded is converted feature using the one-hot encoding method. The shape of Variation feature: (2124, 1952)

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

Let's build a model just like the earlier!

In [37]:

```
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)

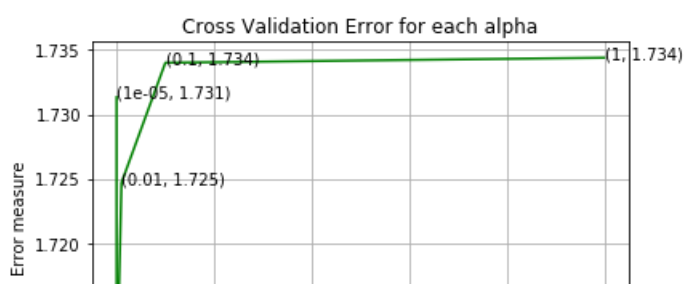
predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
print('For values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

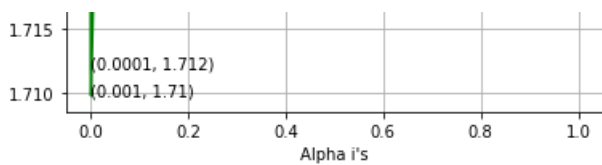
```

```

For values of alpha = 1e-05 The log loss is: 1.731326780723876
For values of alpha = 0.0001 The log loss is: 1.7119671791483242
For values of alpha = 0.001 The log loss is: 1.7097653732990405
For values of alpha = 0.01 The log loss is: 1.7246852812525035
For values of alpha = 0.1 The log loss is: 1.734004742967627
For values of alpha = 1 The log loss is: 1.7343947792547312

```





For values of best alpha = 0.001 The train log loss is: 1.061022199942188
 For values of best alpha = 0.001 The cross validation log loss is: 1.7097653732990405
 For values of best alpha = 0.001 The test log loss is: 1.7204042838973364

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.

This checks Out of Total Data, how much part is Overlapping

In [38]:

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

Q12. How many data points are covered by total 1919 genes in test and cross validation data sets?

Ans

1. In test data 55 out of 665 : 8.270676691729323
2. In cross validation data 60 out of 532 : 11.278195488721805

Observation:

Out of 665 data points in Test Data, 55 are present in Train Data i.e. only 8%(approx)

Out of 532 data points in CV Data, 60 are present in Train Data i.e. only 11%(approx)

- Overlapping is very less. So, this feature is fairly unstable. But, the log-losses are significantly lower than that of Random Model. So, even though this Variation feature is fairly unstable it certainly adds value due to its significant drop in its log-loss. So, keep Variation feature too.

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 [65]:

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

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


In [70]:

```
# 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
train_text_fea_counts
```

Out[70]:

```
array([8.69443322, 8.97552233, 0.03666568, ..., 0.02557271, 0.02778015,
       0.05733632])
```

In [71]:

```
# 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 [72]:

```
dict_list = []
# dict_list =[] contains 9 dictionaries each corresponds to a class
for i in range(1,10):
    cls_text = x_train[y_train['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(x_train)

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 [73]:

```
#response coding of text features
train_text_feature_responseCoding = get_text_responsecoding(x_train)
test_text_feature_responseCoding = get_text_responsecoding(x_test)
cv_text_feature_responseCoding = get_text_responsecoding(x_cv)
```

In [74]:

```
# 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
```

It's always a good practice to Normalize the data after One Hot Encoding

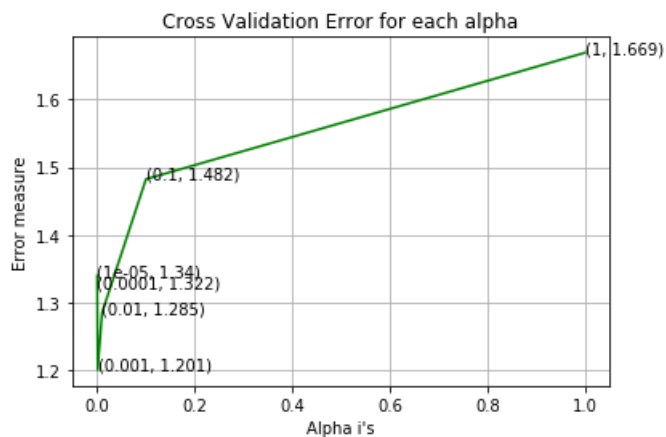
In [75]:

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



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

For values of alpha = 1e-05 The log loss is: 1.3395646099640204
 For values of alpha = 0.0001 The log loss is: 1.3223418174601882
 For values of alpha = 0.001 The log loss is: 1.2006670469026821
 For values of alpha = 0.01 The log loss is: 1.2851733627851996
 For values of alpha = 0.1 The log loss is: 1.4821758596181456
 For values of alpha = 1 The log loss is: 1.668777678490356



For values of best alpha = 0.001 The train log loss is: 0.6642991339986491
 For values of best alpha = 0.001 The cross validation log loss is: 1.2006670469026821
 For values of best alpha = 0.001 The test log loss is: 1.1323377955939227

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

Ans. Yes, it seems like!

In [79]:

```
def get_intersec_text(df):
    df_text_vec = TfidfVectorizer(min_df=3,max_features=1000)
    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 [80]:

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

23.4 % of word of test data appeared in train data
 22.8 % of word of Cross Validation appeared in train data

4. Machine Learning Models

In [81]:

```
#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 [82]:

```
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 [83]:

```
# 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 = TfidfVectorizer()
    var_count_vec = TfidfVectorizer()
    text_count_vec = TfidfVectorizer(min_df=3)

    gene_vec = gene_count_vec.fit(x_train['Gene'])
    var_vec = var_count_vec.fit(x_train['Variation'])
    text_vec = text_count_vec.fit(x_train['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 [84]:

```
# 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(y_train['Class']))

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

cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(y_cv['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 [85]:

```

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

```

```

One hot encoding features :
(number of data points * number of features) in train data = (2124, 54097)
(number of data points * number of features) in test data = (665, 54097)
(number of data points * number of features) in cross validation data = (532, 54097)

```

In [86]:

```

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

In [87]:

```
# 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 probabilities 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_v = sig_clf.predict_proba(test_x_onehotCoding)
```


349 Text feature [100] present in test data point [True]
350 Text feature [clinically] present in test data point [True]
352 Text feature [caused] present in test data point [True]
354 Text feature [confer] present in test data point [True]
356 Text feature [introduction] present in test data point [True]
357 Text feature [sensitivity] present in test data point [True]
358 Text feature [incubated] present in test data point [True]
360 Text feature [significant] present in test data point [True]
361 Text feature [therapy] present in test data point [True]
362 Text feature [mapk] present in test data point [True]
363 Text feature [stimulation] present in test data point [True]
364 Text feature [containing] present in test data point [True]
365 Text feature [mutated] present in test data point [True]
367 Text feature [doses] present in test data point [True]
369 Text feature [comparable] present in test data point [True]
371 Text feature [16] present in test data point [True]
372 Text feature [identify] present in test data point [True]
373 Text feature [disease] present in test data point [True]
374 Text feature [ml] present in test data point [True]
376 Text feature [resistant] present in test data point [True]
377 Text feature [fact] present in test data point [True]
378 Text feature [differences] present in test data point [True]
380 Text feature [greater] present in test data point [True]
382 Text feature [cancer] present in test data point [True]
386 Text feature [signal] present in test data point [True]
390 Text feature [maintained] present in test data point [True]
394 Text feature [regulation] present in test data point [True]
395 Text feature [primary] present in test data point [True]
396 Text feature [strongly] present in test data point [True]
397 Text feature [whole] present in test data point [True]
398 Text feature [seen] present in test data point [True]
399 Text feature [manner] present in test data point [True]
400 Text feature [understanding] present in test data point [True]
401 Text feature [manufacturer] present in test data point [True]
402 Text feature [50] present in test data point [True]
403 Text feature [therefore] present in test data point [True]
406 Text feature [malignant] present in test data point [True]
407 Text feature [transfected] present in test data point [True]
408 Text feature [normal] present in test data point [True]
410 Text feature [independently] present in test data point [True]
412 Text feature [importance] present in test data point [True]
413 Text feature [whose] present in test data point [True]
414 Text feature [17] present in test data point [True]
415 Text feature [increases] present in test data point [True]
417 Text feature [agents] present in test data point [True]
418 Text feature [demonstrate] present in test data point [True]
420 Text feature [investigated] present in test data point [True]
421 Text feature [fold] present in test data point [True]
422 Text feature [entire] present in test data point [True]
424 Text feature [open] present in test data point [True]
427 Text feature [viability] present in test data point [True]
428 Text feature [type] present in test data point [True]
429 Text feature [developed] present in test data point [True]
430 Text feature [4d] present in test data point [True]
431 Text feature [amplification] present in test data point [True]
432 Text feature [prepared] present in test data point [True]
433 Text feature [3d] present in test data point [True]
435 Text feature [threonine] present in test data point [True]
436 Text feature [22] present in test data point [True]
437 Text feature [detect] present in test data point [True]
439 Text feature [associated] present in test data point [True]
441 Text feature [selective] present in test data point [True]
442 Text feature [possible] present in test data point [True]
444 Text feature [biological] present in test data point [True]
446 Text feature [needed] present in test data point [True]
447 Text feature [erk] present in test data point [True]
448 Text feature [decrease] present in test data point [True]
450 Text feature [4c] present in test data point [True]
451 Text feature [line] present in test data point [True]
452 Text feature [mouse] present in test data point [True]
454 Text feature [form] present in test data point [True]
455 Text feature [generation] present in test data point [True]
456 Text feature [able] present in test data point [True]
457 Text feature [green] present in test data point [True]
458 Text feature [response] present in test data point [True]
460 Text feature [pcr] present in test data point [True]
461 Text feature [predominantly] present in test data point [True]

462 Text feature [highest] present in test data point [True]
463 Text feature [sl] present in test data point [True]
464 Text feature [sample] present in test data point [True]
465 Text feature [event] present in test data point [True]
466 Text feature [subjected] present in test data point [True]
467 Text feature [complete] present in test data point [True]
469 Text feature [collected] present in test data point [True]
470 Text feature [molecule] present in test data point [True]
472 Text feature [represented] present in test data point [True]
473 Text feature [11] present in test data point [True]
475 Text feature [would] present in test data point [True]
476 Text feature [example] present in test data point [True]
478 Text feature [plays] present in test data point [True]
480 Text feature [stable] present in test data point [True]
481 Text feature [center] present in test data point [True]
482 Text feature [used] present in test data point [True]
483 Text feature [summary] present in test data point [True]
484 Text feature [control] present in test data point [True]
485 Text feature [long] present in test data point [True]
486 Text feature [low] present in test data point [True]
488 Text feature [explain] present in test data point [True]
489 Text feature [virus] present in test data point [True]
490 Text feature [find] present in test data point [True]
491 Text feature [initially] present in test data point [True]
492 Text feature [observation] present in test data point [True]
497 Text feature [rate] present in test data point [True]
498 Text feature [reduced] present in test data point [True]
499 Text feature [five] present in test data point [True]
503 Text feature [provide] present in test data point [True]
504 Text feature [proteins] present in test data point [True]
505 Text feature [completely] present in test data point [True]
506 Text feature [notably] present in test data point [True]
508 Text feature [located] present in test data point [True]
509 Text feature [transforming] present in test data point [True]
510 Text feature [forms] present in test data point [True]
511 Text feature [coding] present in test data point [True]
514 Text feature [treating] present in test data point [True]
516 Text feature [context] present in test data point [True]
517 Text feature [investigate] present in test data point [True]
519 Text feature [wild] present in test data point [True]
520 Text feature [metastatic] present in test data point [True]
523 Text feature [board] present in test data point [True]
528 Text feature [data] present in test data point [True]
529 Text feature [frequent] present in test data point [True]
530 Text feature [number] present in test data point [True]
531 Text feature [multiple] present in test data point [True]
534 Text feature [stimulated] present in test data point [True]
535 Text feature [28] present in test data point [True]
536 Text feature [factors] present in test data point [True]
538 Text feature [necessary] present in test data point [True]
539 Text feature [account] present in test data point [True]
541 Text feature [animal] present in test data point [True]
544 Text feature [experimental] present in test data point [True]
545 Text feature [2d] present in test data point [True]
547 Text feature [often] present in test data point [True]
549 Text feature [alterations] present in test data point [True]
551 Text feature [since] present in test data point [True]
555 Text feature [nucleotide] present in test data point [True]
556 Text feature [involved] present in test data point [True]
559 Text feature [required] present in test data point [True]
560 Text feature [sufficient] present in test data point [True]
562 Text feature [stage] present in test data point [True]
563 Text feature [properties] present in test data point [True]
566 Text feature [cases] present in test data point [True]
567 Text feature [remained] present in test data point [True]
568 Text feature [date] present in test data point [True]
569 Text feature [shows] present in test data point [True]
570 Text feature [80] present in test data point [True]
573 Text feature [range] present in test data point [True]
578 Text feature [transfection] present in test data point [True]
580 Text feature [et] present in test data point [True]
581 Text feature [known] present in test data point [True]
584 Text feature [constructs] present in test data point [True]
585 Text feature [adenocarcinoma] present in test data point [True]
587 Text feature [early] present in test data point [True]
588 Text feature [upon] present in test data point [True]
589 Text feature [difference] present in test data point [True]

590 Text feature [appears] present in test data point [True]
599 Text feature [involving] present in test data point [True]
603 Text feature [numbers] present in test data point [True]
604 Text feature [al] present in test data point [True]
607 Text feature [acid] present in test data point [True]
608 Text feature [overall] present in test data point [True]
609 Text feature [third] present in test data point [True]
610 Text feature [resistance] present in test data point [True]
611 Text feature [harbored] present in test data point [True]
612 Text feature [subset] present in test data point [True]
613 Text feature [engineered] present in test data point [True]
615 Text feature [effect] present in test data point [True]
616 Text feature [red] present in test data point [True]
617 Text feature [bearing] present in test data point [True]
620 Text feature [23] present in test data point [True]
621 Text feature [gain] present in test data point [True]
622 Text feature [least] present in test data point [True]
627 Text feature [targeting] present in test data point [True]
628 Text feature [demonstrating] present in test data point [True]
629 Text feature [critical] present in test data point [True]
631 Text feature [33] present in test data point [True]
632 Text feature [protein] present in test data point [True]
633 Text feature [exhibited] present in test data point [True]
636 Text feature [toward] present in test data point [True]
637 Text feature [include] present in test data point [True]
638 Text feature [ca] present in test data point [True]
641 Text feature [membrane] present in test data point [True]
645 Text feature [21] present in test data point [True]
646 Text feature [briefly] present in test data point [True]
651 Text feature [status] present in test data point [True]
654 Text feature [separated] present in test data point [True]
655 Text feature [exclusive] present in test data point [True]
656 Text feature [genes] present in test data point [True]
658 Text feature [inhibit] present in test data point [True]
659 Text feature [carried] present in test data point [True]
660 Text feature [formation] present in test data point [True]
661 Text feature [isolated] present in test data point [True]
666 Text feature [stratagene] present in test data point [True]
668 Text feature [extent] present in test data point [True]
670 Text feature [sanger] present in test data point [True]
673 Text feature [90] present in test data point [True]
674 Text feature [ld] present in test data point [True]
683 Text feature [inhibitory] present in test data point [True]
686 Text feature [s4] present in test data point [True]
690 Text feature [particular] present in test data point [True]
692 Text feature [mutually] present in test data point [True]
693 Text feature [phosphorylated] present in test data point [True]
694 Text feature [region] present in test data point [True]
700 Text feature [administered] present in test data point [True]
704 Text feature [groups] present in test data point [True]
707 Text feature [millipore] present in test data point [True]
708 Text feature [expected] present in test data point [True]
709 Text feature [reports] present in test data point [True]
712 Text feature [median] present in test data point [True]
714 Text feature [buffer] present in test data point [True]
716 Text feature [vectors] present in test data point [True]
718 Text feature [many] present in test data point [True]
720 Text feature [model] present in test data point [True]
721 Text feature [right] present in test data point [True]
722 Text feature [vivo] present in test data point [True]
723 Text feature [identical] present in test data point [True]
724 Text feature [adding] present in test data point [True]
728 Text feature [transiently] present in test data point [True]
730 Text feature [majority] present in test data point [True]
731 Text feature [sites] present in test data point [True]
732 Text feature [aggressive] present in test data point [True]
733 Text feature [glutamine] present in test data point [True]
734 Text feature [limited] present in test data point [True]
736 Text feature [frequency] present in test data point [True]
738 Text feature [see] present in test data point [True]
742 Text feature [marked] present in test data point [True]
744 Text feature [corresponding] present in test data point [True]
745 Text feature [erk1] present in test data point [True]
752 Text feature [subsequent] present in test data point [True]
754 Text feature [six] present in test data point [True]
755 Text feature [correlate] present in test data point [True]
756 Text feature [spectrum] present in test data point [True]

757 Text feature [options] present in test data point [True]
761 Text feature [implications] present in test data point [True]
762 Text feature [schematic] present in test data point [True]
764 Text feature [tested] present in test data point [True]
766 Text feature [hypothesized] present in test data point [True]
770 Text feature [egfr] present in test data point [True]
773 Text feature [ic50] present in test data point [True]
775 Text feature [primers] present in test data point [True]
776 Text feature [full] present in test data point [True]
778 Text feature [roles] present in test data point [True]
779 Text feature [formed] present in test data point [True]
781 Text feature [wide] present in test data point [True]
783 Text feature [deletion] present in test data point [True]
784 Text feature [adenocarcinomas] present in test data point [True]
787 Text feature [noted] present in test data point [True]
789 Text feature [amino] present in test data point [True]
790 Text feature [solid] present in test data point [True]
791 Text feature [5a] present in test data point [True]
792 Text feature [use] present in test data point [True]
793 Text feature [staining] present in test data point [True]
794 Text feature [designed] present in test data point [True]
795 Text feature [regulate] present in test data point [True]
797 Text feature [screen] present in test data point [True]
801 Text feature [marker] present in test data point [True]
802 Text feature [36] present in test data point [True]
803 Text feature [ability] present in test data point [True]
809 Text feature [responsible] present in test data point [True]
810 Text feature [mice] present in test data point [True]
811 Text feature [biochemical] present in test data point [True]
813 Text feature [except] present in test data point [True]
814 Text feature [comparison] present in test data point [True]
818 Text feature [blue] present in test data point [True]
819 Text feature [induction] present in test data point [True]
820 Text feature [crystal] present in test data point [True]
824 Text feature [profile] present in test data point [True]
826 Text feature [effectors] present in test data point [True]
828 Text feature [44] present in test data point [True]
833 Text feature [reveals] present in test data point [True]
835 Text feature [density] present in test data point [True]
836 Text feature [6b] present in test data point [True]
837 Text feature [processed] present in test data point [True]
839 Text feature [exon] present in test data point [True]
841 Text feature [potentially] present in test data point [True]
842 Text feature [understand] present in test data point [True]
843 Text feature [experiment] present in test data point [True]
848 Text feature [better] present in test data point [True]
851 Text feature [essential] present in test data point [True]
852 Text feature [characterize] present in test data point [True]
853 Text feature [intrinsic] present in test data point [True]
857 Text feature [events] present in test data point [True]
858 Text feature [figures] present in test data point [True]
859 Text feature [share] present in test data point [True]
860 Text feature [neither] present in test data point [True]
863 Text feature [confers] present in test data point [True]
864 Text feature [closely] present in test data point [True]
865 Text feature [growing] present in test data point [True]
866 Text feature [components] present in test data point [True]
868 Text feature [left] present in test data point [True]
873 Text feature [sds] present in test data point [True]
875 Text feature [recruitment] present in test data point [True]
877 Text feature [inactive] present in test data point [True]
879 Text feature [exons] present in test data point [True]
881 Text feature [exhibit] present in test data point [True]
882 Text feature [5c] present in test data point [True]
884 Text feature [bars] present in test data point [True]
885 Text feature [mg] present in test data point [True]
890 Text feature [panel] present in test data point [True]
894 Text feature [clearly] present in test data point [True]
896 Text feature [best] present in test data point [True]
898 Text feature [smaller] present in test data point [True]
902 Text feature [29] present in test data point [True]
907 Text feature [pattern] present in test data point [True]
908 Text feature [probed] present in test data point [True]
912 Text feature [promising] present in test data point [True]
921 Text feature [37] present in test data point [True]
924 Text feature [genetic] present in test data point [True]
925 Text feature [residues] present in test data point [True]

```

929 Text feature [conclusion] present in test data point [True]
930 Text feature [outcome] present in test data point [True]
931 Text feature [glycine] present in test data point [True]
934 Text feature [basal] present in test data point [True]
936 Text feature [detectable] present in test data point [True]
940 Text feature [changed] present in test data point [True]
943 Text feature [methods] present in test data point [True]
944 Text feature [involve] present in test data point [True]
945 Text feature [view] present in test data point [True]
947 Text feature [larger] present in test data point [True]
949 Text feature [dose] present in test data point [True]
950 Text feature [specifically] present in test data point [True]
951 Text feature [listed] present in test data point [True]
954 Text feature [week] present in test data point [True]
955 Text feature [given] present in test data point [True]
956 Text feature [review] present in test data point [True]
959 Text feature [proximity] present in test data point [True]
960 Text feature [like] present in test data point [True]
961 Text feature [characterization] present in test data point [True]
965 Text feature [contributes] present in test data point [True]
966 Text feature [34] present in test data point [True]
967 Text feature [removed] present in test data point [True]
971 Text feature [nanomolar] present in test data point [True]
973 Text feature [matched] present in test data point [True]
980 Text feature [direct] present in test data point [True]
981 Text feature [case] present in test data point [True]
983 Text feature [provides] present in test data point [True]
985 Text feature [26] present in test data point [True]
987 Text feature [prognosis] present in test data point [True]
991 Text feature [potency] present in test data point [True]
993 Text feature [original] present in test data point [True]
995 Text feature [latter] present in test data point [True]
Out of the top 1000 features 609 are present in query point

```

4.1.1.4. Feature Importance, Correctly classified point

In [90]:

```

test_point_index = 55
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)

```

Predicted Class : 4

Predicted Class Probabilities: [[0.0177 0.2193 0.0179 0.6655 0.0297 0.022 0.0138 0.0065 0.0076]]

Actual Class : 4

```

-----
11 Text feature [proteins] present in test data point [True]
12 Text feature [protein] present in test data point [True]
13 Text feature [activity] present in test data point [True]
19 Text feature [function] present in test data point [True]
20 Text feature [mammalian] present in test data point [True]
23 Text feature [indicated] present in test data point [True]
24 Text feature [pten] present in test data point [True]
25 Text feature [results] present in test data point [True]
26 Text feature [described] present in test data point [True]
27 Text feature [whereas] present in test data point [True]
30 Text feature [shown] present in test data point [True]
31 Text feature [whether] present in test data point [True]
33 Text feature [loss] present in test data point [True]
34 Text feature [type] present in test data point [True]
35 Text feature [also] present in test data point [True]
39 Text feature [levels] present in test data point [True]
40 Text feature [expressed] present in test data point [True]

```

41 Text feature [two] present in test data point [True]
42 Text feature [functions] present in test data point [True]
43 Text feature [see] present in test data point [True]
44 Text feature [related] present in test data point [True]
46 Text feature [either] present in test data point [True]
47 Text feature [indicate] present in test data point [True]
48 Text feature [mutations] present in test data point [True]
50 Text feature [transfected] present in test data point [True]
51 Text feature [wild] present in test data point [True]
52 Text feature [contribute] present in test data point [True]
55 Text feature [associated] present in test data point [True]
56 Text feature [vector] present in test data point [True]
57 Text feature [vivo] present in test data point [True]
58 Text feature [vitro] present in test data point [True]
59 Text feature [30] present in test data point [True]
60 Text feature [may] present in test data point [True]
61 Text feature [although] present in test data point [True]
63 Text feature [using] present in test data point [True]
65 Text feature [thus] present in test data point [True]
66 Text feature [effects] present in test data point [True]
67 Text feature [reduced] present in test data point [True]
69 Text feature [lower] present in test data point [True]
71 Text feature [missense] present in test data point [True]
73 Text feature [suggest] present in test data point [True]
74 Text feature [expression] present in test data point [True]
80 Text feature [cells] present in test data point [True]
81 Text feature [tagged] present in test data point [True]
82 Text feature [previously] present in test data point [True]
84 Text feature [suppressor] present in test data point [True]
85 Text feature [discussion] present in test data point [True]
86 Text feature [three] present in test data point [True]
88 Text feature [functional] present in test data point [True]
90 Text feature [lack] present in test data point [True]
92 Text feature [suggesting] present in test data point [True]
93 Text feature [made] present in test data point [True]
96 Text feature [indicates] present in test data point [True]
98 Text feature [performed] present in test data point [True]
100 Text feature [similar] present in test data point [True]
102 Text feature [addition] present in test data point [True]
103 Text feature [found] present in test data point [True]
106 Text feature [introduction] present in test data point [True]
108 Text feature [analysis] present in test data point [True]
110 Text feature [percentage] present in test data point [True]
113 Text feature [fact] present in test data point [True]
114 Text feature [terminal] present in test data point [True]
116 Text feature [however] present in test data point [True]
118 Text feature [cellular] present in test data point [True]
119 Text feature [yielded] present in test data point [True]
120 Text feature [amount] present in test data point [True]
121 Text feature [generated] present in test data point [True]
123 Text feature [indicating] present in test data point [True]
127 Text feature [suggested] present in test data point [True]
130 Text feature [role] present in test data point [True]
132 Text feature [one] present in test data point [True]
134 Text feature [several] present in test data point [True]
137 Text feature [show] present in test data point [True]
139 Text feature [mm] present in test data point [True]
141 Text feature [required] present in test data point [True]
143 Text feature [mutants] present in test data point [True]
146 Text feature [control] present in test data point [True]
147 Text feature [high] present in test data point [True]
148 Text feature [binding] present in test data point [True]
149 Text feature [essential] present in test data point [True]
151 Text feature [previous] present in test data point [True]
153 Text feature [rather] present in test data point [True]
156 Text feature [co] present in test data point [True]
157 Text feature [experiment] present in test data point [True]
158 Text feature [low] present in test data point [True]
159 Text feature [figure] present in test data point [True]
160 Text feature [reported] present in test data point [True]
162 Text feature [involved] present in test data point [True]
163 Text feature [lysates] present in test data point [True]
164 Text feature [error] present in test data point [True]
165 Text feature [dependent] present in test data point [True]
167 Text feature [critical] present in test data point [True]
171 Text feature [within] present in test data point [True]
172 Text feature [tested] present in test data point [True]

173 Text feature [bars] present in test data point [True]
177 Text feature [average] present in test data point [True]
180 Text feature [cycle] present in test data point [True]
181 Text feature [15] present in test data point [True]
182 Text feature [together] present in test data point [True]
183 Text feature [key] present in test data point [True]
184 Text feature [10] present in test data point [True]
186 Text feature [fully] present in test data point [True]
188 Text feature [mutant] present in test data point [True]
190 Text feature [regions] present in test data point [True]
191 Text feature [human] present in test data point [True]
192 Text feature [respectively] present in test data point [True]
193 Text feature [monitored] present in test data point [True]
194 Text feature [frequently] present in test data point [True]
196 Text feature [full] present in test data point [True]
197 Text feature [dna] present in test data point [True]
201 Text feature [effect] present in test data point [True]
203 Text feature [probably] present in test data point [True]
206 Text feature [mediated] present in test data point [True]
207 Text feature [included] present in test data point [True]
208 Text feature [including] present in test data point [True]
209 Text feature [acids] present in test data point [True]
211 Text feature [phenotype] present in test data point [True]
212 Text feature [specific] present in test data point [True]
215 Text feature [ala] present in test data point [True]
216 Text feature [changes] present in test data point [True]
218 Text feature [mutation] present in test data point [True]
221 Text feature [used] present in test data point [True]
222 Text feature [major] present in test data point [True]
223 Text feature [fig] present in test data point [True]
226 Text feature [except] present in test data point [True]
227 Text feature [well] present in test data point [True]
228 Text feature [provide] present in test data point [True]
229 Text feature [relevant] present in test data point [True]
232 Text feature [test] present in test data point [True]
234 Text feature [due] present in test data point [True]
236 Text feature [affect] present in test data point [True]
237 Text feature [note] present in test data point [True]
239 Text feature [targeting] present in test data point [True]
244 Text feature [data] present in test data point [True]
245 Text feature [derived] present in test data point [True]
246 Text feature [relative] present in test data point [True]
248 Text feature [prepared] present in test data point [True]
251 Text feature [cell] present in test data point [True]
254 Text feature [40] present in test data point [True]
258 Text feature [compared] present in test data point [True]
259 Text feature [anti] present in test data point [True]
263 Text feature [assay] present in test data point [True]
264 Text feature [yellow] present in test data point [True]
266 Text feature [hypothesis] present in test data point [True]
268 Text feature [might] present in test data point [True]
269 Text feature [multiple] present in test data point [True]
272 Text feature [antibody] present in test data point [True]
277 Text feature [properties] present in test data point [True]
280 Text feature [gene] present in test data point [True]
282 Text feature [sequences] present in test data point [True]
283 Text feature [revealed] present in test data point [True]
285 Text feature [table] present in test data point [True]
286 Text feature [measured] present in test data point [True]
287 Text feature [observation] present in test data point [True]
289 Text feature [according] present in test data point [True]
295 Text feature [contrast] present in test data point [True]
296 Text feature [distinct] present in test data point [True]
297 Text feature [level] present in test data point [True]
299 Text feature [consistent] present in test data point [True]
389 Text feature [site] present in test data point [True]
398 Text feature [page] present in test data point [True]
401 Text feature [putative] present in test data point [True]
403 Text feature [single] present in test data point [True]
404 Text feature [obtained] present in test data point [True]
414 Text feature [many] present in test data point [True]
415 Text feature [mammals] present in test data point [True]
416 Text feature [observed] present in test data point [True]
417 Text feature [induced] present in test data point [True]
419 Text feature [representative] present in test data point [True]
421 Text feature [domain] present in test data point [True]
422 Text feature [different] present in test data point [True]

425 Text feature [taken] present in test data point [True]
426 Text feature [majority] present in test data point [True]
427 Text feature [assayed] present in test data point [True]
431 Text feature [tumor] present in test data point [True]
434 Text feature [another] present in test data point [True]
435 Text feature [significant] present in test data point [True]
436 Text feature [display] present in test data point [True]
440 Text feature [methods] present in test data point [True]
441 Text feature [25] present in test data point [True]
444 Text feature [1998] present in test data point [True]
446 Text feature [assessed] present in test data point [True]
448 Text feature [supporting] present in test data point [True]
449 Text feature [five] present in test data point [True]
450 Text feature [present] present in test data point [True]
451 Text feature [negative] present in test data point [True]
454 Text feature [substantially] present in test data point [True]
456 Text feature [based] present in test data point [True]
460 Text feature [50] present in test data point [True]
462 Text feature [furthermore] present in test data point [True]
465 Text feature [distribution] present in test data point [True]
466 Text feature [possibly] present in test data point [True]
469 Text feature [additional] present in test data point [True]
470 Text feature [complex] present in test data point [True]
472 Text feature [presence] present in test data point [True]
476 Text feature [recent] present in test data point [True]
478 Text feature [importance] present in test data point [True]
480 Text feature [ml] present in test data point [True]
481 Text feature [finally] present in test data point [True]
482 Text feature [endogenous] present in test data point [True]
483 Text feature [examine] present in test data point [True]
485 Text feature [predicted] present in test data point [True]
487 Text feature [linked] present in test data point [True]
488 Text feature [western] present in test data point [True]
489 Text feature [highly] present in test data point [True]
490 Text feature [antibodies] present in test data point [True]
491 Text feature [nuclear] present in test data point [True]
492 Text feature [dd] present in test data point [True]
493 Text feature [remarkably] present in test data point [True]
494 Text feature [considered] present in test data point [True]
502 Text feature [would] present in test data point [True]
504 Text feature [without] present in test data point [True]
507 Text feature [normal] present in test data point [True]
508 Text feature [examined] present in test data point [True]
509 Text feature [1a] present in test data point [True]
512 Text feature [proposed] present in test data point [True]
514 Text feature [cloned] present in test data point [True]
521 Text feature [thought] present in test data point [True]
524 Text feature [24] present in test data point [True]
525 Text feature [specifically] present in test data point [True]
526 Text feature [study] present in test data point [True]
528 Text feature [analyses] present in test data point [True]
531 Text feature [et] present in test data point [True]
533 Text feature [times] present in test data point [True]
534 Text feature [large] present in test data point [True]
535 Text feature [investigate] present in test data point [True]
541 Text feature [regulates] present in test data point [True]
542 Text feature [stained] present in test data point [True]
544 Text feature [deletion] present in test data point [True]
545 Text feature [interestingly] present in test data point [True]
548 Text feature [represent] present in test data point [True]
550 Text feature [100] present in test data point [True]
552 Text feature [investigated] present in test data point [True]
554 Text feature [31] present in test data point [True]
555 Text feature [even] present in test data point [True]
559 Text feature [27] present in test data point [True]
561 Text feature [total] present in test data point [True]
564 Text feature [flag] present in test data point [True]
565 Text feature [26] present in test data point [True]
566 Text feature [all] present in test data point [True]
567 Text feature [surrounded] present in test data point [True]
569 Text feature [nonsense] present in test data point [True]
570 Text feature [support] present in test data point [True]
574 Text feature [2000] present in test data point [True]
577 Text feature [arrest] present in test data point [True]
579 Text feature [larger] present in test data point [True]
587 Text feature [1b] present in test data point [True]
588 Text feature [studies] present in test data point [True]

589 Text feature [regulation] present in test data point [True]
592 Text feature [moreover] present in test data point [True]
593 Text feature [showed] present in test data point [True]
595 Text feature [region] present in test data point [True]
597 Text feature [variety] present in test data point [True]
602 Text feature [since] present in test data point [True]
604 Text feature [blue] present in test data point [True]
606 Text feature [sufficient] present in test data point [True]
609 Text feature [37] present in test data point [True]
610 Text feature [confirmed] present in test data point [True]
611 Text feature [involvement] present in test data point [True]
612 Text feature [inactivation] present in test data point [True]
616 Text feature [expected] present in test data point [True]
622 Text feature [via] present in test data point [True]
631 Text feature [disruption] present in test data point [True]
633 Text feature [residue] present in test data point [True]
635 Text feature [unable] present in test data point [True]
637 Text feature [indeed] present in test data point [True]
638 Text feature [central] present in test data point [True]
639 Text feature [limited] present in test data point [True]
641 Text feature [sites] present in test data point [True]
642 Text feature [significantly] present in test data point [True]
643 Text feature [us] present in test data point [True]
645 Text feature [displayed] present in test data point [True]
647 Text feature [independent] present in test data point [True]
650 Text feature [dominant] present in test data point [True]
652 Text feature [supplemented] present in test data point [True]
656 Text feature [3b] present in test data point [True]
657 Text feature [functionally] present in test data point [True]
658 Text feature [conditions] present in test data point [True]
662 Text feature [latter] present in test data point [True]
664 Text feature [likely] present in test data point [True]
668 Text feature [complete] present in test data point [True]
671 Text feature [cancer] present in test data point [True]
675 Text feature [sequence] present in test data point [True]
676 Text feature [growth] present in test data point [True]
680 Text feature [provided] present in test data point [True]
682 Text feature [29] present in test data point [True]
686 Text feature [identified] present in test data point [True]
689 Text feature [upon] present in test data point [True]
691 Text feature [increased] present in test data point [True]
692 Text feature [nucleus] present in test data point [True]
693 Text feature [hereditary] present in test data point [True]
695 Text feature [extent] present in test data point [True]
697 Text feature [shows] present in test data point [True]
700 Text feature [spectrum] present in test data point [True]
701 Text feature [grown] present in test data point [True]
706 Text feature [common] present in test data point [True]
707 Text feature [plates] present in test data point [True]
714 Text feature [20] present in test data point [True]
715 Text feature [16] present in test data point [True]
720 Text feature [substrate] present in test data point [True]
721 Text feature [12] present in test data point [True]
726 Text feature [besides] present in test data point [True]
731 Text feature [conclusions] present in test data point [True]
734 Text feature [suggests] present in test data point [True]
739 Text feature [recently] present in test data point [True]
740 Text feature [top] present in test data point [True]
741 Text feature [genetic] present in test data point [True]
750 Text feature [instability] present in test data point [True]
751 Text feature [active] present in test data point [True]
754 Text feature [altered] present in test data point [True]
757 Text feature [form] present in test data point [True]
761 Text feature [null] present in test data point [True]
764 Text feature [explain] present in test data point [True]
766 Text feature [affinity] present in test data point [True]
773 Text feature [4a] present in test data point [True]
774 Text feature [like] present in test data point [True]
775 Text feature [defective] present in test data point [True]
776 Text feature [example] present in test data point [True]
779 Text feature [1c] present in test data point [True]
780 Text feature [13] present in test data point [True]
781 Text feature [80] present in test data point [True]
782 Text feature [next] present in test data point [True]
785 Text feature [alignment] present in test data point [True]
786 Text feature [explained] present in test data point [True]
788 Text feature [aqarose] present in test data point [True]

794 Text feature [carrying] present in test data point [True]
795 Text feature [resulted] present in test data point [True]
799 Text feature [process] present in test data point [True]
800 Text feature [prevent] present in test data point [True]
801 Text feature [damage] present in test data point [True]
804 Text feature [deleted] present in test data point [True]
810 Text feature [modified] present in test data point [True]
812 Text feature [figs] present in test data point [True]
813 Text feature [3a] present in test data point [True]
814 Text feature [evidenced] present in test data point [True]
816 Text feature [occur] present in test data point [True]
817 Text feature [members] present in test data point [True]
818 Text feature [measure] present in test data point [True]
819 Text feature [hereafter] present in test data point [True]
823 Text feature [work] present in test data point [True]
825 Text feature [bands] present in test data point [True]
831 Text feature [coupled] present in test data point [True]
836 Text feature [first] present in test data point [True]
837 Text feature [efficiently] present in test data point [True]
839 Text feature [wt] present in test data point [True]
841 Text feature [yet] present in test data point [True]
842 Text feature [still] present in test data point [True]
844 Text feature [notion] present in test data point [True]
845 Text feature [provides] present in test data point [True]
846 Text feature [difficult] present in test data point [True]
847 Text feature [demonstrate] present in test data point [True]
848 Text feature [200] present in test data point [True]
849 Text feature [play] present in test data point [True]
850 Text feature [give] present in test data point [True]
853 Text feature [staining] present in test data point [True]
854 Text feature [identical] present in test data point [True]
855 Text feature [especially] present in test data point [True]
861 Text feature [recognized] present in test data point [True]
862 Text feature [similarly] present in test data point [True]
863 Text feature [corresponding] present in test data point [True]
865 Text feature [mouse] present in test data point [True]
866 Text feature [development] present in test data point [True]
867 Text feature [reduction] present in test data point [True]
871 Text feature [entire] present in test data point [True]
873 Text feature [observations] present in test data point [True]
874 Text feature [coding] present in test data point [True]
877 Text feature [regulate] present in test data point [True]
878 Text feature [mechanism] present in test data point [True]
882 Text feature [groups] present in test data point [True]
883 Text feature [selected] present in test data point [True]
886 Text feature [culture] present in test data point [True]
887 Text feature [4b] present in test data point [True]
888 Text feature [point] present in test data point [True]
890 Text feature [dual] present in test data point [True]
892 Text feature [small] present in test data point [True]
893 Text feature [subset] present in test data point [True]
894 Text feature [right] present in test data point [True]
895 Text feature [phenotypes] present in test data point [True]
896 Text feature [disease] present in test data point [True]
904 Text feature [accumulation] present in test data point [True]
909 Text feature [isolated] present in test data point [True]
913 Text feature [medium] present in test data point [True]
914 Text feature [carried] present in test data point [True]
916 Text feature [33] present in test data point [True]
917 Text feature [processes] present in test data point [True]
918 Text feature [six] present in test data point [True]
923 Text feature [1999] present in test data point [True]
925 Text feature [among] present in test data point [True]
926 Text feature [difference] present in test data point [True]
929 Text feature [known] present in test data point [True]
934 Text feature [understanding] present in test data point [True]
941 Text feature [appropriate] present in test data point [True]
943 Text feature [detect] present in test data point [True]
947 Text feature [genes] present in test data point [True]
952 Text feature [santa] present in test data point [True]
954 Text feature [panel] present in test data point [True]
955 Text feature [early] present in test data point [True]
959 Text feature [remaining] present in test data point [True]
960 Text feature [identification] present in test data point [True]
961 Text feature [polymerase] present in test data point [True]
964 Text feature [number] present in test data point [True]
965 Text feature [2b] present in test data point [True]

```

967 Text feature [absence] present in test data point [True]
968 Text feature [35] present in test data point [True]
970 Text feature [locus] present in test data point [True]
974 Text feature [tumorigenesis] present in test data point [True]
976 Text feature [seems] present in test data point [True]
978 Text feature [image] present in test data point [True]
980 Text feature [clear] present in test data point [True]
983 Text feature [applied] present in test data point [True]
987 Text feature [cancers] present in test data point [True]
989 Text feature [32] present in test data point [True]
992 Text feature [monoclonal] present in test data point [True]
993 Text feature [particular] present in test data point [True]
995 Text feature [homozygous] present in test data point [True]
997 Text feature [suppressors] present in test data point [True]
Out of the top 1000 features 416 are present in query point

```

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

In [91]:

```

# 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_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

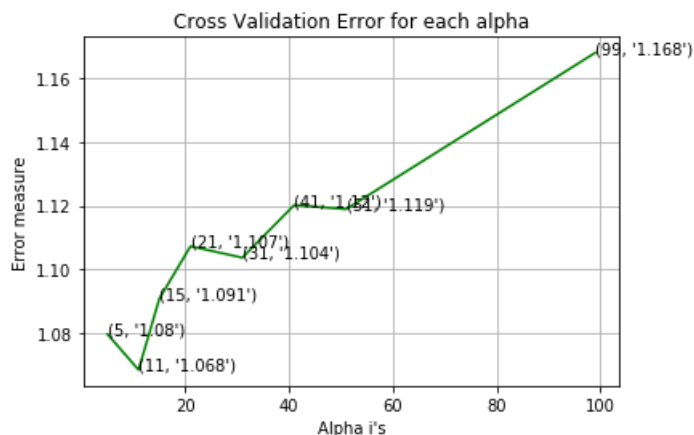
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ',
      alpha[best_alpha],
      "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 5
Log Loss : 1.0795439630377217
for alpha = 11
Log Loss : 1.0683061181112927
for alpha = 15
Log Loss : 1.0906890668476643
for alpha = 21
Log Loss : 1.1072009818093094
for alpha = 31
Log Loss : 1.1036161316699014
for alpha = 41
Log Loss : 1.1201871993944572
for alpha = 51
Log Loss : 1.1188904390996472
for alpha = 99
Log Loss : 1.1681250137822985

```



```

For values of best alpha = 11 The train log loss is: 0.6488621065913742
For values of best alpha = 11 The cross validation log loss is: 1.0683061181112927
For values of best alpha = 11 The test log loss is: 1.0146612573885447

```

4.2.2. Testing the model with best hyper paramters

In [92]:

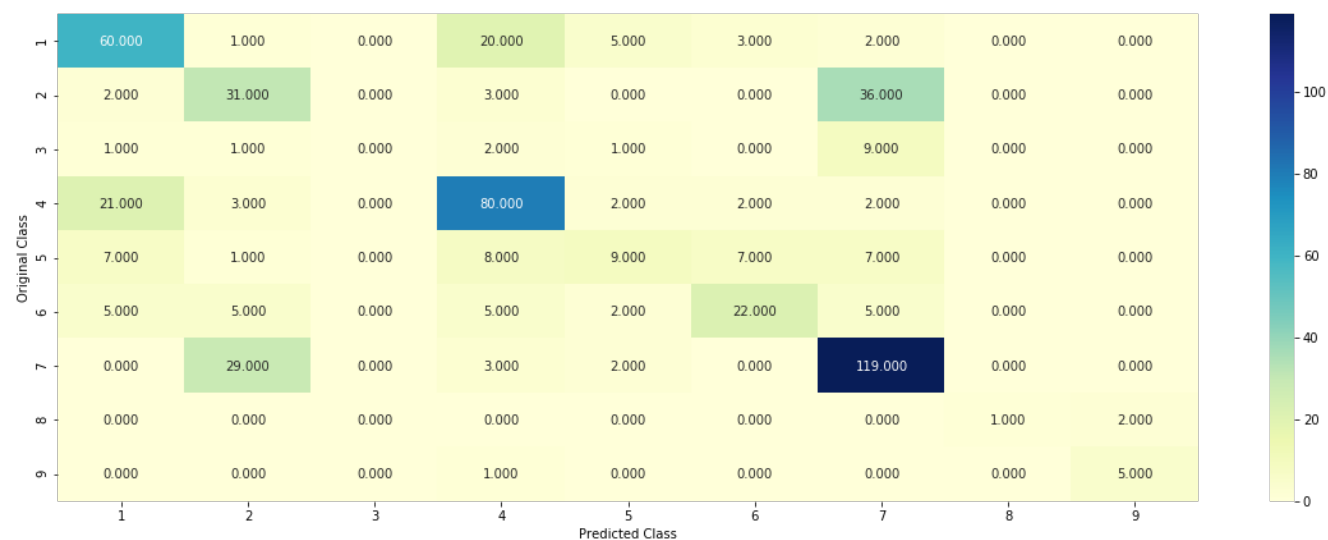
```
# find more about KNeighborsClassifier()
# here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None,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

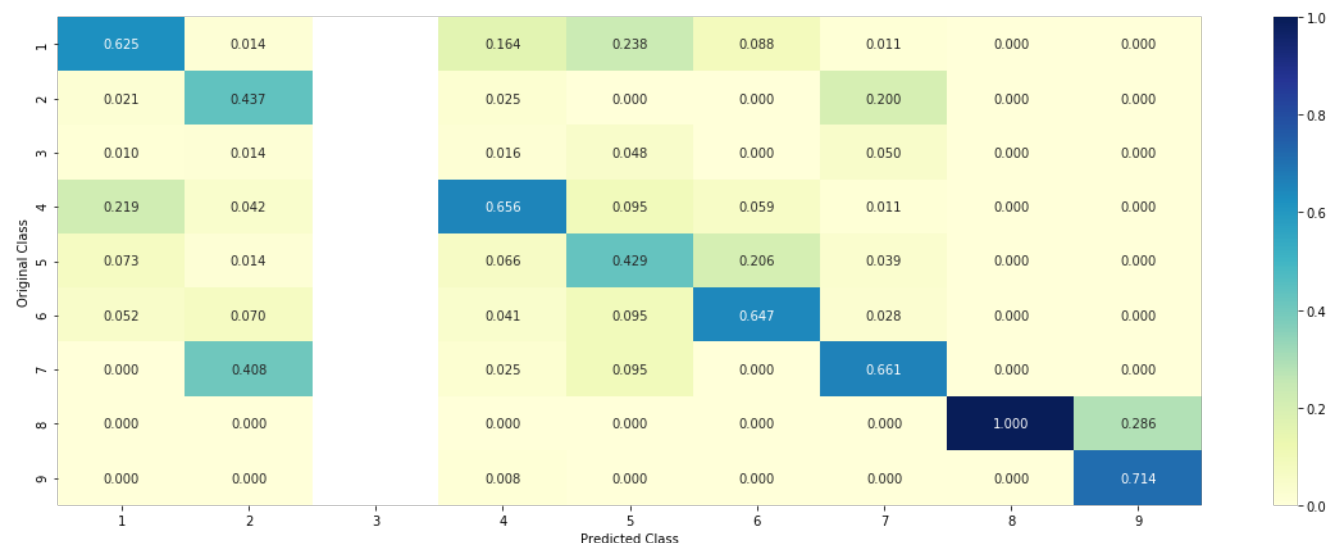
Log loss : 1.0683061181112927

Number of mis-classified points : 0.38533834586466165

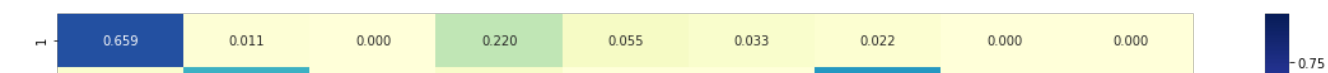
----- Confusion matrix -----

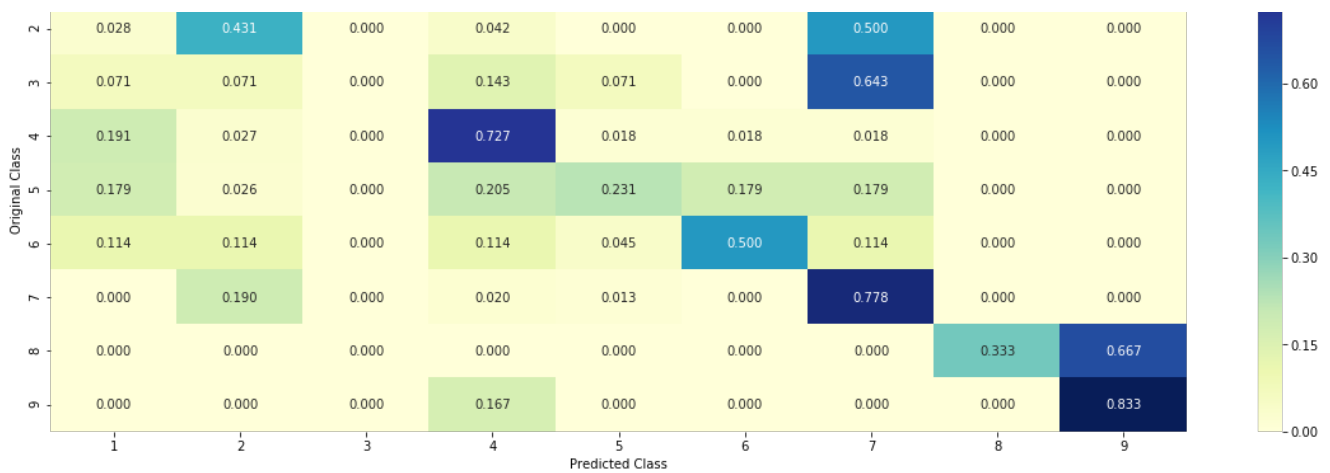


----- Precision matrix (Column Sum=1) -----



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





4.2.3. Sample Query point -1

In [93]:

```
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 : 7

Actual Class : 5

The 11 nearest neighbours of the test points belongs to classes [7 3 7 5 6 7 2 2 7 6 2]

Fequency of nearest points : Counter({7: 4, 2: 3, 6: 2, 3: 1, 5: 1})

4.2.4. Sample Query Point-2

In [94]:

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

test_point_index = 100

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

Predicted Class : 1

Actual Class : 1

the k value for knn is 11 and the nearest neighbours of the test points belongs to classes [1 1 1 1 4 1 1 1 1 1 1]

Fequency of nearest points : Counter({1: 10, 4: 1})

4.3 Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper paramter tuning

In [95]:

```
# 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/
#-----

# 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 probabilitates 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))

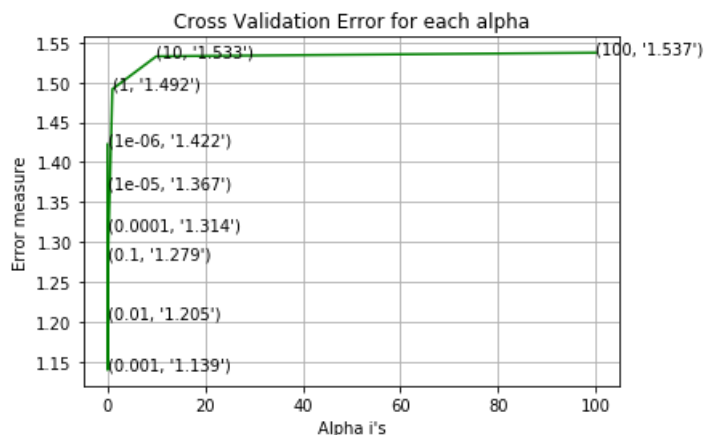
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha],
      "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 1e-06
Log Loss : 1.422444502240648
for alpha = 1e-05
Log Loss : 1.366863075870454
for alpha = 0.0001
Log Loss : 1.3141995015404864
for alpha = 0.001
Log Loss : 1.1394779231665602
for alpha = 0.01
Log Loss : 1.204757853821474
for alpha = 0.1
Log Loss : 1.2787309923144823
for alpha = 1
Log Loss : 1.4915208470155683
for alpha = 10
Log Loss : 1.5327466246827102
for alpha = 100
Log Loss : 1.5372282076365054

```



```

For values of best alpha = 0.001 The train log loss is: 0.5819240279166931
For values of best alpha = 0.001 The cross validation log loss is: 1.1394779231665602
For values of best alpha = 0.001 The test log loss is: 1.0702274828843725

```

4.3.1.2. Testing the model with best hyper paramters

In [96]:

```

# 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-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)
```

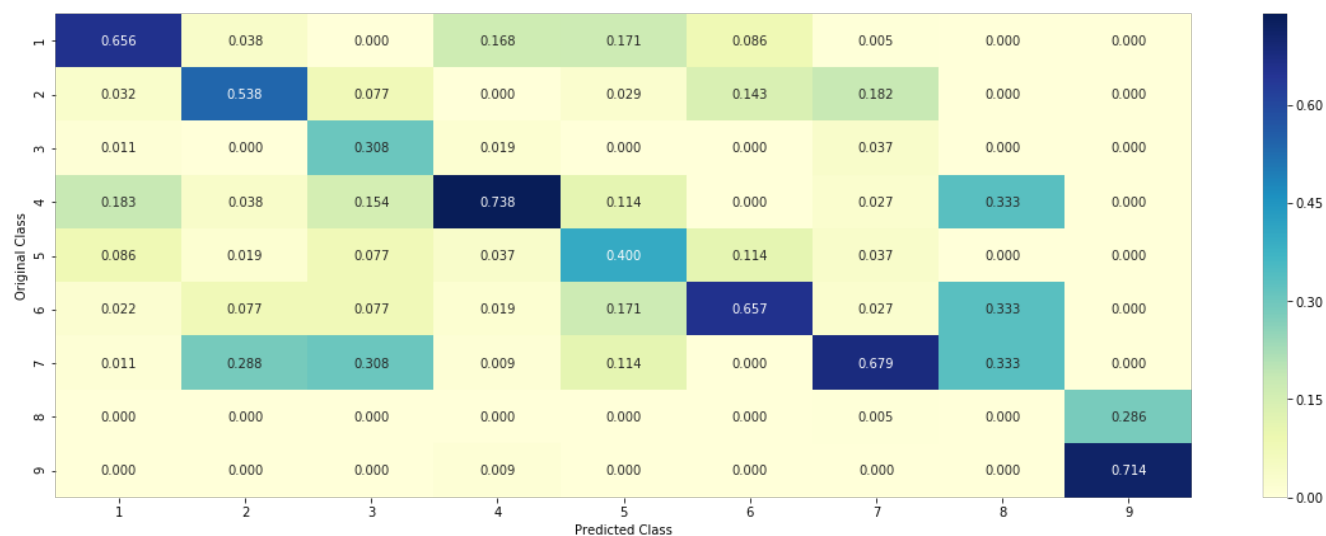
Log loss : 1.1394779231665602

Number of mis-classified points : 0.35902255639097747

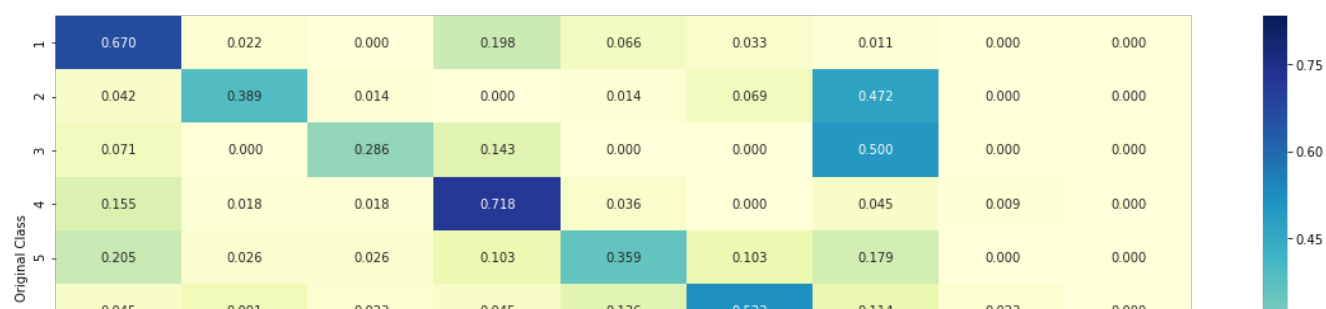
----- Confusion matrix -----

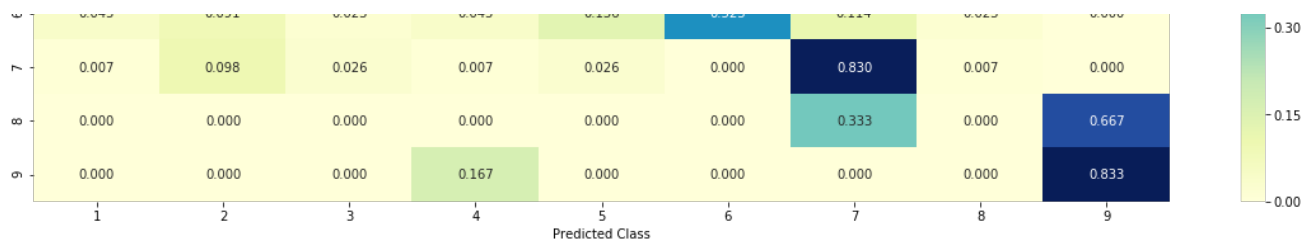


----- Precision matrix (Column Sum=1) -----



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





4.3.1.3. Feature Importance

In [97]:

```
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. Incorrectly Classified point

In [98]:

```
# from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
test_point_index = 1
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)
```

Predicted Class : 7

Predicted Class Probabilities: [[0.025 0.261 0.0079 0.031 0.0358 0.0116 0.5993 0.0175 0.011]]

Actual Class : 5

```
-----
19 Text feature [constitutive] present in test data point [True]
21 Text feature [constitutively] present in test data point [True]
44 Text feature [activated] present in test data point [True]
66 Text feature [technology] present in test data point [True]
100 Text feature [activation] present in test data point [True]
108 Text feature [transforming] present in test data point [True]
144 Text feature [ba] present in test data point [True]
157 Text feature [constants] present in test data point [True]
158 Text feature [downstream] present in test data point [True]
276 Text feature [rap] present in test data point [True]
285 Text feature [anomalous] present in test data point [True]
```

```

205 Text feature [nanomolar] present in test data point [True]
410 Text feature [phosphorylation] present in test data point [True]
434 Text feature [transformation] present in test data point [True]
518 Text feature [braf] present in test data point [True]
530 Text feature [kinase] present in test data point [True]
553 Text feature [activate] present in test data point [True]
627 Text feature [activating] present in test data point [True]
662 Text feature [v12] present in test data point [True]
686 Text feature [doses] present in test data point [True]
699 Text feature [expressing] present in test data point [True]
728 Text feature [erk] present in test data point [True]
789 Text feature [rbd] present in test data point [True]
818 Text feature [balb] present in test data point [True]
834 Text feature [signaling] present in test data point [True]
883 Text feature [erk1] present in test data point [True]
955 Text feature [inhibitor] present in test data point [True]
Out of the top 1000 features 26 are present in query point

```

4.3.1.3.2. Correctly Classified point

In [99]:

```

test_point_index = 55
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)

```

Predicted Class : 4

Predicted Class Probabilities: [[0.0018 0.0191 0.0026 0.9482 0.0082 0.0015 0.0024 0.0135 0.0028]]

Actual Class : 4

```

-----
39 Text feature [shrunken] present in test data point [True]
81 Text feature [devoid] present in test data point [True]
156 Text feature [shen] present in test data point [True]
260 Text feature [instability] present in test data point [True]
265 Text feature [hypophosphorylated] present in test data point [True]
277 Text feature [suppressor] present in test data point [True]
285 Text feature [biallelic] present in test data point [True]
403 Text feature [stk11] present in test data point [True]
414 Text feature [dn] present in test data point [True]
427 Text feature [obliterate] present in test data point [True]
450 Text feature [pten] present in test data point [True]
458 Text feature [homologues] present in test data point [True]
476 Text feature [dysfunctions] present in test data point [True]
494 Text feature [monolayers] present in test data point [True]
612 Text feature [triggering] present in test data point [True]
731 Text feature [senescence] present in test data point [True]
737 Text feature [78k] present in test data point [True]
776 Text feature [apoptosis] present in test data point [True]
799 Text feature [roulston] present in test data point [True]
804 Text feature [mammalian] present in test data point [True]
807 Text feature [irradiation] present in test data point [True]
824 Text feature [fbxw7] present in test data point [True]
923 Text feature [kaufman] present in test data point [True]
949 Text feature [nucleus] present in test data point [True]
Out of the top 1000 features 24 are present in query point

```

Logistic Regression

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter 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/
#-----

# 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)
```

```

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))

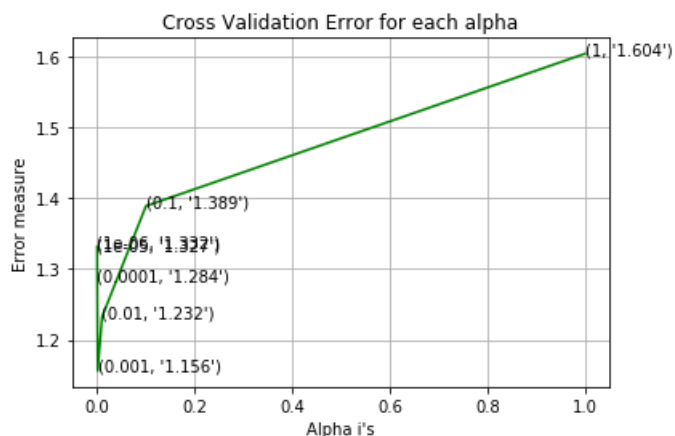
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha],
      "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 1e-06
Log Loss : 1.3315665854936976
for alpha = 1e-05
Log Loss : 1.3269445980972805
for alpha = 0.0001
Log Loss : 1.2839797580400472
for alpha = 0.001
Log Loss : 1.1560628416196026
for alpha = 0.01
Log Loss : 1.2324143067974982
for alpha = 0.1
Log Loss : 1.3890238671036788
for alpha = 1
Log Loss : 1.6042423928376766

```



```

For values of best alpha = 0.001 The train log loss is: 0.5725749905997384
For values of best alpha = 0.001 The cross validation log loss is: 1.1560628416196026
For values of best alpha = 0.001 The test log loss is: 1.0900545958618366

```

4.3.2.2. Testing model with best hyper parameters

In [101]:

```

# 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.

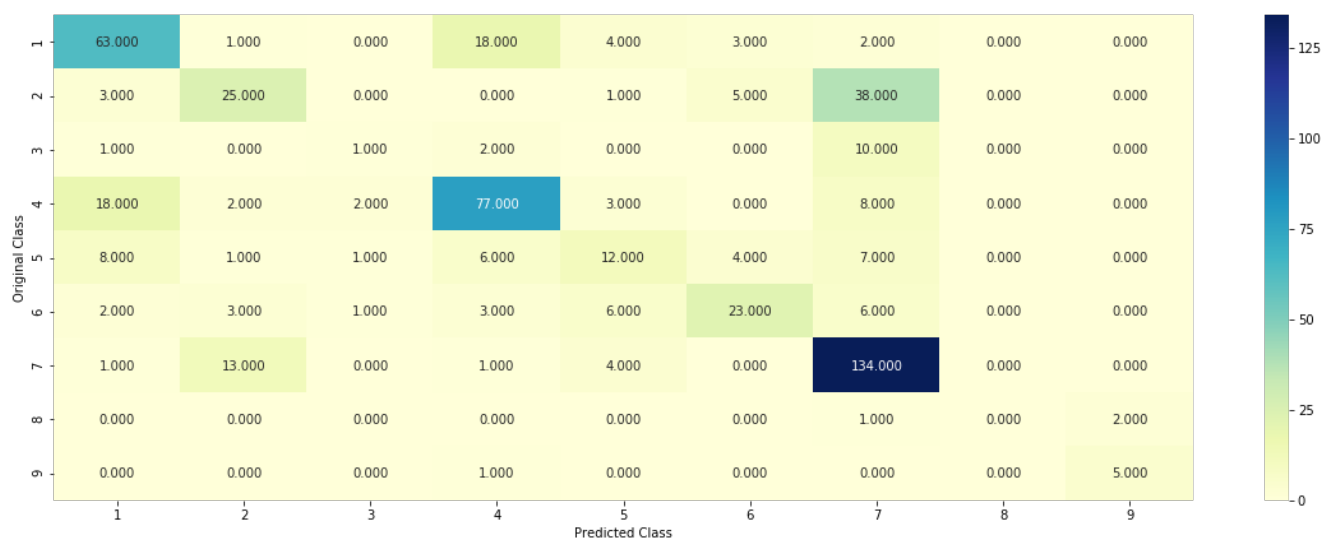
#-----
# video link:
#-----

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)

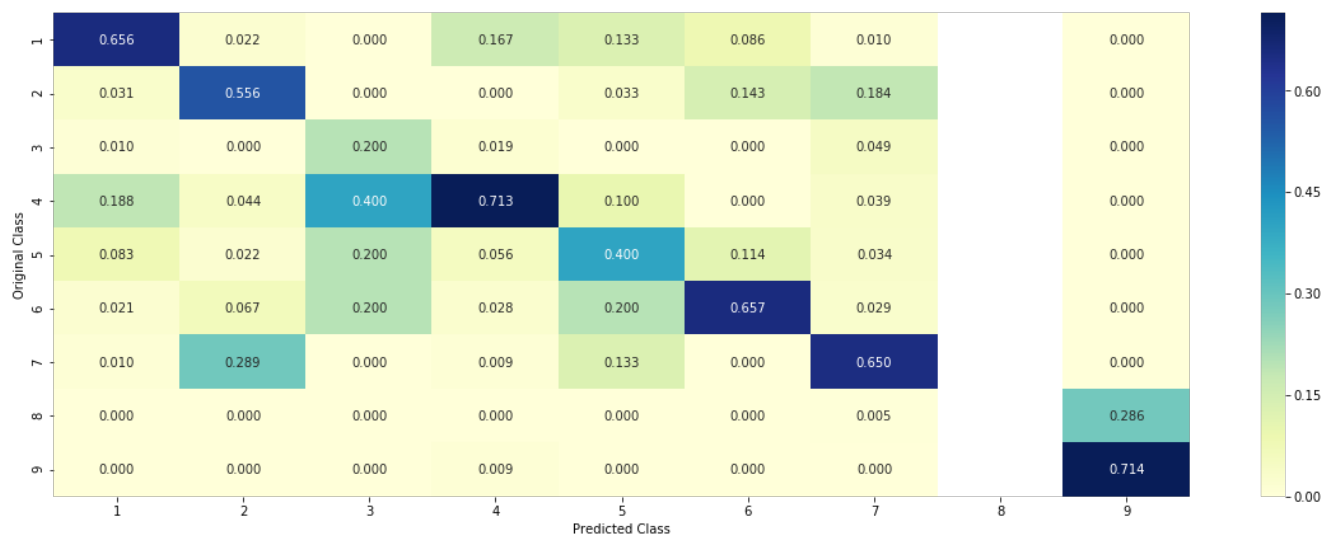
```

Log loss : 1.1560628416196026
Number of mis-classified points : 0.3609022556390977

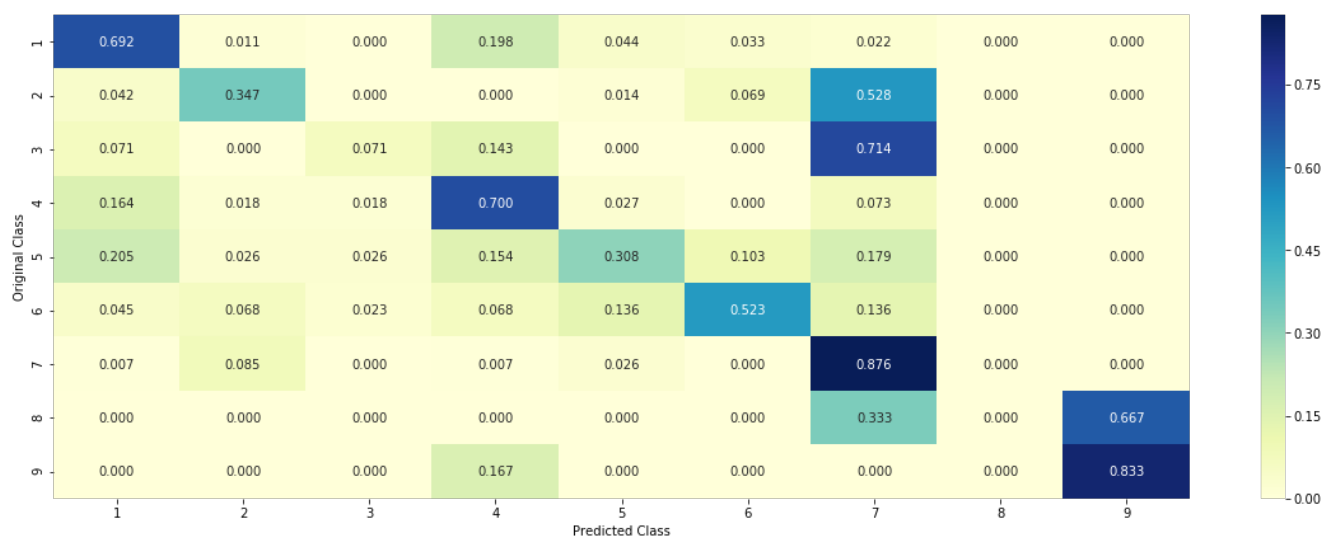
----- Confusion matrix -----



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



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



4.3.2.3. Feature Importance, Correctly Classified point

In [102]:

```
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
test_point_index = 100
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)
```

Predicted Class : 1

Predicted Class Probabilities: [[0.5373 0.0947 0.0062 0.1429 0.0219 0.0177 0.1729 0.0058 0.0007]]

Actual Class : 1

202 Text feature [foxa1] present in test data point [True]

506 Text feature [winged] present in test data point [True]

683 Text feature [gt] present in test data point [True]

897 Text feature [lysate] present in test data point [True]

982 Text feature [stated] present in test data point [True]

Out of the top 1000 features 5 are present in query point

4.3.2.4. Feature Importance, Inorrectly Classified point

In [103]:

```
test_point_index = 1
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)
```

Predicted Class : 7

Predicted Class Probabilities: [[2.900e-02 2.475e-01 1.200e-03 3.260e-02 1.790e-02 7.200e-03 6.579e-01

6.600e-03 1.000e-04]]

Actual Class : 5

66 Text feature [constitutively] present in test data point [True]

67 Text feature [constitutive] present in test data point [True]

100 Text feature [activated] present in test data point [True]

141 Text feature [transforming] present in test data point [True]

147 Text feature [technology] present in test data point [True]

149 Text feature [activation] present in test data point [True]

161 Text feature [ba] present in test data point [True]

183 Text feature [downstream] present in test data point [True]

251 Text feature [phosphorylation] present in test data point [True]

263 Text feature [braf] present in test data point [True]

270 Text feature [expressing] present in test data point [True]

275 Text feature [rap] present in test data point [True]

288 Text feature [transformation] present in test data point [True]

289 Text feature [activate] present in test data point [True]

315 Text feature [kinase] present in test data point [True]

322 Text feature [nanomolar] present in test data point [True]

347 Text feature [activating] present in test data point [True]

```

377 Text feature [activating] present in test data point [True]
382 Text feature [constants] present in test data point [True]
472 Text feature [doses] present in test data point [True]
520 Text feature [signaling] present in test data point [True]
549 Text feature [erk] present in test data point [True]
577 Text feature [inhibitor] present in test data point [True]
638 Text feature [v600e] present in test data point [True]
655 Text feature [v12] present in test data point [True]
680 Text feature [factor] present in test data point [True]
773 Text feature [threonine] present in test data point [True]
814 Text feature [proliferation] present in test data point [True]
817 Text feature [rbd] present in test data point [True]
854 Text feature [erk1] present in test data point [True]
865 Text feature [elevated] present in test data point [True]
869 Text feature [oncogenic] present in test data point [True]
935 Text feature [genotyping] present in test data point [True]
948 Text feature [virus] present in test data point [True]
950 Text feature [lung] present in test data point [True]
959 Text feature [lipid] present in test data point [True]
993 Text feature [adenocarcinoma] present in test data point [True]
Out of the top 1000 features 36 are present in query point

```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

In [104]:

```

# 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)

```



```

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', r
andom_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))

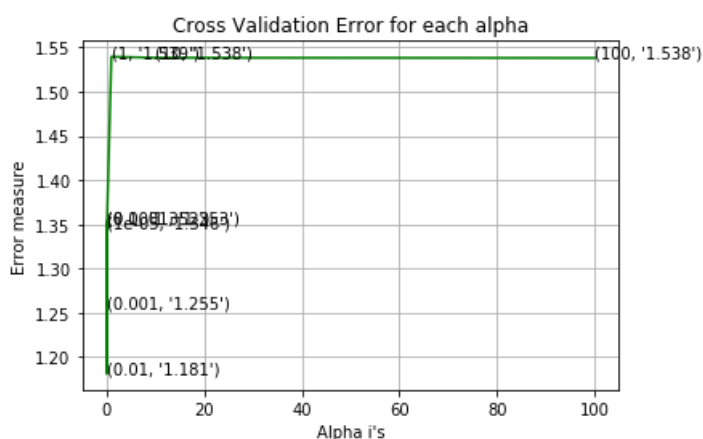
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha],
      "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for C = 1e-05
Log Loss : 1.3464374703389008
for C = 0.0001
Log Loss : 1.3525059896876905
for C = 0.001
Log Loss : 1.2551704837220767
for C = 0.01
Log Loss : 1.180915417541493
for C = 0.1
Log Loss : 1.3523247955700861
for C = 1
Log Loss : 1.539453769919501
for C = 10
Log Loss : 1.538323431956391
for C = 100
Log Loss : 1.5380580187277482

```



For values of best alpha = 0.01 The train log loss is: 0.7113656549491312

For values of best alpha = 0.01 The train log loss is: 0.7113000047451912
 For values of best alpha = 0.01 The cross validation log loss is: 1.180915417541493
 For values of best alpha = 0.01 The test log loss is: 1.113728582819541

4.4.2. Testing model with best hyper parameters

In [105]:

```
# 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)

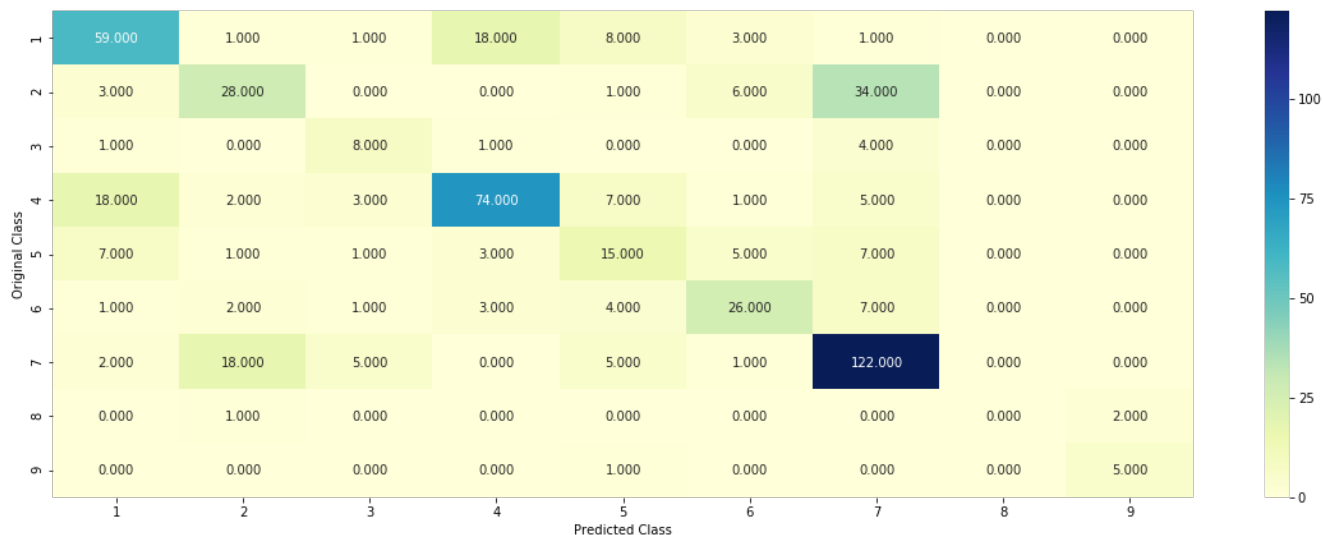
# Some of methods of SVM()
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
# predict(X) Perform classification on samples in X.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

# 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)
```

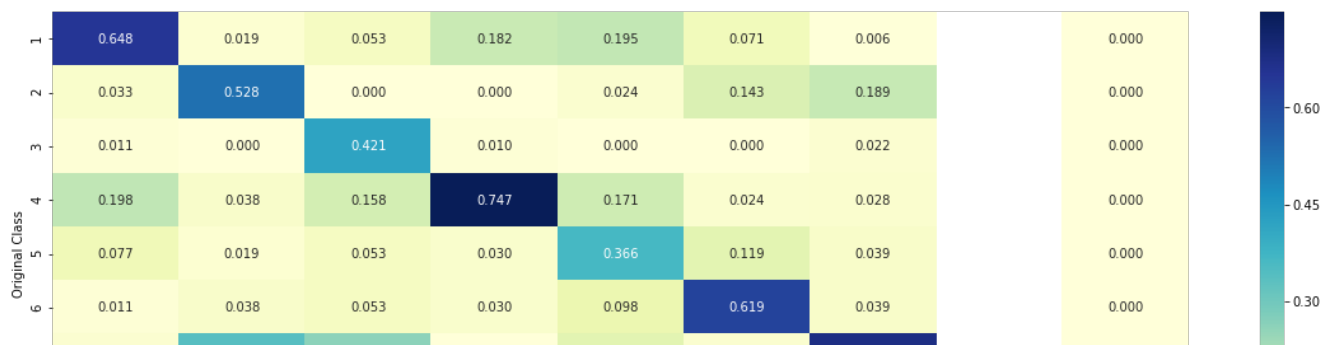
Log loss : 1.180915417541493

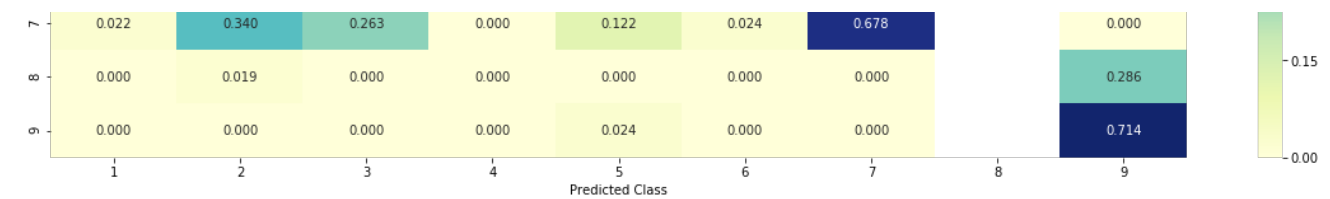
Number of mis-classified points : 0.36654135338345867

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

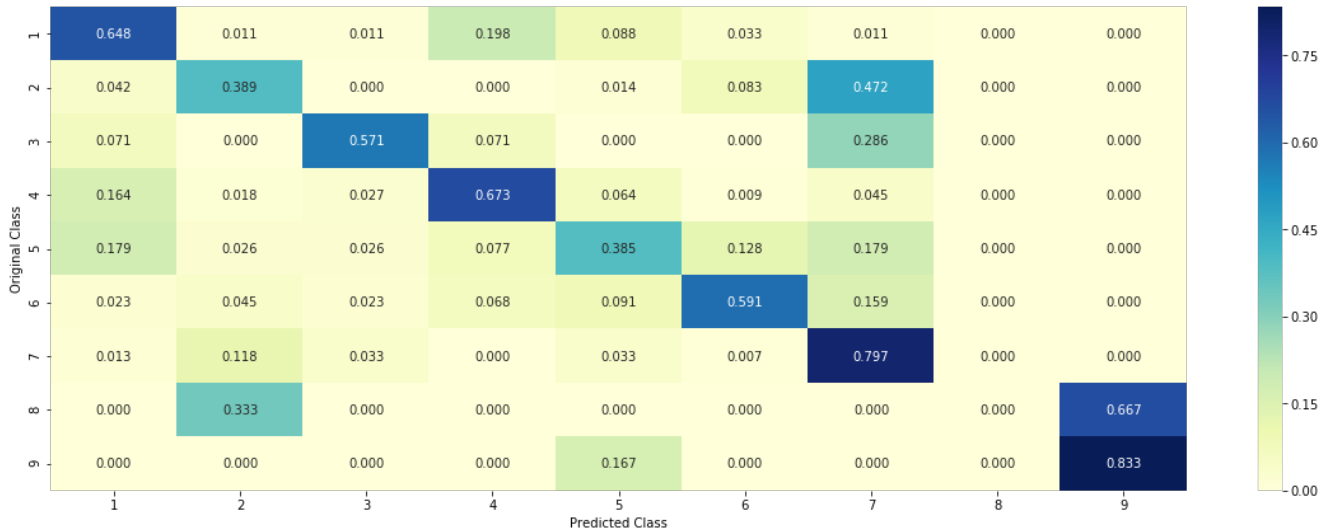


----- Precision matrix (Column Sum=1) -----





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



4.3.3. Feature Importance

4.3.3.1. For Incorrectly classified point

In [106]:

```
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 = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0287 0.3543 0.007 0.0429 0.0333 0.0163 0.5028 0.0116 0.0031]]

Actual Class : 5

```
-----
21 Text feature [constitutive] present in test data point [True]
24 Text feature [constitutively] present in test data point [True]
31 Text feature [activation] present in test data point [True]
37 Text feature [activated] present in test data point [True]
38 Text feature [ba] present in test data point [True]
47 Text feature [expressing] present in test data point [True]
72 Text feature [technology] present in test data point [True]
74 Text feature [transforming] present in test data point [True]
77 Text feature [activate] present in test data point [True]
95 Text feature [nanomolar] present in test data point [True]
112 Text feature [doses] present in test data point [True]
120 Text feature [activating] present in test data point [True]
138 Text feature [phosphorylation] present in test data point [True]
```

```

149 Text feature [transformation] present in test data point [True]
161 Text feature [downstream] present in test data point [True]
176 Text feature [rap] present in test data point [True]
182 Text feature [kinase] present in test data point [True]
241 Text feature [inhibitor] present in test data point [True]
285 Text feature [signaling] present in test data point [True]
291 Text feature [proliferation] present in test data point [True]
306 Text feature [factor] present in test data point [True]
313 Text feature [broadened] present in test data point [True]
345 Text feature [independent] present in test data point [True]
354 Text feature [constants] present in test data point [True]
450 Text feature [oncogenic] present in test data point [True]
468 Text feature [absence] present in test data point [True]
525 Text feature [erk] present in test data point [True]
700 Text feature [spectrometric] present in test data point [True]
713 Text feature [virus] present in test data point [True]
815 Text feature [cc] present in test data point [True]
827 Text feature [inhibited] present in test data point [True]
831 Text feature [v600e] present in test data point [True]
876 Text feature [egf] present in test data point [True]
906 Text feature [erk1] present in test data point [True]
933 Text feature [braf] present in test data point [True]
939 Text feature [intrinsic] present in test data point [True]
940 Text feature [3b] present in test data point [True]
975 Text feature [superposition] present in test data point [True]
990 Text feature [pathways] present in test data point [True]
Out of the top 1000 features 39 are present in query point

```

4.3.3.2. For Incorrectly classified point

In [107]:

```

test_point_index = 31
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)

```

Predicted Class : 2

Predicted Class Probabilities: [[0.0546 0.6336 0.0031 0.0249 0.024 0.0113 0.2367 0.0089 0.0029]]

Actual Class : 7

```

-----
21 Text feature [achieved] present in test data point [True]
25 Text feature [surgical] present in test data point [True]
43 Text feature [swung] present in test data point [True]
44 Text feature [months] present in test data point [True]
47 Text feature [phenyl] present in test data point [True]
87 Text feature [fkbp12] present in test data point [True]
95 Text feature [liquid] present in test data point [True]
97 Text feature [benchmark] present in test data point [True]
113 Text feature [childhood1] present in test data point [True]
138 Text feature [fk506] present in test data point [True]
154 Text feature [subgrouping] present in test data point [True]
164 Text feature [clipped] present in test data point [True]
170 Text feature [bsmbi] present in test data point [True]
183 Text feature [optimization] present in test data point [True]
211 Text feature [peptidyl] present in test data point [True]
218 Text feature [bmp4] present in test data point [True]
226 Text feature [wrote] present in test data point [True]
241 Text feature [max] present in test data point [True]
246 Text feature [white] present in test data point [True]
371 Text feature [qdot] present in test data point [True]
382 Text feature [0103] present in test data point [True]
393 Text feature [bre] present in test data point [True]
455 Text feature [mainstay] present in test data point [True]

```

```

460 Text feature [amplifications] present in test data point [True]
483 Text feature [pressure] present in test data point [True]
679 Text feature [ossification] present in test data point [True]
685 Text feature [bmpr] present in test data point [True]
691 Text feature [416] present in test data point [True]
735 Text feature [alk2] present in test data point [True]
753 Text feature [myogenic] present in test data point [True]
806 Text feature [r206h] present in test data point [True]
812 Text feature [time] present in test data point [True]
826 Text feature [qiagen] present in test data point [True]
840 Text feature [paraffin] present in test data point [True]
846 Text feature [milk] present in test data point [True]
855 Text feature [episodes] present in test data point [True]
885 Text feature [hawkins] present in test data point [True]
908 Text feature [skeletal] present in test data point [True]
910 Text feature [regimen] present in test data point [True]
927 Text feature [extended] present in test data point [True]
945 Text feature [refmac5] present in test data point [True]
947 Text feature [v6] present in test data point [True]
959 Text feature [limit] present in test data point [True]
961 Text feature [fatal2] present in test data point [True]
966 Text feature [spss] present in test data point [True]
983 Text feature [myb] present in test data point [True]
997 Text feature [histopathological] present in test data point [True]
Out of the top 1000 features 47 are present in query point

```

4.5 Random Forest Classifier

4.5.1. Hyper paramter tuning (With One hot Encoding)

In [108]:

```

# -----
# 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 : 1.2471948767446823
for n_estimators = 100 and max depth = 10
Log Loss : 1.1963460101138519
for n_estimators = 200 and max depth = 5
Log Loss : 1.2277320725569603
for n_estimators = 200 and max depth = 10
Log Loss : 1.1875378620338506
for n_estimators = 500 and max depth = 5
Log Loss : 1.2256012242967114
for n_estimators = 500 and max depth = 10
Log Loss : 1.183540609943275
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2205011826373053
for n_estimators = 1000 and max depth = 10
Log Loss : 1.1785151254006876
for n_estimators = 2000 and max depth = 5
Log Loss : 1.2164603261276896
for n_estimators = 2000 and max depth = 10
Log Loss : 1.1773659265570025
For values of best estimator = 2000 The train log loss is: 0.6571307143356205
For values of best estimator = 2000 The cross validation log loss is: 1.1773659265570025
For values of best estimator = 2000 The test log loss is: 1.1334958117701477

```

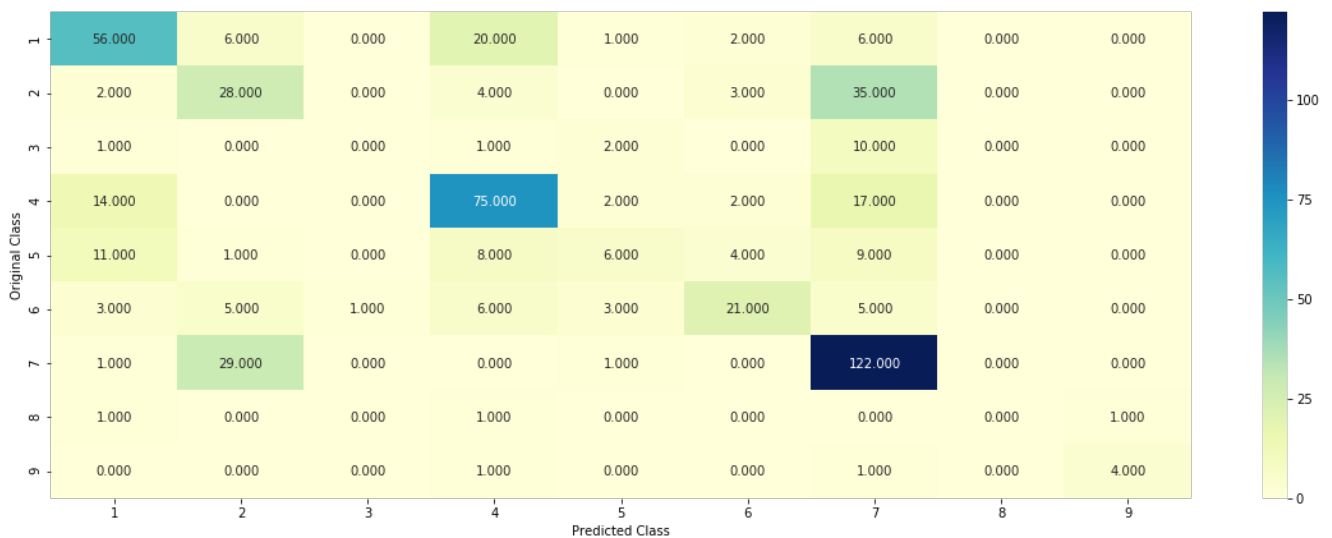
4.5.2. Testing model with best hyper parameters (One Hot Encoding)

In [109]:

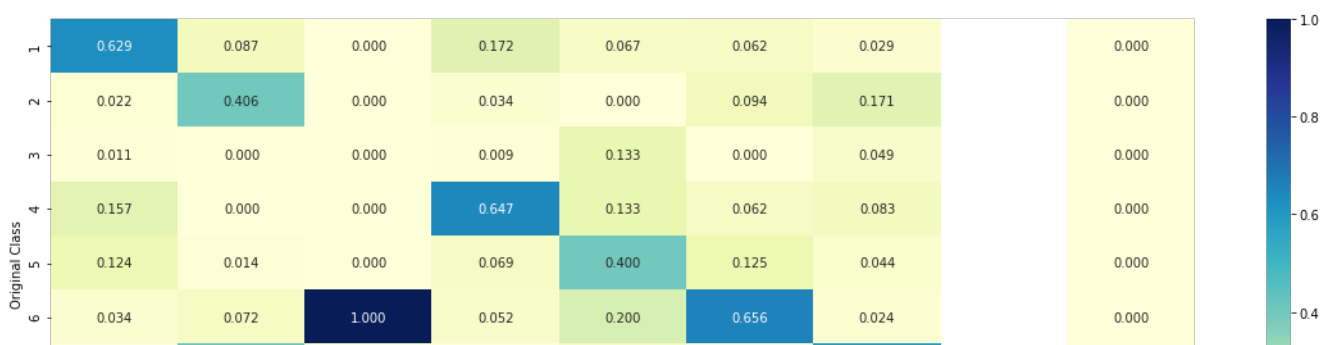
Log loss : 1.1773659265570025

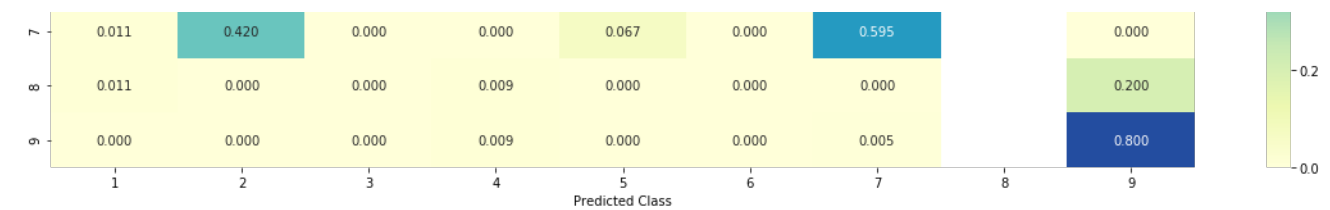
Number of mis-classified points : 0.41353383458646614

```
----- Confusion matrix -----
```

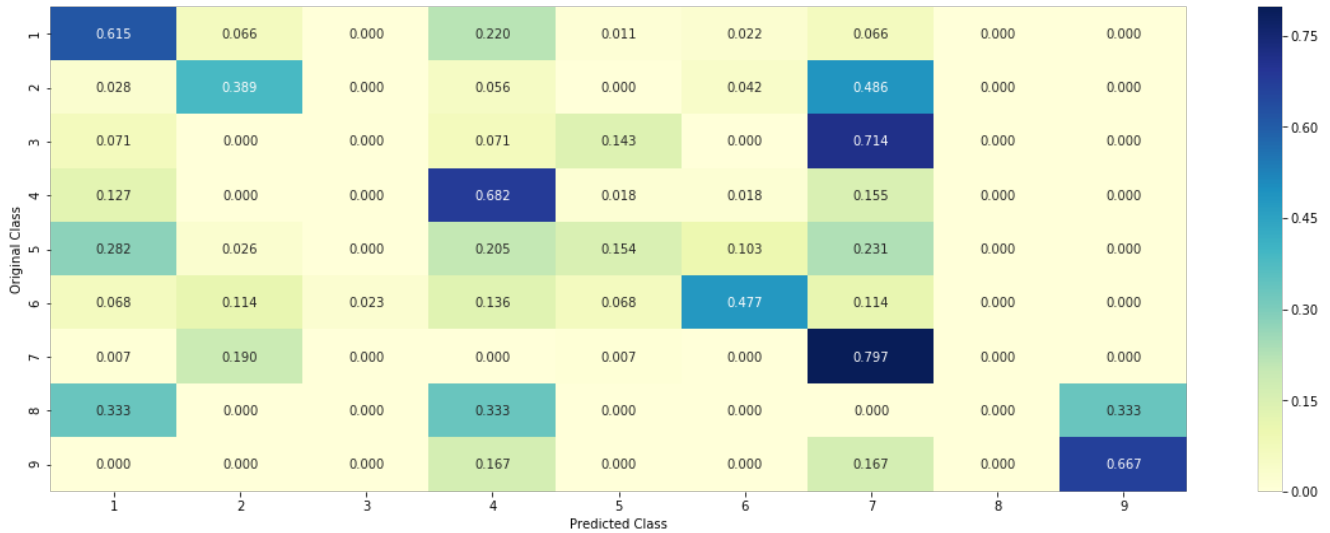


```
----- Precision matrix (Columm Sum=1) -----
```





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



4.5.3. Feature Importance

In []:

```
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)
```

4.5.3.1. Correctly Classified point

In [113]:

```
test_point_index = 10
no_feature = 1000
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],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)
```

Predicted Class : 4

Predicted Class Probabilities: [[0.2387 0.0691 0.0204 0.4616 0.0569 0.0452 0.0941 0.0058 0.0083]]

Actual Class : 4

```
-----
1 Text feature [activating] present in test data point [True]
2 Text feature [activation] present in test data point [True]
6 Text feature [oncogenic] present in test data point [True]
7 Text feature [phosphorylation] present in test data point [True]
9 Text feature [constitutive] present in test data point [True]
10 Text feature [missense] present in test data point [True]
12 Text feature [growth] present in test data point [True]
```


12 Text feature [growth] present in test data point [True]
13 Text feature [signaling] present in test data point [True]
14 Text feature [suppressor] present in test data point [True]
15 Text feature [function] present in test data point [True]
19 Text feature [inhibition] present in test data point [True]
20 Text feature [loss] present in test data point [True]
23 Text feature [proliferation] present in test data point [True]
27 Text feature [therapy] present in test data point [True]
31 Text feature [treated] present in test data point [True]
33 Text feature [months] present in test data point [True]
35 Text feature [expressing] present in test data point [True]
38 Text feature [cell] present in test data point [True]
40 Text feature [functional] present in test data point [True]
41 Text feature [cells] present in test data point [True]
42 Text feature [oncogene] present in test data point [True]
50 Text feature [protein] present in test data point [True]
53 Text feature [resistance] present in test data point [True]
56 Text feature [clinical] present in test data point [True]
57 Text feature [stability] present in test data point [True]
59 Text feature [response] present in test data point [True]
67 Text feature [variant] present in test data point [True]
73 Text feature [lines] present in test data point [True]
81 Text feature [proteins] present in test data point [True]
88 Text feature [null] present in test data point [True]
89 Text feature [expression] present in test data point [True]
93 Text feature [pathway] present in test data point [True]
95 Text feature [factor] present in test data point [True]
96 Text feature [active] present in test data point [True]
97 Text feature [survival] present in test data point [True]
101 Text feature [independent] present in test data point [True]
103 Text feature [il] present in test data point [True]
110 Text feature [functions] present in test data point [True]
111 Text feature [use] present in test data point [True]
115 Text feature [kit] present in test data point [True]
117 Text feature [resistant] present in test data point [True]
118 Text feature [expected] present in test data point [True]
119 Text feature [classify] present in test data point [True]
121 Text feature [sensitive] present in test data point [True]
129 Text feature [acquired] present in test data point [True]
137 Text feature [retained] present in test data point [True]
139 Text feature [tumors] present in test data point [True]
140 Text feature [activity] present in test data point [True]
141 Text feature [patient] present in test data point [True]
148 Text feature [atp] present in test data point [True]
150 Text feature [nuclear] present in test data point [True]
152 Text feature [mutants] present in test data point [True]
156 Text feature [presence] present in test data point [True]
157 Text feature [likely] present in test data point [True]
158 Text feature [affect] present in test data point [True]
162 Text feature [affected] present in test data point [True]
164 Text feature [mechanism] present in test data point [True]
166 Text feature [dna] present in test data point [True]
167 Text feature [sequence] present in test data point [True]
170 Text feature [mutant] present in test data point [True]
173 Text feature [e3] present in test data point [True]
179 Text feature [molecular] present in test data point [True]
181 Text feature [ring] present in test data point [True]
182 Text feature [enhanced] present in test data point [True]
184 Text feature [pathways] present in test data point [True]
187 Text feature [assay] present in test data point [True]
188 Text feature [ability] present in test data point [True]
192 Text feature [folding] present in test data point [True]
197 Text feature [uncertain] present in test data point [True]
203 Text feature [used] present in test data point [True]
206 Text feature [responses] present in test data point [True]
208 Text feature [interaction] present in test data point [True]
209 Text feature [hours] present in test data point [True]
214 Text feature [deletion] present in test data point [True]
216 Text feature [bind] present in test data point [True]
217 Text feature [lung] present in test data point [True]
223 Text feature [model] present in test data point [True]
224 Text feature [contrast] present in test data point [True]
227 Text feature [binding] present in test data point [True]
230 Text feature [21] present in test data point [True]
232 Text feature [primary] present in test data point [True]
235 Text feature [oncogenes] present in test data point [True]
240 Text feature [transfected] present in test data point [True]
245 Text feature [20] present in test data point [True]

245 Text feature [3b] present in test data point [True]
246 Text feature [11] present in test data point [True]
247 Text feature [analysis] present in test data point [True]
248 Text feature [antibodies] present in test data point [True]
253 Text feature [14] present in test data point [True]
255 Text feature [recently] present in test data point [True]
256 Text feature [mutation] present in test data point [True]
257 Text feature [results] present in test data point [True]
262 Text feature [increased] present in test data point [True]
263 Text feature [ligase] present in test data point [True]
264 Text feature [majority] present in test data point [True]
266 Text feature [core] present in test data point [True]
268 Text feature [predictions] present in test data point [True]
274 Text feature [localization] present in test data point [True]
283 Text feature [gain] present in test data point [True]
284 Text feature [substrate] present in test data point [True]
289 Text feature [inactivated] present in test data point [True]
290 Text feature [18] present in test data point [True]
291 Text feature [study] present in test data point [True]
293 Text feature [reduced] present in test data point [True]
294 Text feature [purified] present in test data point [True]
295 Text feature [within] present in test data point [True]
297 Text feature [experiments] present in test data point [True]
298 Text feature [domains] present in test data point [True]
301 Text feature [wild] present in test data point [True]
306 Text feature [site] present in test data point [True]
310 Text feature [gene] present in test data point [True]
311 Text feature [purposes] present in test data point [True]
312 Text feature [human] present in test data point [True]
313 Text feature [target] present in test data point [True]
318 Text feature [absence] present in test data point [True]
319 Text feature [driven] present in test data point [True]
321 Text feature [studies] present in test data point [True]
323 Text feature [vitro] present in test data point [True]
324 Text feature [length] present in test data point [True]
325 Text feature [probability] present in test data point [True]
329 Text feature [levels] present in test data point [True]
334 Text feature [control] present in test data point [True]
335 Text feature [indicated] present in test data point [True]
340 Text feature [one] present in test data point [True]
341 Text feature [subtle] present in test data point [True]
344 Text feature [type] present in test data point [True]
347 Text feature [showed] present in test data point [True]
348 Text feature [genes] present in test data point [True]
349 Text feature [found] present in test data point [True]
353 Text feature [several] present in test data point [True]
355 Text feature [sequencing] present in test data point [True]
357 Text feature [43] present in test data point [True]
360 Text feature [total] present in test data point [True]
362 Text feature [data] present in test data point [True]
364 Text feature [expressed] present in test data point [True]
365 Text feature [addition] present in test data point [True]
366 Text feature [previously] present in test data point [True]
367 Text feature [known] present in test data point [True]
372 Text feature [epidermal] present in test data point [True]
374 Text feature [evidence] present in test data point [True]
375 Text feature [mutations] present in test data point [True]
382 Text feature [lysates] present in test data point [True]
383 Text feature [12] present in test data point [True]
384 Text feature [13] present in test data point [True]
387 Text feature [level] present in test data point [True]
391 Text feature [1a] present in test data point [True]
394 Text feature [co] present in test data point [True]
400 Text feature [anti] present in test data point [True]
403 Text feature [full] present in test data point [True]
404 Text feature [including] present in test data point [True]
406 Text feature [ongoing] present in test data point [True]
408 Text feature [whether] present in test data point [True]
410 Text feature [fell] present in test data point [True]
412 Text feature [tagged] present in test data point [True]
413 Text feature [genetic] present in test data point [True]
416 Text feature [reported] present in test data point [True]
417 Text feature [observed] present in test data point [True]
421 Text feature [role] present in test data point [True]
422 Text feature [cultured] present in test data point [True]
424 Text feature [well] present in test data point [True]
426 Text feature [40] present in test data point [True]
427 Text feature [10] present in test data point [True]

427 Text feature [region] present in test data point [True]
431 Text feature [putative] present in test data point [True]
433 Text feature [shown] present in test data point [True]
435 Text feature [interact] present in test data point [True]
436 Text feature [transcriptional] present in test data point [True]
437 Text feature [also] present in test data point [True]
442 Text feature [second] present in test data point [True]
443 Text feature [approach] present in test data point [True]
444 Text feature [suggesting] present in test data point [True]
445 Text feature [subcellular] present in test data point [True]
446 Text feature [testing] present in test data point [True]
449 Text feature [relative] present in test data point [True]
450 Text feature [24] present in test data point [True]
451 Text feature [tumor] present in test data point [True]
454 Text feature [lost] present in test data point [True]
456 Text feature [viability] present in test data point [True]
457 Text feature [next] present in test data point [True]
459 Text feature [directed] present in test data point [True]
467 Text feature [specific] present in test data point [True]
468 Text feature [similar] present in test data point [True]
469 Text feature [domain] present in test data point [True]
472 Text feature [occur] present in test data point [True]
474 Text feature [given] present in test data point [True]
475 Text feature [repeat] present in test data point [True]
476 Text feature [transcription] present in test data point [True]
480 Text feature [specific] present in test data point [True]
486 Text feature [set] present in test data point [True]
487 Text feature [19] present in test data point [True]
488 Text feature [ubiquitin] present in test data point [True]
489 Text feature [26] present in test data point [True]
492 Text feature [complete] present in test data point [True]
495 Text feature [structural] present in test data point [True]
498 Text feature [based] present in test data point [True]
499 Text feature [confirmed] present in test data point [True]
500 Text feature [antibody] present in test data point [True]
501 Text feature [structure] present in test data point [True]
502 Text feature [although] present in test data point [True]
504 Text feature [analyzed] present in test data point [True]
505 Text feature [western] present in test data point [True]
508 Text feature [ha] present in test data point [True]
510 Text feature [tested] present in test data point [True]
512 Text feature [culture] present in test data point [True]
513 Text feature [mutated] present in test data point [True]
514 Text feature [may] present in test data point [True]
516 Text feature [folded] present in test data point [True]
517 Text feature [discovery] present in test data point [True]
519 Text feature [methods] present in test data point [True]
520 Text feature [interestingly] present in test data point [True]
521 Text feature [2a] present in test data point [True]
527 Text feature [15] present in test data point [True]
528 Text feature [35] present in test data point [True]
529 Text feature [17] present in test data point [True]
533 Text feature [findings] present in test data point [True]
536 Text feature [cancer] present in test data point [True]
537 Text feature [42] present in test data point [True]
538 Text feature [remaining] present in test data point [True]
539 Text feature [incubated] present in test data point [True]
541 Text feature [mutagenesis] present in test data point [True]
546 Text feature [typically] present in test data point [True]
548 Text feature [residue] present in test data point [True]
549 Text feature [substitution] present in test data point [True]
551 Text feature [partial] present in test data point [True]
552 Text feature [25] present in test data point [True]
558 Text feature [discussion] present in test data point [True]
560 Text feature [database] present in test data point [True]
561 Text feature [dependent] present in test data point [True]
562 Text feature [first] present in test data point [True]
563 Text feature [regulatory] present in test data point [True]
567 Text feature [acid] present in test data point [True]
569 Text feature [figure] present in test data point [True]
571 Text feature [introduction] present in test data point [True]
573 Text feature [identified] present in test data point [True]
574 Text feature [table] present in test data point [True]
575 Text feature [mechanisms] present in test data point [True]
578 Text feature [compared] present in test data point [True]
579 Text feature [washed] present in test data point [True]
580 Text feature [10] present in test data point [True]
...

583 Text feature [significance] present in test data point [True]
586 Text feature [possible] present in test data point [True]
587 Text feature [none] present in test data point [True]
591 Text feature [pcr] present in test data point [True]
592 Text feature [using] present in test data point [True]
593 Text feature [regulation] present in test data point [True]
596 Text feature [however] present in test data point [True]
597 Text feature [minutes] present in test data point [True]
602 Text feature [significant] present in test data point [True]
603 Text feature [partially] present in test data point [True]
605 Text feature [common] present in test data point [True]
610 Text feature [consistent] present in test data point [True]
613 Text feature [novel] present in test data point [True]
618 Text feature [targets] present in test data point [True]
620 Text feature [new] present in test data point [True]
625 Text feature [residues] present in test data point [True]
626 Text feature [truncated] present in test data point [True]
627 Text feature [30] present in test data point [True]
628 Text feature [ml] present in test data point [True]
629 Text feature [containing] present in test data point [True]
630 Text feature [degradation] present in test data point [True]
634 Text feature [number] present in test data point [True]
635 Text feature [reporter] present in test data point [True]
636 Text feature [relevant] present in test data point [True]
641 Text feature [demonstrated] present in test data point [True]
646 Text feature [difference] present in test data point [True]
647 Text feature [plasmid] present in test data point [True]
648 Text feature [line] present in test data point [True]
650 Text feature [presented] present in test data point [True]
651 Text feature [via] present in test data point [True]
653 Text feature [defined] present in test data point [True]
654 Text feature [provide] present in test data point [True]
658 Text feature [subset] present in test data point [True]
664 Text feature [phenotypes] present in test data point [True]
668 Text feature [competitive] present in test data point [True]
669 Text feature [need] present in test data point [True]
670 Text feature [somatic] present in test data point [True]
671 Text feature [mediated] present in test data point [True]
672 Text feature [interactions] present in test data point [True]
673 Text feature [confirm] present in test data point [True]
674 Text feature [performed] present in test data point [True]
675 Text feature [system] present in test data point [True]
676 Text feature [involved] present in test data point [True]
677 Text feature [obtained] present in test data point [True]
678 Text feature [38] present in test data point [True]
686 Text feature [20] present in test data point [True]
688 Text feature [suppression] present in test data point [True]
691 Text feature [still] present in test data point [True]
695 Text feature [rate] present in test data point [True]
696 Text feature [measured] present in test data point [True]
698 Text feature [22] present in test data point [True]
700 Text feature [approximately] present in test data point [True]
701 Text feature [better] present in test data point [True]
704 Text feature [3a] present in test data point [True]
708 Text feature [times] present in test data point [True]
709 Text feature [identification] present in test data point [True]
710 Text feature [consequences] present in test data point [True]
712 Text feature [four] present in test data point [True]
715 Text feature [frequently] present in test data point [True]
716 Text feature [analyses] present in test data point [True]
717 Text feature [mutational] present in test data point [True]
718 Text feature [resulting] present in test data point [True]
719 Text feature [tissue] present in test data point [True]
721 Text feature [might] present in test data point [True]
722 Text feature [genomic] present in test data point [True]
723 Text feature [noted] present in test data point [True]
724 Text feature [able] present in test data point [True]
726 Text feature [show] present in test data point [True]
728 Text feature [cycle] present in test data point [True]
729 Text feature [characteristics] present in test data point [True]
730 Text feature [half] present in test data point [True]
731 Text feature [hcl] present in test data point [True]
734 Text feature [revealed] present in test data point [True]
737 Text feature [determine] present in test data point [True]
740 Text feature [sequences] present in test data point [True]
743 Text feature [medium] present in test data point [True]
752 Text feature [mrna] present in test data point [True]

754 Text feature [among] present in test data point [True]
756 Text feature [remains] present in test data point [True]
760 Text feature [16] present in test data point [True]
762 Text feature [across] present in test data point [True]
763 Text feature [acids] present in test data point [True]
764 Text feature [selected] present in test data point [True]
765 Text feature [due] present in test data point [True]
767 Text feature [mek1] present in test data point [True]
768 Text feature [distinct] present in test data point [True]
769 Text feature [reflect] present in test data point [True]
771 Text feature [directly] present in test data point [True]
774 Text feature [cannot] present in test data point [True]
775 Text feature [60] present in test data point [True]
776 Text feature [inactivate] present in test data point [True]
777 Text feature [many] present in test data point [True]
779 Text feature [evaluate] present in test data point [True]
782 Text feature [include] present in test data point [True]
785 Text feature [lack] present in test data point [True]
786 Text feature [23] present in test data point [True]
791 Text feature [elevated] present in test data point [True]
793 Text feature [37] present in test data point [True]
794 Text feature [complex] present in test data point [True]
795 Text feature [range] present in test data point [True]
796 Text feature [multiple] present in test data point [True]
797 Text feature [three] present in test data point [True]
804 Text feature [position] present in test data point [True]
809 Text feature [constructs] present in test data point [True]
813 Text feature [derived] present in test data point [True]
815 Text feature [without] present in test data point [True]
816 Text feature [types] present in test data point [True]
817 Text feature [increase] present in test data point [True]
822 Text feature [promote] present in test data point [True]
826 Text feature [transfection] present in test data point [True]
828 Text feature [least] present in test data point [True]
829 Text feature [per] present in test data point [True]
830 Text feature [endogenous] present in test data point [True]
831 Text feature [two] present in test data point [True]
832 Text feature [fig] present in test data point [True]
834 Text feature [2c] present in test data point [True]
835 Text feature [like] present in test data point [True]
838 Text feature [suggests] present in test data point [True]
839 Text feature [34] present in test data point [True]
841 Text feature [agents] present in test data point [True]
843 Text feature [wt] present in test data point [True]
844 Text feature [fold] present in test data point [True]
846 Text feature [change] present in test data point [True]
849 Text feature [cancers] present in test data point [True]
850 Text feature [unable] present in test data point [True]
856 Text feature [included] present in test data point [True]
857 Text feature [genome] present in test data point [True]
859 Text feature [described] present in test data point [True]
864 Text feature [associated] present in test data point [True]
865 Text feature [together] present in test data point [True]
868 Text feature [upon] present in test data point [True]
870 Text feature [evaluated] present in test data point [True]
871 Text feature [characterized] present in test data point [True]
875 Text feature [proteasome] present in test data point [True]
876 Text feature [frame] present in test data point [True]
879 Text feature [either] present in test data point [True]
881 Text feature [reduce] present in test data point [True]
883 Text feature [thus] present in test data point [True]
885 Text feature [examined] present in test data point [True]
888 Text feature [sample] present in test data point [True]
892 Text feature [mouse] present in test data point [True]
895 Text feature [buffer] present in test data point [True]
898 Text feature [strong] present in test data point [True]
899 Text feature [mapping] present in test data point [True]
900 Text feature [2b] present in test data point [True]
902 Text feature [150] present in test data point [True]
906 Text feature [1b] present in test data point [True]
909 Text feature [distribution] present in test data point [True]
912 Text feature [dtt] present in test data point [True]
914 Text feature [mm] present in test data point [True]
919 Text feature [amino] present in test data point [True]
921 Text feature [63] present in test data point [True]
922 Text feature [established] present in test data point [True]
926 Text feature [association] present in test data point [True]

```

927 Text feature [represent] present in test data point [True]
928 Text feature [work] present in test data point [True]
931 Text feature [observation] present in test data point [True]
935 Text feature [result] present in test data point [True]
937 Text feature [could] present in test data point [True]
941 Text feature [generated] present in test data point [True]
944 Text feature [required] present in test data point [True]
945 Text feature [alter] present in test data point [True]
946 Text feature [form] present in test data point [True]
951 Text feature [phenotype] present in test data point [True]
952 Text feature [respectively] present in test data point [True]
953 Text feature [stably] present in test data point [True]
954 Text feature [comparison] present in test data point [True]
955 Text feature [determined] present in test data point [True]
956 Text feature [recognition] present in test data point [True]
957 Text feature [species] present in test data point [True]
960 Text feature [commonly] present in test data point [True]
961 Text feature [panel] present in test data point [True]
971 Text feature [impaired] present in test data point [True]
975 Text feature [rather] present in test data point [True]
981 Text feature [materials] present in test data point [True]
984 Text feature [occurring] present in test data point [True]
986 Text feature [appropriate] present in test data point [True]
987 Text feature [promoter] present in test data point [True]
988 Text feature [research] present in test data point [True]
994 Text feature [apoptosis] present in test data point [True]
997 Text feature [general] present in test data point [True]
998 Text feature [critical] present in test data point [True]
Out of the top 1000 features 419 are present in query point

```

4.5.3.2. Inorrectly Classified point

In [111]:

```

test_point_index = 15
no_feature = 1000
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actuall Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_impfeature_names(indices[:no_feature],
                      x_test['TEXT'].iloc[test_point_index],
                      x_test['Gene'].iloc[test_point_index],
                      x_test['Variation'].iloc[test_point_index],
                      no_feature)

```

Predicted Class : 7

Predicted Class Probabilities: [[0.0488 0.2297 0.0185 0.0287 0.1008 0.0335 0.5304 0.0052 0.0044]]

Actuall Class : 2

```

-----
0 Text feature [kinase] present in test data point [True]
1 Text feature [activating] present in test data point [True]
2 Text feature [activation] present in test data point [True]
3 Text feature [inhibitors] present in test data point [True]
4 Text feature [inhibitor] present in test data point [True]
5 Text feature [activated] present in test data point [True]
7 Text feature [phosphorylation] present in test data point [True]
8 Text feature [tyrosine] present in test data point [True]
9 Text feature [constitutive] present in test data point [True]
10 Text feature [missense] present in test data point [True]
11 Text feature [erk] present in test data point [True]
12 Text feature [growth] present in test data point [True]
13 Text feature [signaling] present in test data point [True]
14 Text feature [suppressor] present in test data point [True]
16 Text feature [treatment] present in test data point [True]
19 Text feature [inhibition] present in test data point [True]
20 Text feature [loss] present in test data point [True]
22 Text feature [akt] present in test data point [True]
26 Text feature [ba] present in test data point [True]
27 Text feature [therapy] present in test data point [True]
29 Text feature [downstream] present in test data point [True]
31 Text feature [factors] present in test data point [True]

```

31 Text feature [treated] present in test data point [True]
32 Text feature [kinases] present in test data point [True]
33 Text feature [months] present in test data point [True]
35 Text feature [expressing] present in test data point [True]
36 Text feature [trials] present in test data point [True]
37 Text feature [drug] present in test data point [True]
38 Text feature [cell] present in test data point [True]
39 Text feature [therapeutic] present in test data point [True]
40 Text feature [functional] present in test data point [True]
41 Text feature [cells] present in test data point [True]
43 Text feature [mek] present in test data point [True]
47 Text feature [variants] present in test data point [True]
50 Text feature [protein] present in test data point [True]
51 Text feature [deleterious] present in test data point [True]
52 Text feature [f3] present in test data point [True]
55 Text feature [neutral] present in test data point [True]
56 Text feature [clinical] present in test data point [True]
59 Text feature [response] present in test data point [True]
62 Text feature [patients] present in test data point [True]
63 Text feature [brca1] present in test data point [True]
67 Text feature [variant] present in test data point [True]
70 Text feature [classified] present in test data point [True]
71 Text feature [ras] present in test data point [True]
73 Text feature [lines] present in test data point [True]
74 Text feature [efficacy] present in test data point [True]
76 Text feature [nsc1c] present in test data point [True]
77 Text feature [dose] present in test data point [True]
78 Text feature [amplification] present in test data point [True]
79 Text feature [potential] present in test data point [True]
80 Text feature [serum] present in test data point [True]
81 Text feature [proteins] present in test data point [True]
83 Text feature [harboring] present in test data point [True]
84 Text feature [pathogenicity] present in test data point [True]
86 Text feature [daily] present in test data point [True]
89 Text feature [expression] present in test data point [True]
93 Text feature [pathway] present in test data point [True]
100 Text feature [pten] present in test data point [True]
101 Text feature [independent] present in test data point [True]
103 Text feature [il] present in test data point [True]
108 Text feature [brca2] present in test data point [True]
112 Text feature [effective] present in test data point [True]
115 Text feature [kit] present in test data point [True]
121 Text feature [sensitive] present in test data point [True]
122 Text feature [pi3k] present in test data point [True]
124 Text feature [mapk] present in test data point [True]
127 Text feature [combined] present in test data point [True]
139 Text feature [tumors] present in test data point [True]
140 Text feature [activity] present in test data point [True]
141 Text feature [patient] present in test data point [True]
152 Text feature [mutants] present in test data point [True]
156 Text feature [presence] present in test data point [True]
157 Text feature [likely] present in test data point [True]
159 Text feature [benefit] present in test data point [True]
162 Text feature [affected] present in test data point [True]
165 Text feature [murine] present in test data point [True]
166 Text feature [dna] present in test data point [True]
167 Text feature [sequence] present in test data point [True]
170 Text feature [mutant] present in test data point [True]
178 Text feature [median] present in test data point [True]
179 Text feature [molecular] present in test data point [True]
184 Text feature [pathways] present in test data point [True]
187 Text feature [assay] present in test data point [True]
188 Text feature [ability] present in test data point [True]
189 Text feature [available] present in test data point [True]
190 Text feature [raf] present in test data point [True]
193 Text feature [therapies] present in test data point [True]
195 Text feature [classification] present in test data point [True]
203 Text feature [used] present in test data point [True]
206 Text feature [responses] present in test data point [True]
207 Text feature [hybridization] present in test data point [True]
211 Text feature [weeks] present in test data point [True]
215 Text feature [biopsy] present in test data point [True]
217 Text feature [lung] present in test data point [True]
218 Text feature [clinically] present in test data point [True]
219 Text feature [terminal] present in test data point [True]
224 Text feature [contrast] present in test data point [True]
226 Text feature [64] present in test data point [True]
227 Text feature [64] present in test data point [True]

227 Text feature [binding] present in test data point [True]
230 Text feature [21] present in test data point [True]
241 Text feature [plasma] present in test data point [True]
245 Text feature [3b] present in test data point [True]
246 Text feature [11] present in test data point [True]
247 Text feature [analysis] present in test data point [True]
253 Text feature [14] present in test data point [True]
254 Text feature [strand] present in test data point [True]
255 Text feature [recently] present in test data point [True]
256 Text feature [mutation] present in test data point [True]
257 Text feature [results] present in test data point [True]
264 Text feature [majority] present in test data point [True]
265 Text feature [small] present in test data point [True]
266 Text feature [core] present in test data point [True]
277 Text feature [iarc] present in test data point [True]
278 Text feature [harbored] present in test data point [True]
285 Text feature [alk] present in test data point [True]
287 Text feature [refractory] present in test data point [True]
290 Text feature [18] present in test data point [True]
291 Text feature [study] present in test data point [True]
294 Text feature [purified] present in test data point [True]
295 Text feature [within] present in test data point [True]
297 Text feature [experiments] present in test data point [True]
299 Text feature [case] present in test data point [True]
301 Text feature [wild] present in test data point [True]
306 Text feature [site] present in test data point [True]
310 Text feature [gene] present in test data point [True]
312 Text feature [human] present in test data point [True]
313 Text feature [target] present in test data point [True]
314 Text feature [nhgri] present in test data point [True]
315 Text feature [vector] present in test data point [True]
316 Text feature [50] present in test data point [True]
319 Text feature [driven] present in test data point [True]
321 Text feature [studies] present in test data point [True]
322 Text feature [individuals] present in test data point [True]
323 Text feature [vitro] present in test data point [True]
324 Text feature [length] present in test data point [True]
327 Text feature [assessment] present in test data point [True]
329 Text feature [levels] present in test data point [True]
334 Text feature [control] present in test data point [True]
340 Text feature [one] present in test data point [True]
343 Text feature [controls] present in test data point [True]
344 Text feature [type] present in test data point [True]
348 Text feature [genes] present in test data point [True]
349 Text feature [found] present in test data point [True]
351 Text feature [paraffin] present in test data point [True]
355 Text feature [sequencing] present in test data point [True]
357 Text feature [43] present in test data point [True]
360 Text feature [total] present in test data point [True]
362 Text feature [data] present in test data point [True]
364 Text feature [expressed] present in test data point [True]
365 Text feature [addition] present in test data point [True]
366 Text feature [previously] present in test data point [True]
367 Text feature [known] present in test data point [True]
370 Text feature [old] present in test data point [True]
371 Text feature [disease] present in test data point [True]
373 Text feature [highly] present in test data point [True]
374 Text feature [evidence] present in test data point [True]
375 Text feature [mutations] present in test data point [True]
377 Text feature [32] present in test data point [True]
379 Text feature [harbor] present in test data point [True]
383 Text feature [12] present in test data point [True]
384 Text feature [13] present in test data point [True]
391 Text feature [1a] present in test data point [True]
394 Text feature [co] present in test data point [True]
403 Text feature [full] present in test data point [True]
404 Text feature [including] present in test data point [True]
405 Text feature [stat3] present in test data point [True]
406 Text feature [ongoing] present in test data point [True]
409 Text feature [group] present in test data point [True]
411 Text feature [currently] present in test data point [True]
413 Text feature [genetic] present in test data point [True]
415 Text feature [truncation] present in test data point [True]
416 Text feature [reported] present in test data point [True]
418 Text feature [qiagen] present in test data point [True]
419 Text feature [review] present in test data point [True]
421 Text feature [role] present in test data point [True]
...

423 Text feature [basis] present in test data point [True]
424 Text feature [well] present in test data point [True]
426 Text feature [40] present in test data point [True]
429 Text feature [year] present in test data point [True]
434 Text feature [fusion] present in test data point [True]
436 Text feature [transcriptional] present in test data point [True]
437 Text feature [also] present in test data point [True]
438 Text feature [institutional] present in test data point [True]
439 Text feature [exon] present in test data point [True]
442 Text feature [second] present in test data point [True]
443 Text feature [approach] present in test data point [True]
449 Text feature [relative] present in test data point [True]
450 Text feature [24] present in test data point [True]
451 Text feature [tumor] present in test data point [True]
452 Text feature [individual] present in test data point [True]
453 Text feature [identify] present in test data point [True]
455 Text feature [developed] present in test data point [True]
456 Text feature [viability] present in test data point [True]
457 Text feature [next] present in test data point [True]
464 Text feature [predisposition] present in test data point [True]
467 Text feature [suggest] present in test data point [True]
468 Text feature [similar] present in test data point [True]
469 Text feature [domain] present in test data point [True]
474 Text feature [given] present in test data point [True]
476 Text feature [transcription] present in test data point [True]
477 Text feature [kras] present in test data point [True]
480 Text feature [specific] present in test data point [True]
481 Text feature [days] present in test data point [True]
482 Text feature [46] present in test data point [True]
484 Text feature [example] present in test data point [True]
486 Text feature [set] present in test data point [True]
487 Text feature [19] present in test data point [True]
489 Text feature [26] present in test data point [True]
490 Text feature [certain] present in test data point [True]
491 Text feature [sustained] present in test data point [True]
492 Text feature [complete] present in test data point [True]
494 Text feature [overall] present in test data point [True]
498 Text feature [based] present in test data point [True]
499 Text feature [confirmed] present in test data point [True]
502 Text feature [although] present in test data point [True]
503 Text feature [important] present in test data point [True]
504 Text feature [analyzed] present in test data point [True]
506 Text feature [33] present in test data point [True]
507 Text feature [detection] present in test data point [True]
511 Text feature [biopsies] present in test data point [True]
513 Text feature [mutated] present in test data point [True]
514 Text feature [may] present in test data point [True]
517 Text feature [discovery] present in test data point [True]
518 Text feature [progressed] present in test data point [True]
519 Text feature [methods] present in test data point [True]
520 Text feature [interestingly] present in test data point [True]
521 Text feature [2a] present in test data point [True]
527 Text feature [15] present in test data point [True]
529 Text feature [17] present in test data point [True]
533 Text feature [findings] present in test data point [True]
535 Text feature [cases] present in test data point [True]
536 Text feature [cancer] present in test data point [True]
537 Text feature [42] present in test data point [True]
540 Text feature [regions] present in test data point [True]
543 Text feature [samples] present in test data point [True]
547 Text feature [members] present in test data point [True]
552 Text feature [25] present in test data point [True]
555 Text feature [institute] present in test data point [True]
556 Text feature [lipid] present in test data point [True]
558 Text feature [discussion] present in test data point [True]
560 Text feature [database] present in test data point [True]
561 Text feature [dependent] present in test data point [True]
562 Text feature [first] present in test data point [True]
563 Text feature [regulatory] present in test data point [True]
566 Text feature [present] present in test data point [True]
567 Text feature [acid] present in test data point [True]
568 Text feature [received] present in test data point [True]
569 Text feature [figure] present in test data point [True]
571 Text feature [introduction] present in test data point [True]
572 Text feature [braf] present in test data point [True]
573 Text feature [identified] present in test data point [True]
574 Text feature [table] present in test data point [True]

577 Text feature [31] present in test data point [True]
578 Text feature [compared] present in test data point [True]
580 Text feature [10] present in test data point [True]
581 Text feature [49] present in test data point [True]
582 Text feature [specimens] present in test data point [True]
583 Text feature [significance] present in test data point [True]
586 Text feature [possible] present in test data point [True]
588 Text feature [pretreatment] present in test data point [True]
590 Text feature [prediction] present in test data point [True]
591 Text feature [pcr] present in test data point [True]
592 Text feature [using] present in test data point [True]
593 Text feature [regulation] present in test data point [True]
595 Text feature [28] present in test data point [True]
596 Text feature [however] present in test data point [True]
597 Text feature [minutes] present in test data point [True]
602 Text feature [significant] present in test data point [True]
605 Text feature [common] present in test data point [True]
608 Text feature [forms] present in test data point [True]
610 Text feature [consistent] present in test data point [True]
613 Text feature [novel] present in test data point [True]
617 Text feature [final] present in test data point [True]
618 Text feature [targets] present in test data point [True]
620 Text feature [new] present in test data point [True]
622 Text feature [investigated] present in test data point [True]
623 Text feature [step] present in test data point [True]
627 Text feature [30] present in test data point [True]
628 Text feature [ml] present in test data point [True]
631 Text feature [related] present in test data point [True]
632 Text feature [fish] present in test data point [True]
634 Text feature [number] present in test data point [True]
641 Text feature [demonstrated] present in test data point [True]
650 Text feature [presented] present in test data point [True]
653 Text feature [defined] present in test data point [True]
654 Text feature [provide] present in test data point [True]
656 Text feature [another] present in test data point [True]
657 Text feature [53] present in test data point [True]
659 Text feature [according] present in test data point [True]
663 Text feature [detected] present in test data point [True]
669 Text feature [need] present in test data point [True]
670 Text feature [somatic] present in test data point [True]
673 Text feature [confirm] present in test data point [True]
674 Text feature [performed] present in test data point [True]
675 Text feature [system] present in test data point [True]
676 Text feature [involved] present in test data point [True]
677 Text feature [obtained] present in test data point [True]
678 Text feature [38] present in test data point [True]
679 Text feature [sequenced] present in test data point [True]
682 Text feature [dominant] present in test data point [True]
686 Text feature [20] present in test data point [True]
687 Text feature [calculated] present in test data point [True]
689 Text feature [less] present in test data point [True]
692 Text feature [average] present in test data point [True]
693 Text feature [98] present in test data point [True]
694 Text feature [population] present in test data point [True]
695 Text feature [rate] present in test data point [True]
698 Text feature [22] present in test data point [True]
700 Text feature [approximately] present in test data point [True]
704 Text feature [3a] present in test data point [True]
709 Text feature [identification] present in test data point [True]
710 Text feature [consequences] present in test data point [True]
716 Text feature [analyses] present in test data point [True]
717 Text feature [mutational] present in test data point [True]
718 Text feature [resulting] present in test data point [True]
719 Text feature [tissue] present in test data point [True]
721 Text feature [might] present in test data point [True]
722 Text feature [genomic] present in test data point [True]
723 Text feature [noted] present in test data point [True]
724 Text feature [able] present in test data point [True]
728 Text feature [cycle] present in test data point [True]
729 Text feature [characteristics] present in test data point [True]
732 Text feature [sites] present in test data point [True]
733 Text feature [frozen] present in test data point [True]
734 Text feature [revealed] present in test data point [True]
735 Text feature [effects] present in test data point [True]
736 Text feature [additional] present in test data point [True]
737 Text feature [determine] present in test data point [True]
740 Text feature [sequences] present in test data point [True]

746 Text feature [29] present in test data point [True]
748 Text feature [value] present in test data point [True]
750 Text feature [allele] present in test data point [True]
752 Text feature [mrna] present in test data point [True]
760 Text feature [16] present in test data point [True]
762 Text feature [across] present in test data point [True]
764 Text feature [selected] present in test data point [True]
765 Text feature [due] present in test data point [True]
768 Text feature [distinct] present in test data point [True]
771 Text feature [directly] present in test data point [True]
775 Text feature [60] present in test data point [True]
781 Text feature [furthermore] present in test data point [True]
784 Text feature [27] present in test data point [True]
786 Text feature [23] present in test data point [True]
787 Text feature [primers] present in test data point [True]
792 Text feature [previous] present in test data point [True]
793 Text feature [37] present in test data point [True]
795 Text feature [range] present in test data point [True]
796 Text feature [multiple] present in test data point [True]
797 Text feature [three] present in test data point [True]
800 Text feature [amplified] present in test data point [True]
801 Text feature [tables] present in test data point [True]
808 Text feature [report] present in test data point [True]
809 Text feature [constructs] present in test data point [True]
813 Text feature [derived] present in test data point [True]
815 Text feature [without] present in test data point [True]
816 Text feature [types] present in test data point [True]
821 Text feature [frequent] present in test data point [True]
824 Text feature [single] present in test data point [True]
829 Text feature [per] present in test data point [True]
831 Text feature [two] present in test data point [True]
832 Text feature [fig] present in test data point [True]
834 Text feature [2c] present in test data point [True]
836 Text feature [routine] present in test data point [True]
837 Text feature [rapid] present in test data point [True]
839 Text feature [34] present in test data point [True]
841 Text feature [agents] present in test data point [True]
842 Text feature [55] present in test data point [True]
843 Text feature [wt] present in test data point [True]
851 Text feature [47] present in test data point [True]
852 Text feature [3d] present in test data point [True]
853 Text feature [fh] present in test data point [True]
856 Text feature [included] present in test data point [True]
857 Text feature [genome] present in test data point [True]
859 Text feature [described] present in test data point [True]
860 Text feature [exons] present in test data point [True]
862 Text feature [version] present in test data point [True]
864 Text feature [associated] present in test data point [True]
868 Text feature [upon] present in test data point [True]
870 Text feature [evaluated] present in test data point [True]
871 Text feature [characterized] present in test data point [True]
876 Text feature [frame] present in test data point [True]
878 Text feature [account] present in test data point [True]
879 Text feature [either] present in test data point [True]
883 Text feature [thus] present in test data point [True]
885 Text feature [examined] present in test data point [True]
888 Text feature [sample] present in test data point [True]
890 Text feature [flanking] present in test data point [True]
893 Text feature [event] present in test data point [True]
898 Text feature [strong] present in test data point [True]
899 Text feature [mapping] present in test data point [True]
902 Text feature [150] present in test data point [True]
906 Text feature [1b] present in test data point [True]
909 Text feature [distribution] present in test data point [True]
911 Text feature [indicating] present in test data point [True]
913 Text feature [agent] present in test data point [True]
915 Text feature [polymerase] present in test data point [True]
916 Text feature [selection] present in test data point [True]
918 Text feature [extracted] present in test data point [True]
920 Text feature [100] present in test data point [True]
927 Text feature [represent] present in test data point [True]
929 Text feature [occurrence] present in test data point [True]
930 Text feature [000] present in test data point [True]
935 Text feature [result] present in test data point [True]
936 Text feature [36] present in test data point [True]
938 Text feature [reports] present in test data point [True]
942 Text feature [rna] present in test data point [True]

```

944 Text feature [required] present in test data point [True]
949 Text feature [diagnosis] present in test data point [True]
952 Text feature [respectively] present in test data point [True]
953 Text feature [stably] present in test data point [True]
954 Text feature [comparison] present in test data point [True]
961 Text feature [panel] present in test data point [True]
964 Text feature [size] present in test data point [True]
965 Text feature [94] present in test data point [True]
967 Text feature [empty] present in test data point [True]
972 Text feature [regression] present in test data point [True]
973 Text feature [targeting] present in test data point [True]
974 Text feature [45] present in test data point [True]
976 Text feature [since] present in test data point [True]
982 Text feature [considered] present in test data point [True]
984 Text feature [occurring] present in test data point [True]
986 Text feature [appropriate] present in test data point [True]
987 Text feature [promoter] present in test data point [True]
988 Text feature [research] present in test data point [True]
993 Text feature [address] present in test data point [True]
Out of the top 1000 features 425 are present in query point

```

Random Forest

4.5.3. Hyper paramter tuning (With Response Coding)

In [114]:

```

# -----
# 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-forests-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))
'''
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 : 2.1914634908646633
for n_estimators = 10 and max depth = 3
Log Loss : 1.814854557293049
for n_estimators = 10 and max depth = 5
Log Loss : 1.6934378563564338
for n_estimators = 10 and max depth = 10
Log Loss : 1.5701267912497772
for n_estimators = 50 and max depth = 2
Log Loss : 1.7673885193948395
for n_estimators = 50 and max depth = 3
Log Loss : 1.5237879636990221
for n_estimators = 50 and max depth = 5
Log Loss : 1.4272678347121501
for n_estimators = 50 and max depth = 10
Log Loss : 1.582095678039206
for n_estimators = 100 and max depth = 2
Log Loss : 1.5626563305034231
for n_estimators = 100 and max depth = 3
Log Loss : 1.4897072688726447
for n_estimators = 100 and max depth = 5
Log Loss : 1.3731624240113367
for n_estimators = 100 and max depth = 10
Log Loss : 1.5867130269500642
for n_estimators = 200 and max depth = 2
Log Loss : 1.5866307458292317
for n_estimators = 200 and max depth = 3
Log Loss : 1.4952439171466034
for n_estimators = 200 and max depth = 5
Log Loss : 1.3905987800902242

```

```

Log Loss : 1.59059978000602343
for n_estimators = 200 and max depth = 10
Log Loss : 1.6680545399543352
for n_estimators = 500 and max depth = 2
Log Loss : 1.6600233011525787
for n_estimators = 500 and max depth = 3
Log Loss : 1.5578060598530976
for n_estimators = 500 and max depth = 5
Log Loss : 1.3853471757148181
for n_estimators = 500 and max depth = 10
Log Loss : 1.723532391497242
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6598440797424854
for n_estimators = 1000 and max depth = 3
Log Loss : 1.574917995146402
for n_estimators = 1000 and max depth = 5
Log Loss : 1.3952537809813668
for n_estimators = 1000 and max depth = 10
Log Loss : 1.744862473011676
For values of best alpha = 100 The train log loss is: 0.0591896777045458
For values of best alpha = 100 The cross validation log loss is: 1.3731624240113367
For values of best alpha = 100 The test log loss is: 1.309827971558602

```

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

In [115]:

```

# -----
# 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).

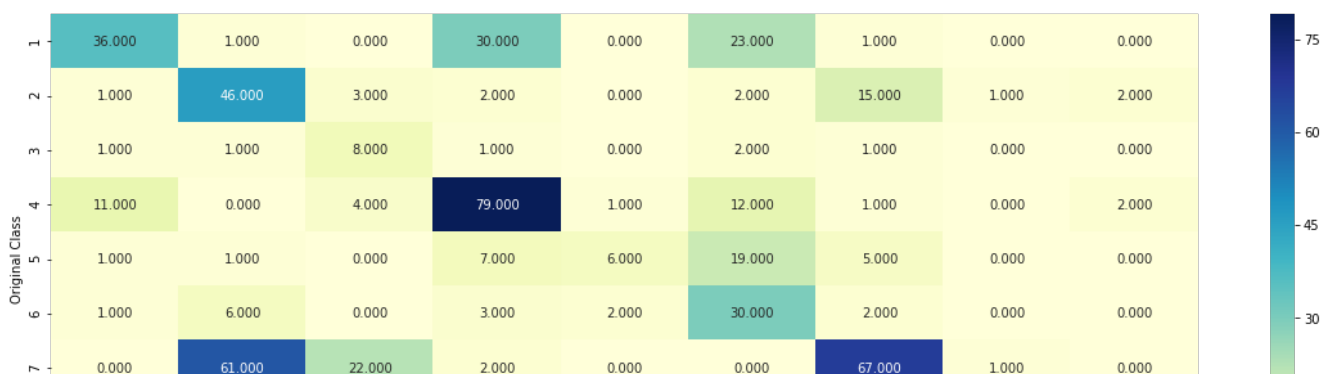
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-fores
t-and-their-construction-2/
# -----

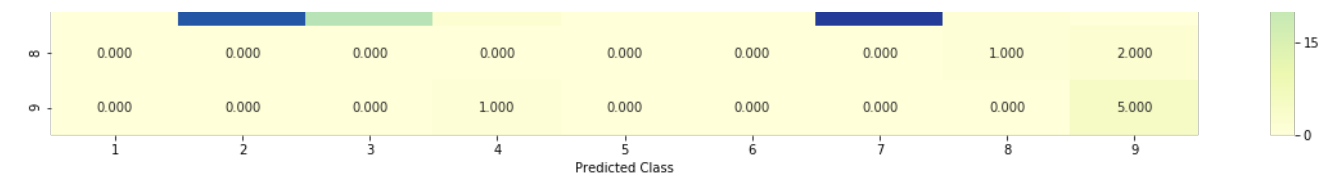
clf = RandomForestClassifier(max_depth=max_depth[int(best_alpha%4)],
n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_features='auto', random_state=42)
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y,cv_x_responseCoding,cv_y, clf)

```

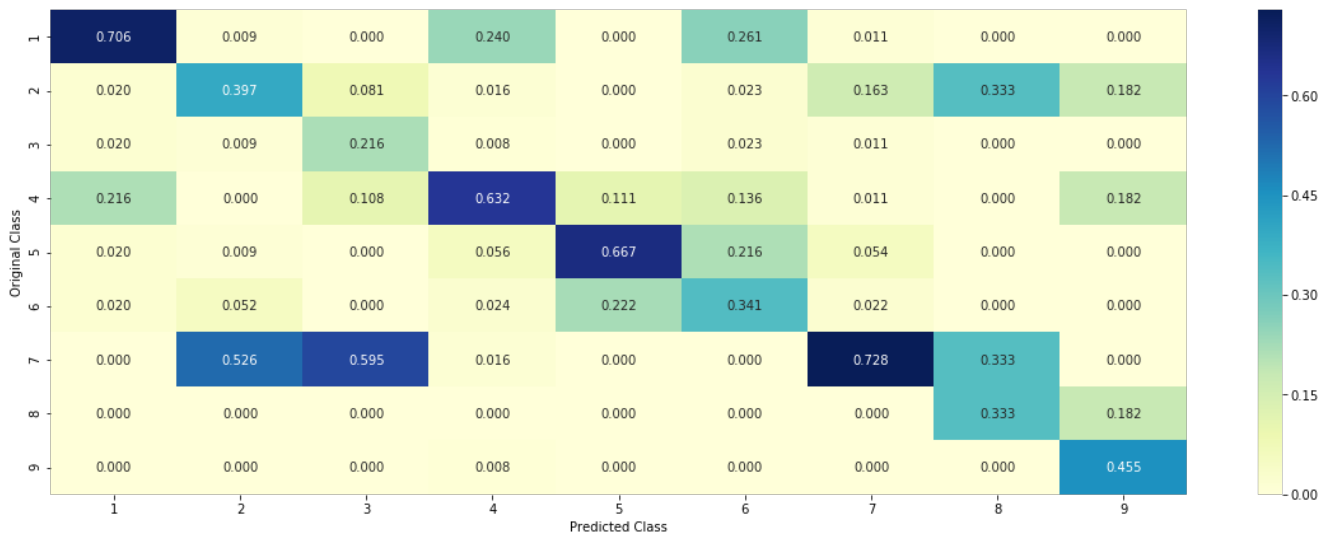
Log loss : 1.373162424011337
Number of mis-classified points : 0.4774436090225564

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

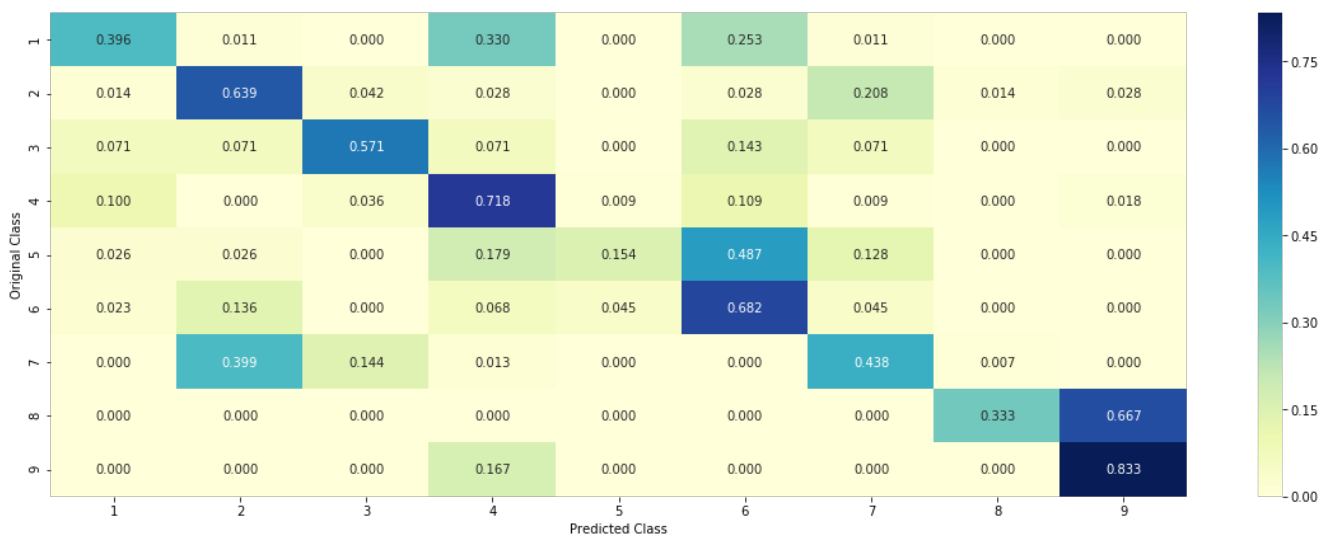




----- Precision matrix (Column Sum=1) -----



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



4.5.5. Feature Importance

4.5.5.1. Correctly Classified point

In [116]:

```
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 = 100
no_feature = 1000
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 : 1

Predicted Class Probabilities: [[0.9762 0.0014 0.0016 0.0099 0.0011 0.0034 0.0017 0.0028 0.0018]]

Actual Class : 1

4.5.5.2. Incorrectly Classified point

In [117]:

```

test_point_index = 31
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 : 2

Predicted Class Probabilities: [[0.029 0.5391 0.0409 0.0357 0.0142 0.0419 0.0333 0.2027 0.0631]]

Actual Class : 7

4.7 Stack the models

4.7.1 testing with hyper parameter tuning

In [118]:

```

# 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/
#-----

# 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/
# -----

# 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=0.01, 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=1000)
clf3.fit(train_x_onehotCoding, train_y)
sig_clf3 = CalibratedClassifierCV(clf3, method="sigmoid")

sig_clf1.fit(train_x_onehotCoding, train_y)
print("Logistic Regression : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf1.predict_proba(cv_x_onehotCoding))))
sig_clf2.fit(train_x_onehotCoding, train_y)
print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf2.predict_proba(cv_x_onehotCoding))))
sig_clf3.fit(train_x_onehotCoding, train_y)
print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predict_proba(cv_x_onehotCoding))))
print("-"*50)
alpha = [0.0001, 0.001, 0.01, 0.1, 1, 10]
best_alpha = 999
for i in alpha:
    lr = LogisticRegression(C=i)
    sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_probas=True)
    sclf.fit(train_x_onehotCoding, train_y)
    print("Stacking 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: 1.14  
Support vector machines : Log Loss: 1.21  
Naive Bayes : Log Loss: 1.25
```

```
-----  
Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.173  
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 1.997  
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.420  
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.113  
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.203  
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.386
```

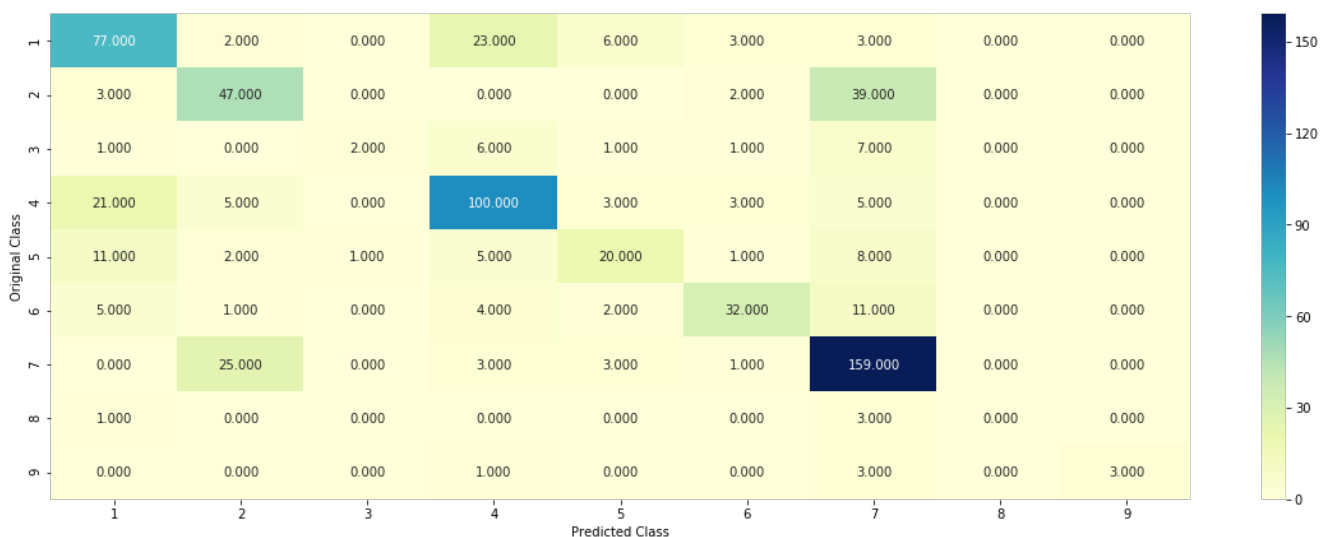
4.7.2 testing the model with the best hyper parameters

In [119]:

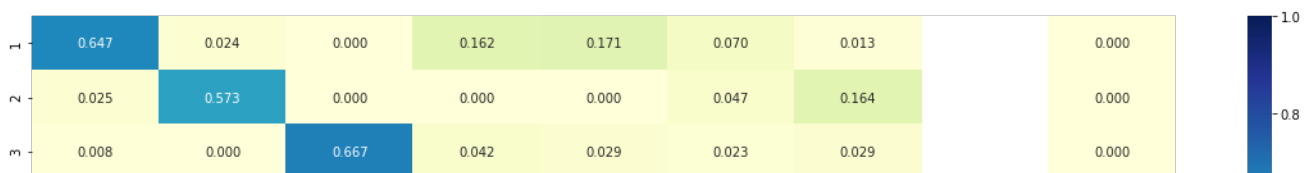
```
lr = LogisticRegression(C=0.1)  
scf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_proba=  
s=True)  
scf.fit(train_x_onehotCoding, train_y)  
  
log_error = log_loss(train_y, scf.predict_proba(train_x_onehotCoding))  
print("Log loss (train) on the stacking classifier :",log_error)  
  
log_error = log_loss(cv_y, scf.predict_proba(cv_x_onehotCoding))  
print("Log loss (CV) on the stacking classifier :",log_error)  
  
log_error = log_loss(test_y, scf.predict_proba(test_x_onehotCoding))  
print("Log loss (test) on the stacking classifier :",log_error)  
  
print("Number of missclassified point :", np.count_nonzero((scf.predict(test_x_onehotCoding)-  
test_y))/test_y.shape[0])  
plot_confusion_matrix(test_y=test_y, predict_y=scf.predict(test_x_onehotCoding))
```

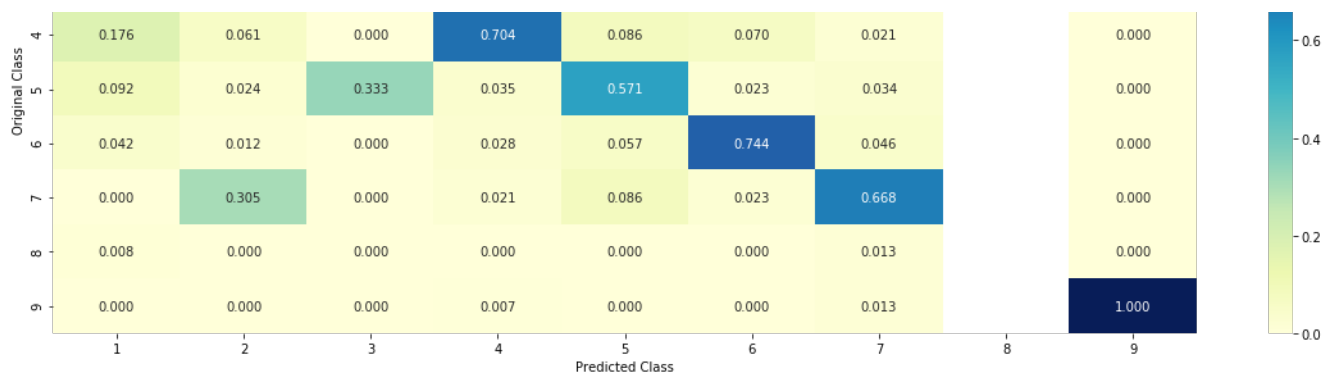
```
Log loss (train) on the stacking classifier : 0.6270217314244309  
Log loss (CV) on the stacking classifier : 1.1128146028708639  
Log loss (test) on the stacking classifier : 1.0738888234178559  
Number of missclassified point : 0.3383458646616541
```

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

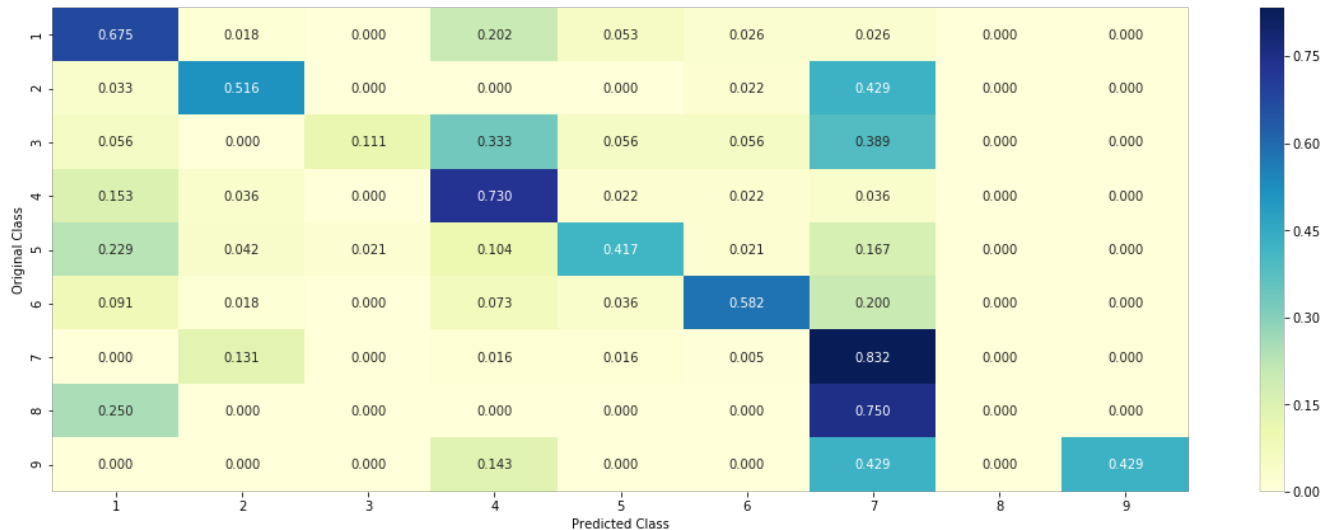


----- Precision matrix (Column Sum=1) -----





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



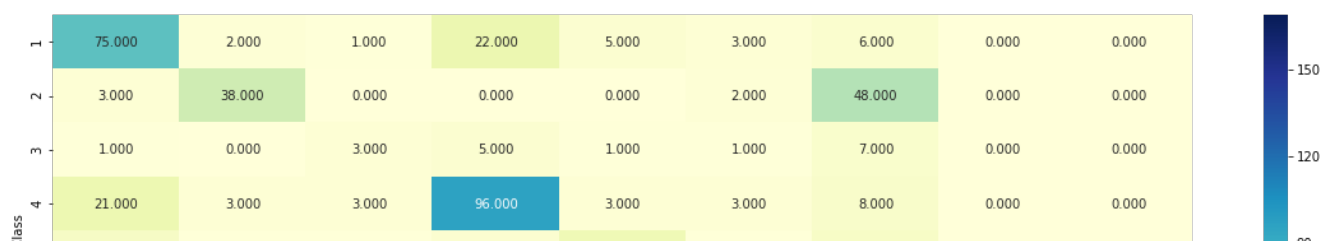
4.7.3 Maximum Voting classifier

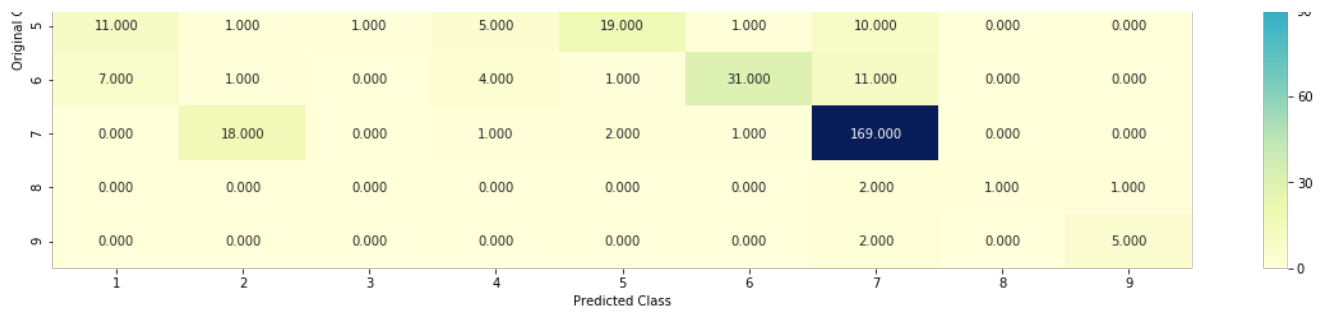
In [120]:

```
#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))
```

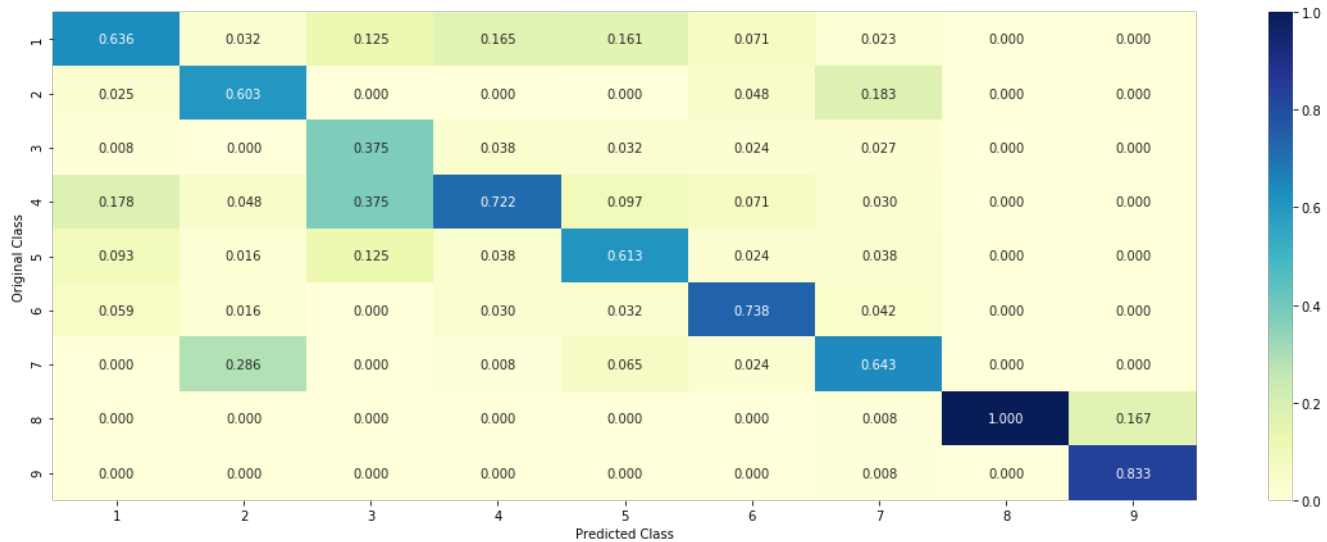
Log loss (train) on the VotingClassifier : 0.7036953034241531
Log loss (CV) on the VotingClassifier : 1.0944254902631223
Log loss (test) on the VotingClassifier : 1.0286822528052888
Number of missclassified point : 0.34285714285714286

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

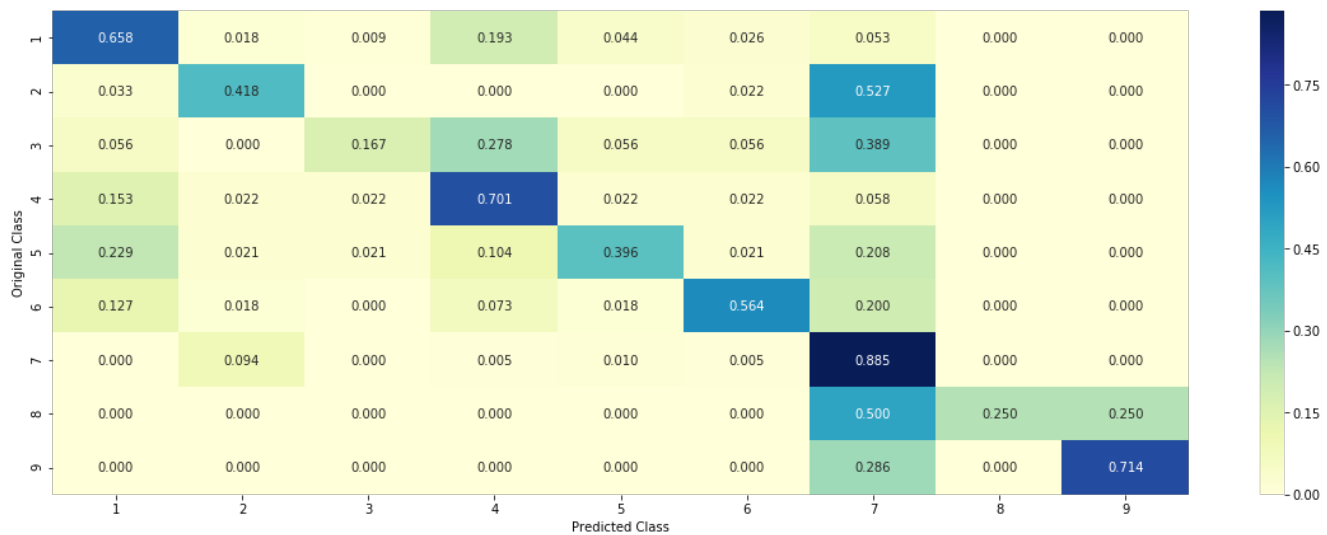




Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Lets summarize above models before proceeding with the feature engineering approach.

In [207]:

```
from prettytable import PrettyTable
ptable = PrettyTable()
ptable.title = "*** Model Summary *** [Performance Metric: Log-Loss]"
ptable.field_names=["Model Name","Train","CV","Test","% Misclassified Points"]
ptable.add_row(["Naive Bayes","0.92","1.24","1.17","40"])
ptable.add_row(["KNN","0.64","1.07","1.01","38"])
ptable.add_row(["Logistic Regression With Class balancing","0.58","1.14","1.07","36"])
ptable.add_row(["Logistic Regression Without Class balancing","0.57","1.15","1.09","36"])
```

```

ptable.add_row(["Linear SVM","0.71","1.18","1.11","36"])
ptable.add_row(["Random Forest Classifier With One hot Encoding","0.65","1.17","1.13","41"])
ptable.add_row(["Random Forest Classifier With Response Coding","0.05","1.37","1.31","48"])
ptable.add_row(["Stack Models:LR+NB+SVM","0.63","1.11","1.07","33"])
ptable.add_row(["Maximum Voting classifier","0.70","1.09","1.03","34"])
print(ptable)
print()

```

Model Name	Train	CV	Test	% Misclassified Points
Naive Bayes	0.92	1.24	1.17	40
KNN	0.64	1.07	1.01	38
Logistic Regression With Class balancing	0.58	1.14	1.07	36
Logistic Regression Without Class balancing	0.57	1.15	1.09	36
Linear SVM	0.71	1.18	1.11	36
Random Forest Classifier With One hot Encoding	0.65	1.17	1.13	41
Random Forest Classifier With Response Coding	0.05	1.37	1.31	48
Stack Models:LR+NB+SVM	0.63	1.11	1.07	33
Maximum Voting classifier	0.70	1.09	1.03	34

Summary:

From Pretty table we can say that 'Logistic Regression With Class balancing' is best choice . So, in further tasks we will use LR with Class Balancing as out model

Task 2:

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

Logistic Regression With Class Balancing

Gene Feature

In [122]:

```

#response-coding of the Gene feature
# alpha is used for laplace smoothing
alpha = 1

# train gene feature
train_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", x_train))

# test gene feature
test_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", x_test))

# cross validation gene feature
cv_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", x_cv))

```

In [123]:

```

# one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer(ngram_range=(1, 2))
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(x_train['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(x_test['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(x_cv['Gene'])

# don't forget to normalize every feature
train_gene_feature_onehotCoding = normalize(train_gene_feature_onehotCoding, axis=0)
test_gene_feature_onehotCoding = normalize(test_gene_feature_onehotCoding, axis=0)
cv_gene_feature_onehotCoding = normalize(cv_gene_feature_onehotCoding, axis=0)

```

Variation Feature

In [124]:

```
# alpha is used for laplace smoothing
alpha = 1

# train gene feature
train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_train))

# test gene feature
test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_test))

# cross validation gene feature
cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_cv))
```

In [125]:

```
# one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer(ngram_range=(1, 2))
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(x_train['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(x_test['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(x_cv['Variation'])

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

Text Feature (Using CountVectorizer-->unigrams,bigrams)

In [126]:

```
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = CountVectorizer(min_df=3,ngram_range=(1, 2))
train_text_feature_onehotCoding = text_vectorizer.fit_transform(x_train['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 : 749979

In [127]:

```
#response coding of text features
train_text_feature_responseCoding = get_text_responsecoding(x_train)
test_text_feature_responseCoding = get_text_responsecoding(x_test)
cv_text_feature_responseCoding = get_text_responsecoding(x_cv)

# 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
```

It is always a good programming practice to Normalize after you OneHotEncode or ResponseCode

In [128]:

```

# 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(x_test['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(x_cv['TEXT'])
# don't forget to normalize every feature
cv_text_feature_onehotCoding = normalize(cv_text_feature_onehotCoding, axis=0)

```

Stack above three features

In [129]:

```

# 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(y_train['Class']))

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

cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(y_cv['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 [130]:

```

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

```

```

One hot encoding features :
(number of data points * number of features) in train data = (2124, 752271)
(number of data points * number of features) in test data = (665, 752271)
(number of data points * number of features) in cross validation data = (532, 752271)

```

In [131]:

```

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

```

```

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

```

Logistic Regression with Class Balancing

In [132]:

```

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 : 1.5647275723106386
for alpha = 1e-05
Log Loss : 1.5528832912793056
for alpha = 0.0001
Log Loss : 1.561910204182532
for alpha = 0.001
Log Loss : 1.4953342898900472

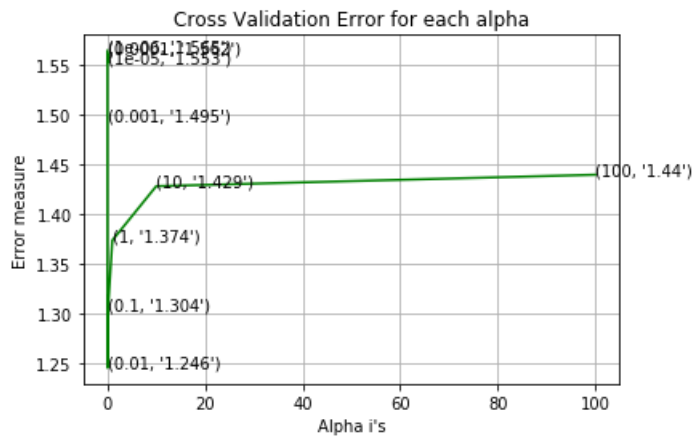
```



```

for alpha = 0.01
Log Loss : 1.246029518056866
for alpha = 0.1
Log Loss : 1.3038345911082367
for alpha = 1
Log Loss : 1.374137220824016
for alpha = 10
Log Loss : 1.4285912301357453
for alpha = 100
Log Loss : 1.4400977902868737

```



For values of best alpha = 0.01 The train log loss is: 0.8591839927993891
 For values of best alpha = 0.01 The cross validation log loss is: 1.246029518056866
 For values of best alpha = 0.01 The test log loss is: 1.2053930532171568

In [133]:

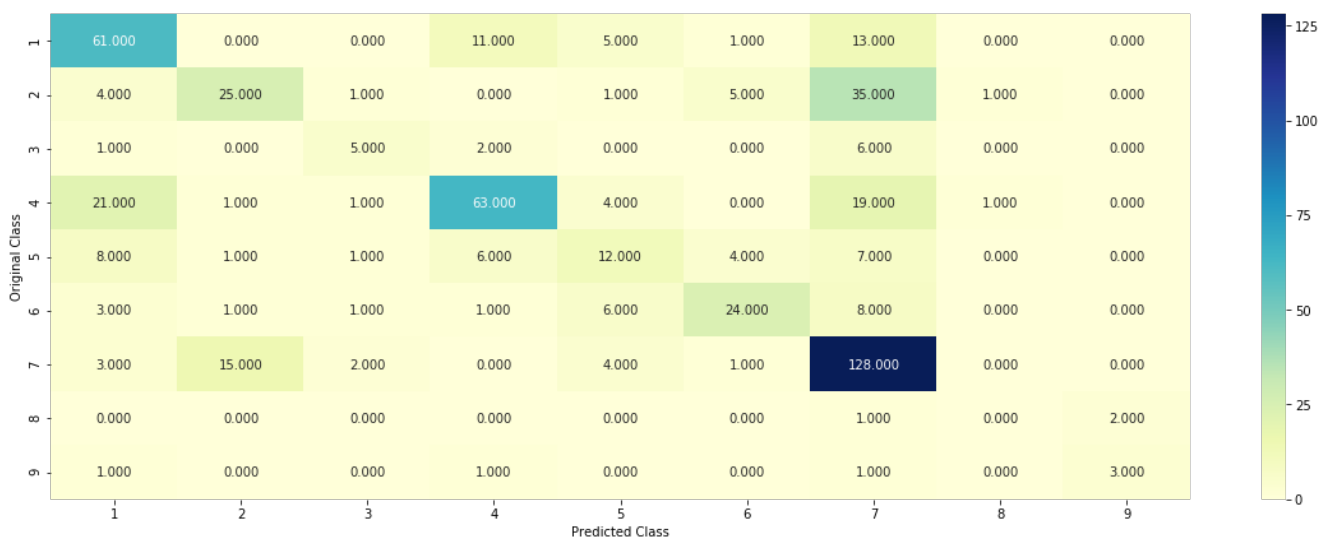
```

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)

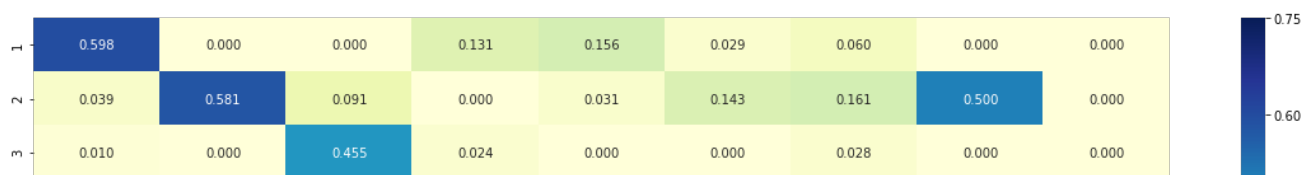
```

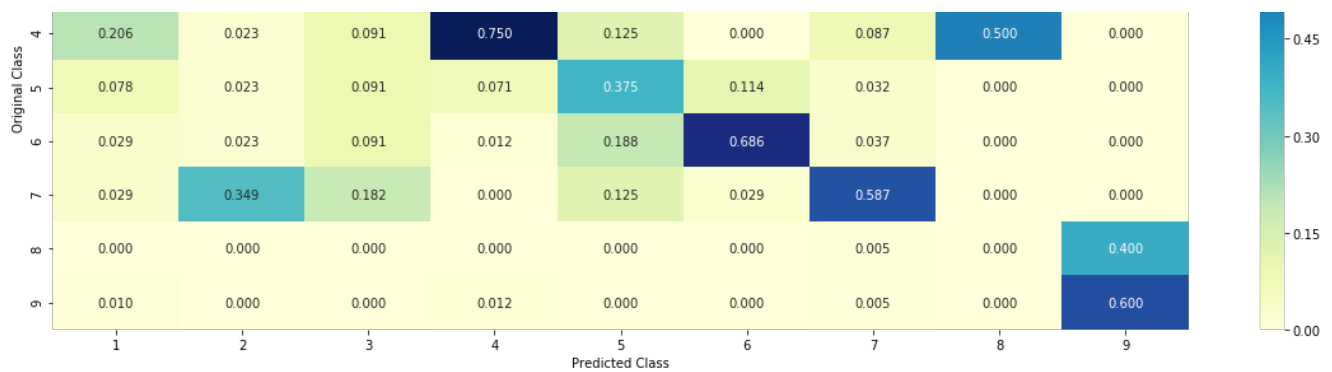
Log loss : 1.246029518056866
 Number of mis-classified points : 0.3966165413533835

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

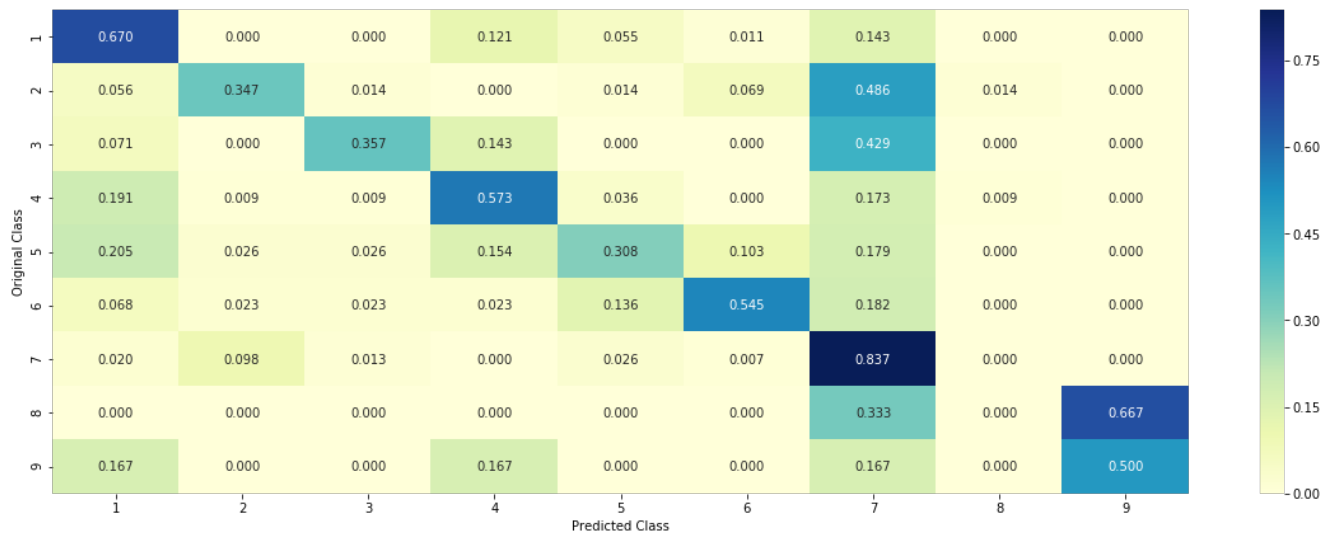


----- Precision matrix (Column Sum=1) -----





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



The application of Unigrams and Bigrams didn't help much in minimising the log-loss

Task-3:

Try any of the feature engineering techniques discussed in the course to reduce the CV and test log-loss to a value less than 1.0

Gene Feature

In [134]:

```
result = pd.merge(data_variants, data_text, on='ID', how='left')
result.loc[result['TEXT'].isnull(), 'TEXT'] = result['Gene'] + ' ' + result['Variation']
y_true = result['Class'].values
result.Gene = result.Gene.str.replace('\s+', '_')
result.Variation = result.Variation.str.replace('\s+', '_')

x_train, x_test, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_size=0.2)
x_train, x_cv, y_train, y_cv = train_test_split(x_train, y_train, stratify=y_train, test_size=0.2)
```

In [135]:

```
# get_gv_fea_dict: Get Gene variation Feature Dict
def get_gv_fea_dict(alpha, feature, df):
    value_count = x_train[feature].value_counts()
    gv_dict = dict()
    for i, denominator in value_count.items():
        vec = []
        for k in range(1, 10):
            cls_cnt = x_train.loc[(x_train['Class']==k) & (x_train[feature]==i)]
            vec.append((cls_cnt.shape[0] + alpha*10) / (denominator + 90*alpha))
        gv_dict[i] = vec
```

```

        return gv_dict

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    gv_dict = get_gv_fea_dict(alpha, feature, df)
    value_count = x_train[feature].value_counts()
    gv_fea = []
    for index, row in df.iterrows():
        if row[feature] in dict(value_count).keys():
            gv_fea.append(gv_dict[row[feature]])
        else:
            gv_fea.append([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
    return gv_fea

```

In [136]:

```

#response-coding of the Gene feature
# alpha is used for laplace smoothing
alpha = 1

# train gene feature
train_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", x_train))

# test gene feature
test_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", x_test))

# cross validation gene feature
cv_gene_feature_responseCoding = np.array(get_gv_feature(alpha, "Gene", x_cv))

```

In [137]:

```

# one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer()
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(x_train['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(x_test['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(x_cv['Gene'])

```

Variation Feature

In [138]:

```

# alpha is used for laplace smoothing
alpha = 1

# train gene feature
train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_train))

# test gene feature
test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_test))

# cross validation gene feature
cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", x_cv))

```

In [139]:

```

# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer()
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(x_train['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(x_test['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(x_cv['Variation'])

```

Text Feature

In [140]:

```

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

```

```

        return dictionary

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

```

In [141]:

```

# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer()
train_text_feature_onehotCoding = text_vectorizer.fit_transform(x_train['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 : 125643

In [142]:

```

dict_list = []
# dict_list=[] contains 9 dictionaries each corresponds to a class
for i in range(1,10):
    cls_text = x_train[x_train['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(x_train)

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 [143]:

```

#response coding of text features
train_text_feature_responseCoding = get_text_responsecoding(x_train)
test_text_feature_responseCoding = get_text_responsecoding(x_test)
cv_text_feature_responseCoding = get_text_responsecoding(x_cv)

# 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 Team(axis=1, / / cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.sum(axis=1)).T
```

In [144]:

```
test_text_feature_onehotCoding = text_vectorizer.transform(x_test['TEXT'])
cv_text_feature_onehotCoding = text_vectorizer.transform(x_cv['TEXT'])
```

Feature Engineering -1

Merging Gene+Variation into 1 single List . Lets see what happens

In [145]:

```
# Collecting all the genes and variations data into a single list
gene_variation = []

for gene in data_variants['Gene'].values:
    gene_variation.append(gene)

for variation in data_variants['Variation'].values:
    gene_variation.append(variation)
```

In [154]:

```
len(gene_variation)
```

Out[154]:

6642

In [146]:

```
tfidfVectorizer = TfidfVectorizer(max_features=1000)
text2 = tfidfVectorizer.fit_transform(gene_variation)
gene_variation_features = tfidfVectorizer.get_feature_names()

train_text = tfidfVectorizer.transform(x_train['TEXT'])
test_text = tfidfVectorizer.transform(x_test['TEXT'])
cv_text = tfidfVectorizer.transform(x_cv['TEXT'])
```

Stack above three features

In [147]:

```
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)
)

# Adding the train_text feature
train_x_onehotCoding = hstack((train_gene_var_onehotCoding, train_text))
train_x_onehotCoding = hstack((train_x_onehotCoding, train_text_feature_onehotCoding)).tocsr()
train_y = np.array(list(x_train['Class']))

# Adding the test_text feature
test_x_onehotCoding = hstack((test_gene_var_onehotCoding, test_text))
test_x_onehotCoding = hstack((test_x_onehotCoding, test_text_feature_onehotCoding)).tocsr()
test_y = np.array(list(x_test['Class']))

# Adding the cv_text feature
cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text))
cv_x_onehotCoding = hstack((cv_x_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(x_cv['Class']))

train_gene_var_responseCoding =
hstack((train_gene_feature_responseCoding, train_variation_feature_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 [148]:

```

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

```

One hot encoding features :

(number of data points * number of features) in train data = (2124, 128842)

(number of data points * number of features) in test data = (665, 128842)

(number of data points * number of features) in cross validation data = (532, 128842)

In [149]:

```

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)

Logistic Regression with Class Balancing

In [150]:

```

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', ran
dom_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))

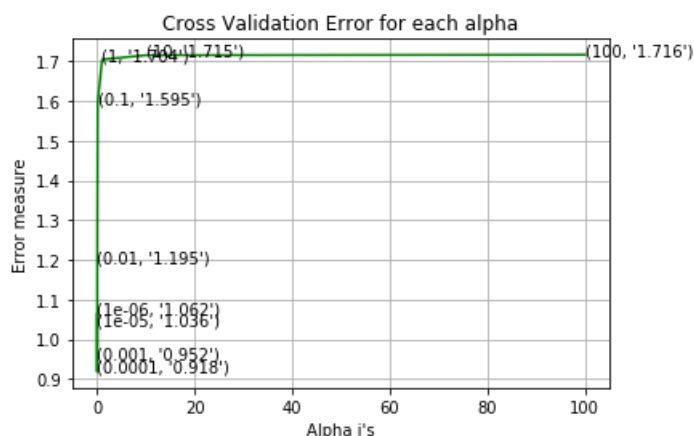
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha], "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 1e-06
Log Loss : 1.062493016875227
for alpha = 1e-05
Log Loss : 1.0361959585884248
for alpha = 0.0001
Log Loss : 0.9180397838001831
for alpha = 0.001
Log Loss : 0.9519728023692211
for alpha = 0.01
Log Loss : 1.1950116711977852
for alpha = 0.1
Log Loss : 1.5945993765858266
for alpha = 1
Log Loss : 1.7044001469110754
for alpha = 10
Log Loss : 1.715307477076851
for alpha = 100
Log Loss : 1.7164540344577839

```



```

For values of best alpha = 0.0001 The train log loss is: 0.45376574789211205
For values of best alpha = 0.0001 The cross validation log loss is: 0.9180397838001831
For values of best alpha = 0.0001 The test log loss is: 1.0137284061198686

```

In [151]:

```

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

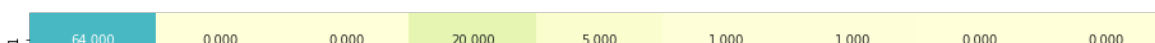
```

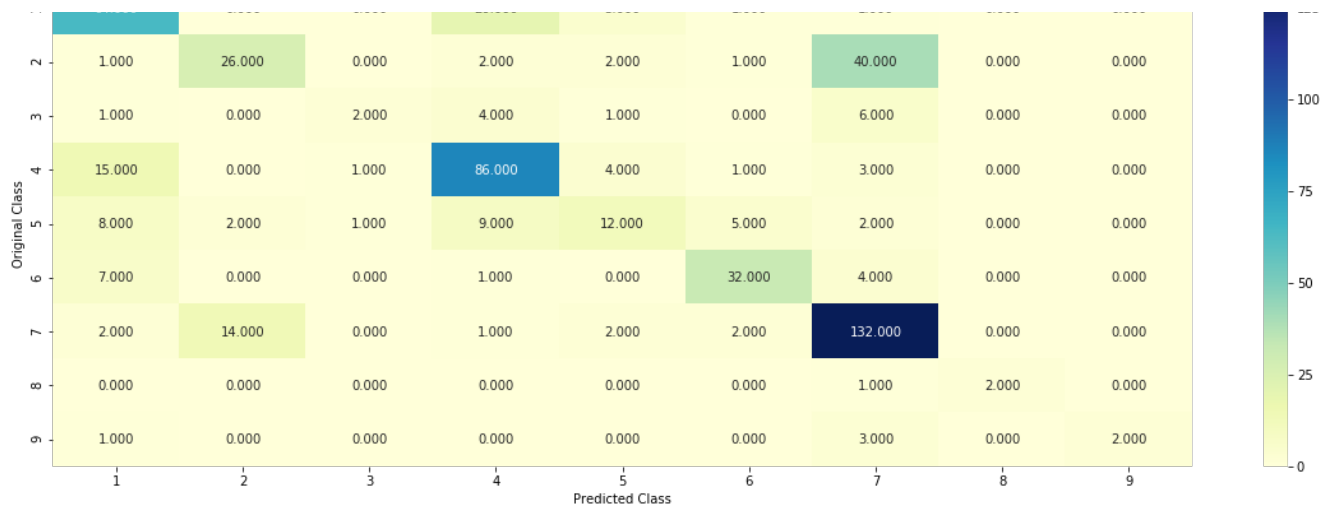
```

Log loss : 0.9180397838001831
Number of mis-classified points : 0.32706766917293234

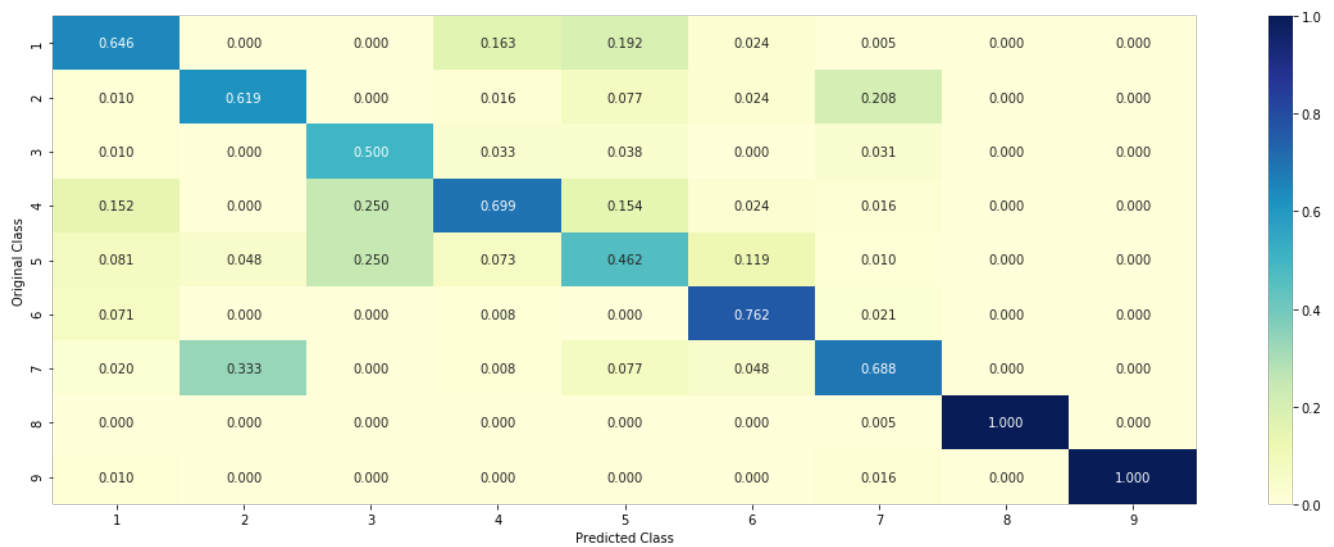
```

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

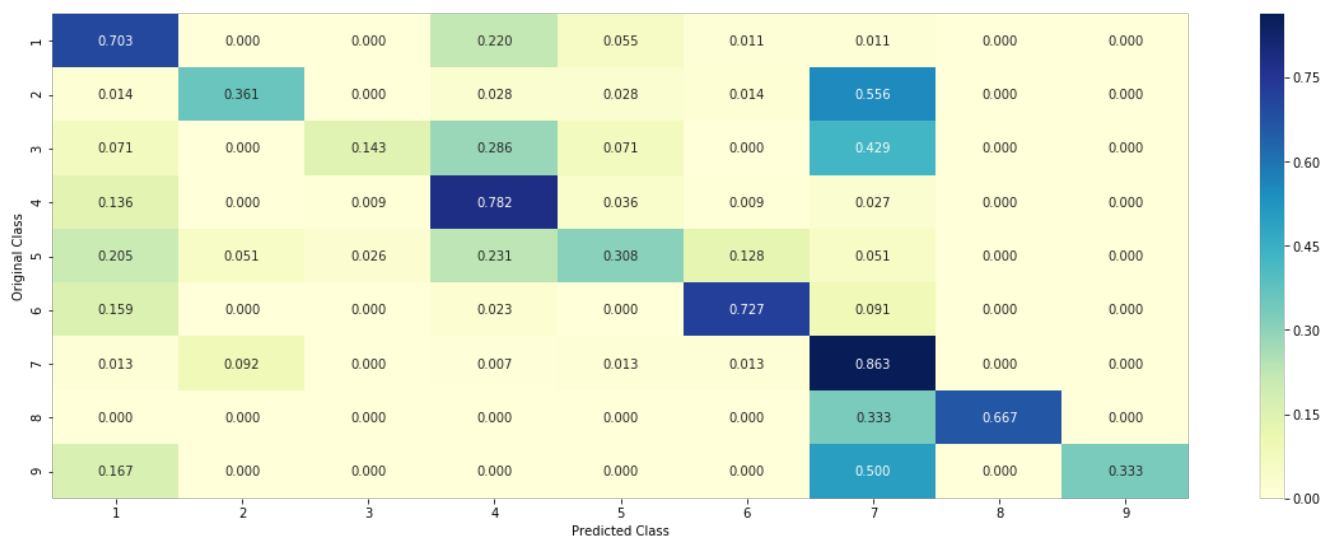




Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Applying this Feature Engineering it gave us better results , since we were able to minimise the log-loss<1.

Feature Engineering-2

Lets join the list of Genes+Variation and see what happens

Let's join the list of Gene + Variation and see what happens

In []:

```
# Collecting all the genes and variations data into a single list
gene_variation = []

for gene in data_variants['Gene'].values:
    gene_variation.append(gene)

for variation in data_variants['Variation'].values:
    gene_variation.append(variation)
```

In [155]:

```
gene_variation1=gene_variation+gene_variation
```

In [156]:

```
len(gene_variation1)
```

Out[156]:

13284

In [157]:

```
tfidfVectorizer = TfidfVectorizer(max_features=1000)
text2 = tfidfVectorizer.fit_transform(gene_variation1)
gene_variation_features = tfidfVectorizer.get_feature_names()

train_text = tfidfVectorizer.transform(x_train['TEXT'])
test_text = tfidfVectorizer.transform(x_test['TEXT'])
cv_text = tfidfVectorizer.transform(x_cv['TEXT'])
```

Stack above three features

In [159]:

```
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)
)

# Adding the train_text feature
train_x_onehotCoding = hstack((train_gene_var_onehotCoding, train_text))
train_x_onehotCoding = hstack((train_x_onehotCoding, train_text_feature_onehotCoding)).tocsr()
train_y = np.array(list(x_train['Class']))

# Adding the test_text feature
test_x_onehotCoding = hstack((test_gene_var_onehotCoding, test_text))
test_x_onehotCoding = hstack((test_x_onehotCoding, test_text_feature_onehotCoding)).tocsr()
test_y = np.array(list(x_test['Class']))

# Adding the cv_text feature
cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text))
cv_x_onehotCoding = hstack((cv_x_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(x_cv['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 [160]:

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

```
One hot encoding features :
(number of data points * number of features) in train data = (2124, 128842)
(number of data points * number of features) in test data = (665, 128842)
(number of data points * number of features) in cross validation data = (532, 128842)
```

In [161]:

```
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)
```

Logistic Regression with Class Balancing

In [162]:

```
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', ran
dom_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))

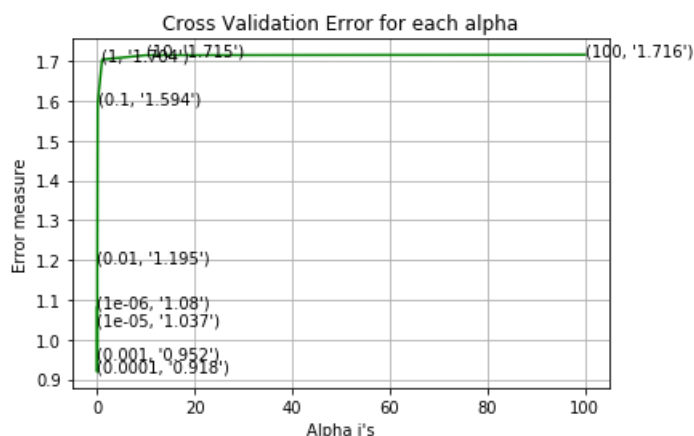
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha], "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 1e-06
Log Loss : 1.080455258961893
for alpha = 1e-05
Log Loss : 1.0373738312276255
for alpha = 0.0001
Log Loss : 0.9184997067865897
for alpha = 0.001
Log Loss : 0.9520260974391295
for alpha = 0.01
Log Loss : 1.194640174281467
for alpha = 0.1
Log Loss : 1.5941327331438089
for alpha = 1
Log Loss : 1.7039768825961057
for alpha = 10
Log Loss : 1.7148748150394522
for alpha = 100
Log Loss : 1.716020189439239

```



```

For values of best alpha = 0.0001 The train log loss is: 0.4539616012897407
For values of best alpha = 0.0001 The cross validation log loss is: 0.9184997067865897
For values of best alpha = 0.0001 The test log loss is: 1.0142290465465913

```

In [163]:

```

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

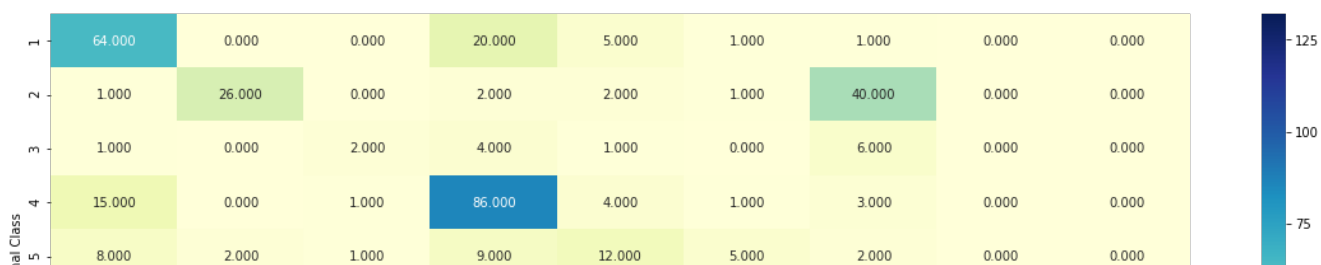
```

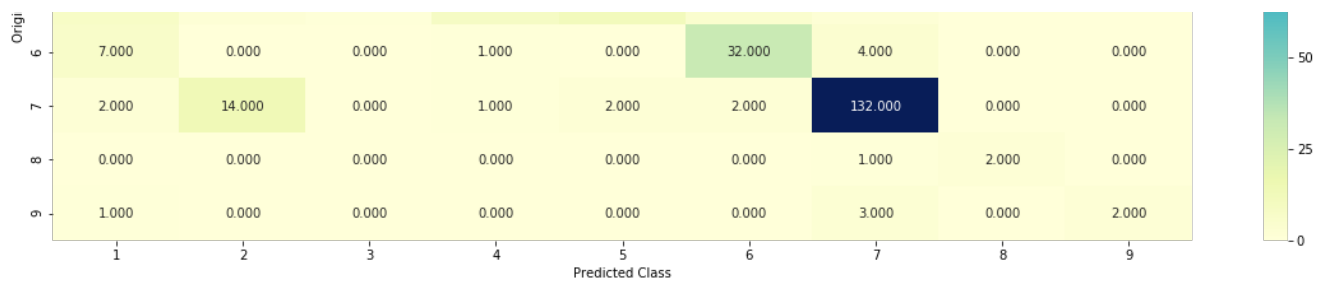
```

Log loss : 0.9184997067865897
Number of mis-classified points : 0.32706766917293234

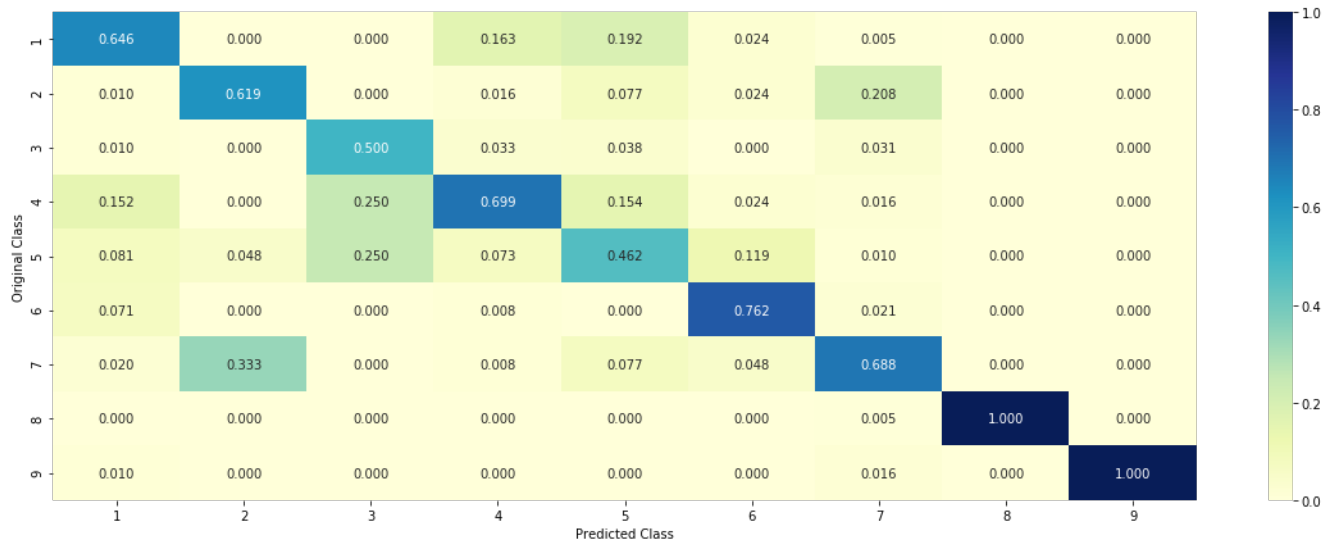
```

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

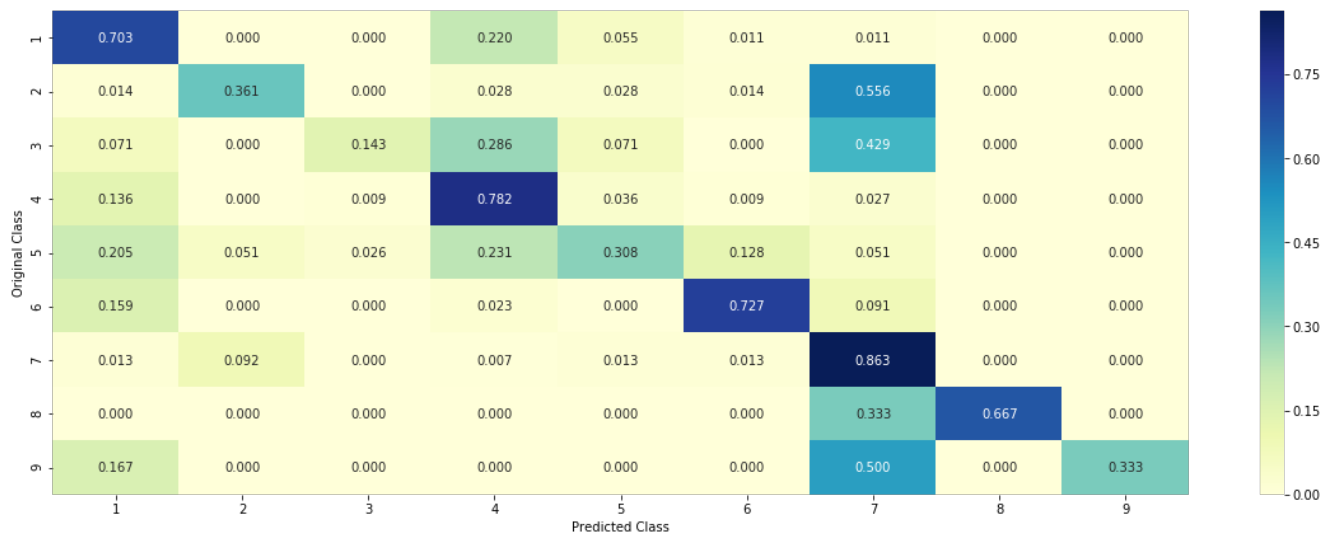




Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Using this Feature Engineering we got $\log_loss < 1$. And, Feature Engineering 2 was better than Feature Engineering 1 in case of Log-loss minimisation

Feature Engineering -3

Since, Gene Feature was found Stable than Variation feature, so, lets take Gene feature multiple times

In [192]:

```
# Collecting all the genes and variations data into a single list
gene_variation = []

for gene in data_variants['Gene'].values:
```

```
gene_variation.append(gene)
```

In [193]:

```
for gene in data_variants['Gene'].values:  
    gene_variation.append(gene)
```

In [194]:

```
for gene in data_variants['Gene'].values:  
    gene_variation.append(gene)
```

In [195]:

```
for gene in data_variants['Gene'].values:  
    gene_variation.append(gene)
```

In [196]:

```
for gene in data_variants['Gene'].values:  
    gene_variation.append(gene)
```

In [197]:

```
for gene in data_variants['Gene'].values:  
    gene_variation.append(gene)
```

In [198]:

```
for gene in data_variants['Gene'].values:  
    gene_variation.append(gene)
```

In [199]:

```
len(gene_variation)
```

Out[199]:

23247

In [200]:

```
tfidfVectorizer = TfidfVectorizer(max_features=1000)  
text2 = tfidfVectorizer.fit_transform(gene_variation)  
gene_variation_features = tfidfVectorizer.get_feature_names()  
  
train_text = tfidfVectorizer.transform(x_train['TEXT'])  
test_text = tfidfVectorizer.transform(x_test['TEXT'])  
cv_text = tfidfVectorizer.transform(x_cv['TEXT'])
```

Stack above three features

In [201]:

```
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)  
)  
  
# Adding the train_text feature  
train_x_onehotCoding = hstack((train_gene_var_onehotCoding, train_text))  
train_x_onehotCoding = hstack((train_x_onehotCoding, train_text_feature_onehotCoding)).tocsr()  
train_y = np.array(list(x_train['Class']))  
  
# Adding the test text feature
```

```
# Adding the test_text feature
test_x_onehotCoding = hstack((test_gene_var_onehotCoding, test_text))
test_x_onehotCoding = hstack((test_x_onehotCoding, test_text_feature_onehotCoding)).tocsr()
test_y = np.array(list(x_test['Class']))

# Adding the cv_text feature
cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text))
cv_x_onehotCoding = hstack((cv_x_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(x_cv['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 [202]:

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

```
One hot encoding features :
(number of data points * number of features) in train data = (2124, 128105)
(number of data points * number of features) in test data = (665, 128105)
(number of data points * number of features) in cross validation data = (532, 128105)
```

In [203]:

```
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)
```

Logistic Regression with Class Balancing

In [204]:

```
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))

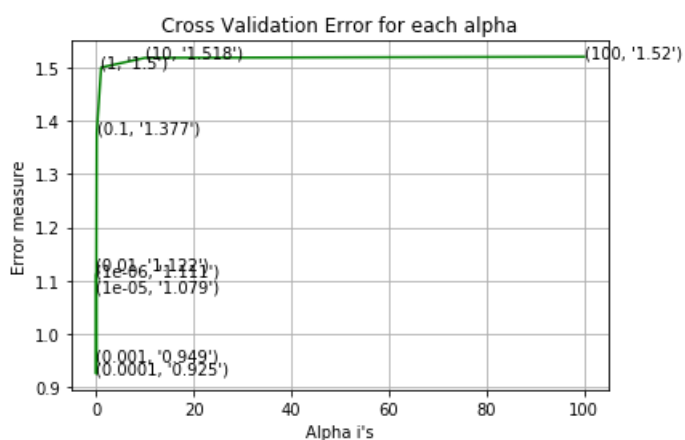
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha], "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 1e-06
Log Loss : 1.1107768026197307
for alpha = 1e-05
Log Loss : 1.078795318635448
for alpha = 0.0001
Log Loss : 0.9250584300681107
for alpha = 0.001
Log Loss : 0.9487667019335568
for alpha = 0.01
Log Loss : 1.1215608913386945
for alpha = 0.1
Log Loss : 1.377171033176285
for alpha = 1
Log Loss : 1.499588802967235
for alpha = 10
Log Loss : 1.5176087910831313
for alpha = 100
Log Loss : 1.5197119443381557

```



```

For values of best alpha = 0.0001 The train log loss is: 0.4548204167853577
For values of best alpha = 0.0001 The cross validation log loss is: 0.9250584300681107
For values of best alpha = 0.0001 The test log loss is: 0.9990770077743751

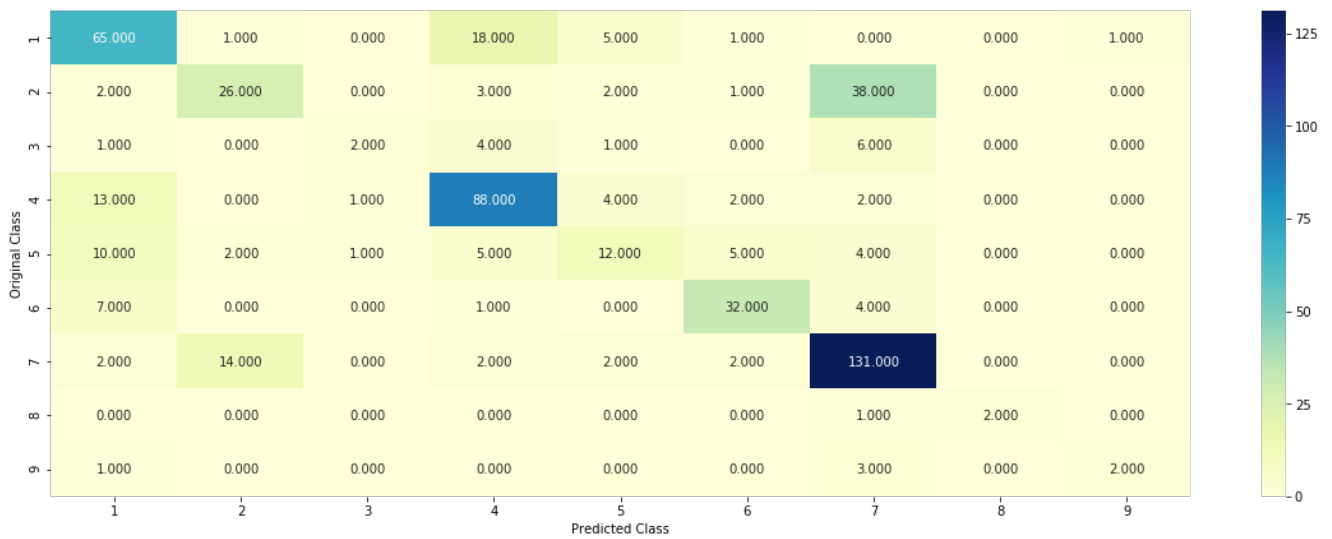
```

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

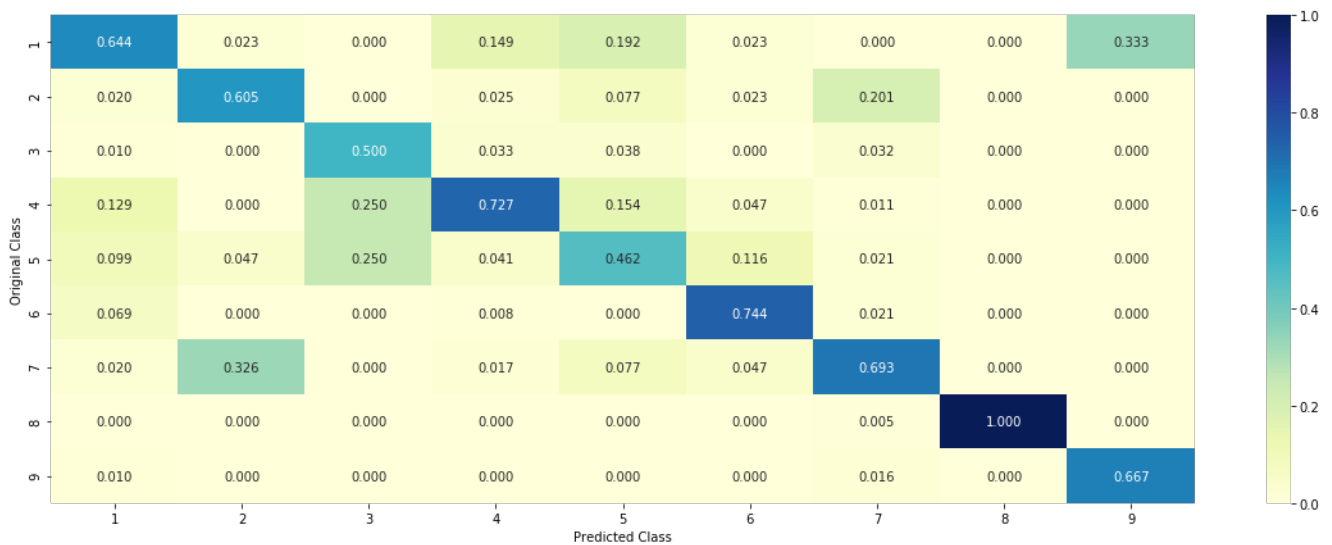
Log loss : 0.9250584300681107

Number of mis-classified points : 0.3233082706766917

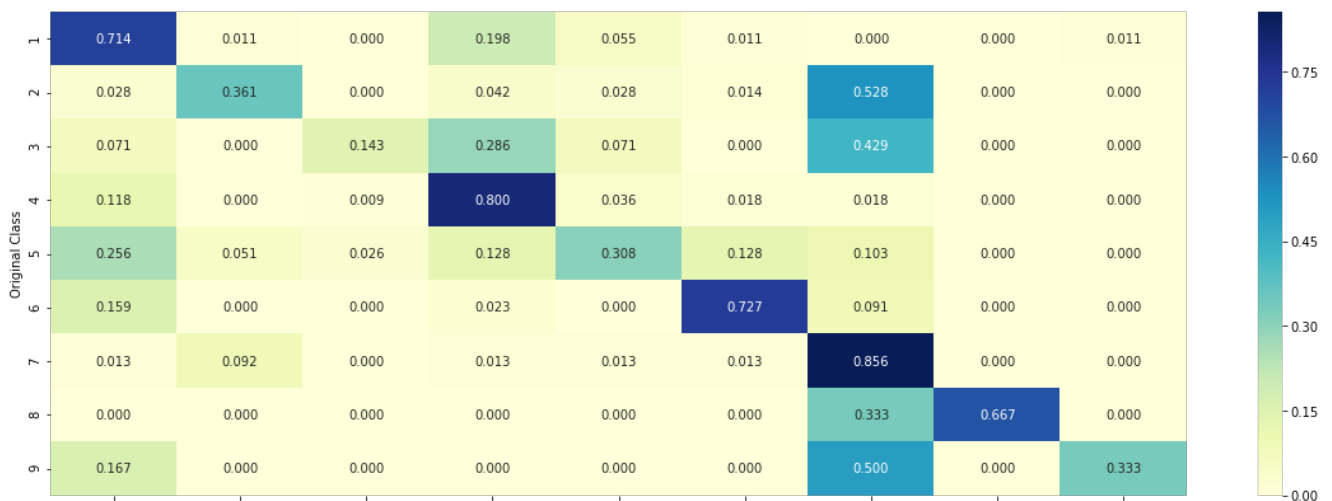
----- Confusion matrix -----



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



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



This was bad than Feature Engineering-2

Summarising using a PRETTY TABLE

In [208]:

```
from prettytable import PrettyTable
ptable = PrettyTable()
ptable.title = "*** Model Summary *** [Performance Metric: Log-Loss]"
ptable.field_names=["Model Name","Train","CV","Test","% Misclassified Points"]
ptable.add_row(["Naive Bayes","0.92","1.24","1.17","40"])
ptable.add_row(["KNN","0.64","1.07","1.01","38"])
ptable.add_row(["Logistic Regression With Class balancing","0.58","1.14","1.07","36"])
ptable.add_row(["Logistic Regression Without Class balancing","0.57","1.15","1.09","36"])
ptable.add_row(["Linear SVM","0.71","1.18","1.11","36"])
ptable.add_row(["Random Forest Classifier With One hot Encoding","0.65","1.17","1.13","41"])
ptable.add_row(["Random Forest Classifier With Response Coding","0.05","1.37","1.31","48"])
ptable.add_row(["Stack Models:LR+NB+SVM","0.63","1.11","1.07","33"])
ptable.add_row(["Maximum Voting classifier","0.70","1.09","1.03","34"])
ptable.add_row(["LR with Class Balancing (Unigrams and Bigrams)","0.86","1.24","1.21","39"])
ptable.add_row(["LR with Class Balancing (Feature Engineering-1)","0.45","0.92","1.01","33"])
ptable.add_row(["LR with Class Balancing (Feature Engineering-2)","0.45","0.91","1.01","33"])
ptable.add_row(["LR with Class Balancing (Feature Engineering-3)","0.45","0.91","0.99","33"])
print(ptable)
print()
```

Model Name	Train	CV	Test	% Misclassified Points
Naive Bayes	0.92	1.24	1.17	40
KNN	0.64	1.07	1.01	38
Logistic Regression With Class balancing	0.58	1.14	1.07	36
Logistic Regression Without Class balancing	0.57	1.15	1.09	36
Linear SVM	0.71	1.18	1.11	36
Random Forest Classifier With One hot Encoding	0.65	1.17	1.13	41
Random Forest Classifier With Response Coding	0.05	1.37	1.31	48
Stack Models:LR+NB+SVM	0.63	1.11	1.07	33
Maximum Voting classifier	0.70	1.09	1.03	34
LR with Class Balancing (Unigrams and Bigrams)	0.86	1.24	1.21	39
LR with Class Balancing (Feature Engineering-1)	0.45	0.92	1.01	33
LR with Class Balancing (Feature Engineering-2)	0.45	0.91	1.01	33
LR with Class Balancing (Feature Engineering-3)	0.45	0.91	0.99	33