

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 [14]:

```
import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import numpy as np
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import SGDClassifier
from imblearn.over_sampling import SMOTE
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC

# from sklearn.cross_validation import StratifiedKFold
from sklearn.model_selection import StratifiedKFold

from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score

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

from mlxtend.classifier import StackingClassifier

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

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

In [2]:

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

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

Out[2]:

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

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

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

3.1.2. Reading Text Data

In [3]:

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

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

Out[3]:

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

3.1.3. Preprocessing of text

In [4]:

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

def nlp_preprocessing(total_text, index, column):
    if type(total_text) is not int:
        string = ""
        # replace every special char with space
        total_text = re.sub('[^a-zA-Z0-9\\n]', ' ', total_text)
        # replace multiple spaces with single space
        total_text = re.sub('\\s+', ' ', total_text)
        # converting all the chars into lower-case.
        total_text = total_text.lower()

        for word in total_text.split():
            if word not in stop_words:
                string = string + word + " "
```

```
# 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 : 236.8552379 seconds
```

In [6]:

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

Out[6]:

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

In [7]:

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

Out[7]:

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

In [8]:

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

In [9]:

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

Out[9]:

```
out[0].
```

	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]:
```

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

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

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

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

```
In [11]:
```

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

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

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

```
In [12]:
```

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

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

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

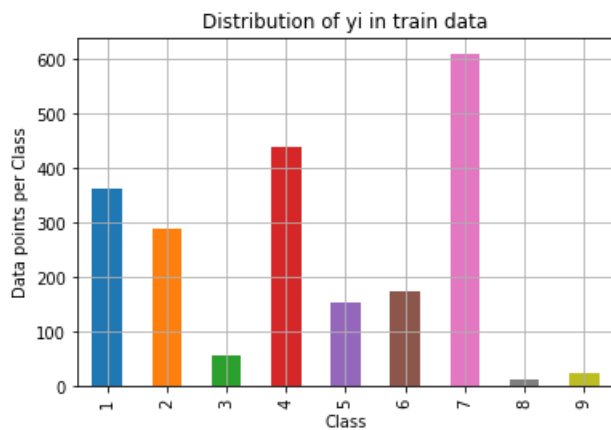
print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of y_i in test data')
plt.grid()
```

```
plt.show()

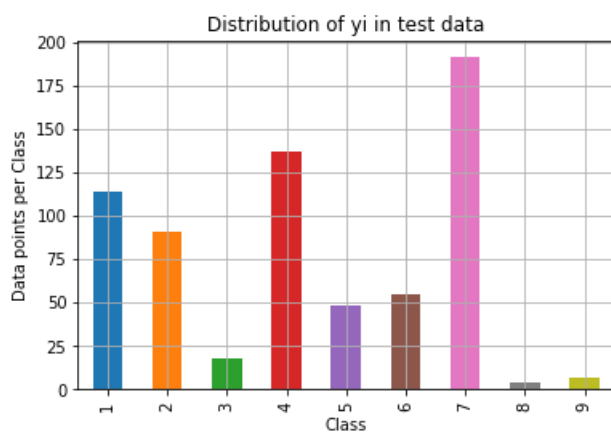
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

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

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

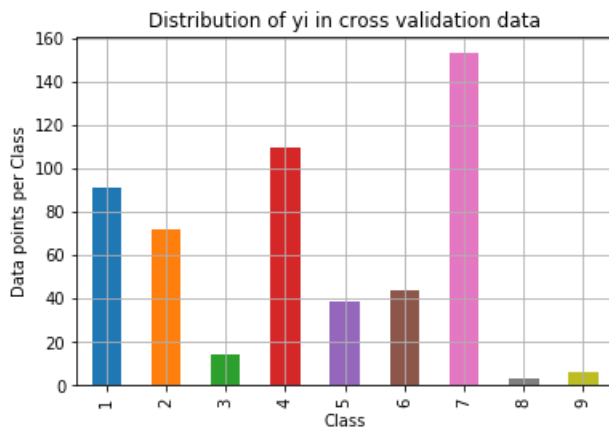


Number of data points in class 7 : 609 (28.672 %)
 Number of data points in class 4 : 439 (20.669 %)
 Number of data points in class 1 : 363 (17.09 %)
 Number of data points in class 2 : 289 (13.606 %)
 Number of data points in class 6 : 176 (8.286 %)
 Number of data points in class 5 : 155 (7.298 %)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



Number of data points in class 7 : 191 (28.722 %)
 Number of data points in class 4 : 137 (20.602 %)
 Number of data points in class 1 : 114 (17.143 %)
 Number of data points in class 2 : 91 (13.684 %)

Number of data points in class 6 : 55 (8.271 %)
 Number of data points in class 5 : 48 (7.218 %)
 Number of data points in class 3 : 18 (2.707 %)
 Number of data points in class 9 : 7 (1.053 %)
 Number of data points in class 8 : 4 (0.602 %)



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

3.2 Prediction using a 'Random' Model

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

In [13]:

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

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

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

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

    labels = [1,2,3,4,5,6,7,8,9]
```



```

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

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

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

```

In [14]:

```

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

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

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

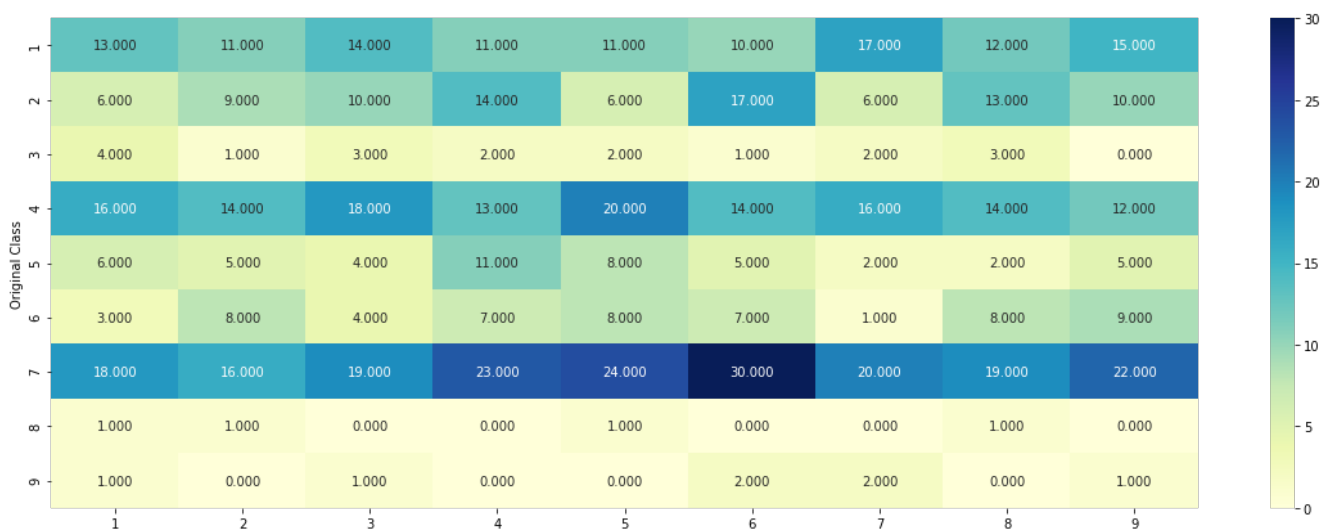
predicted_y = np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

```

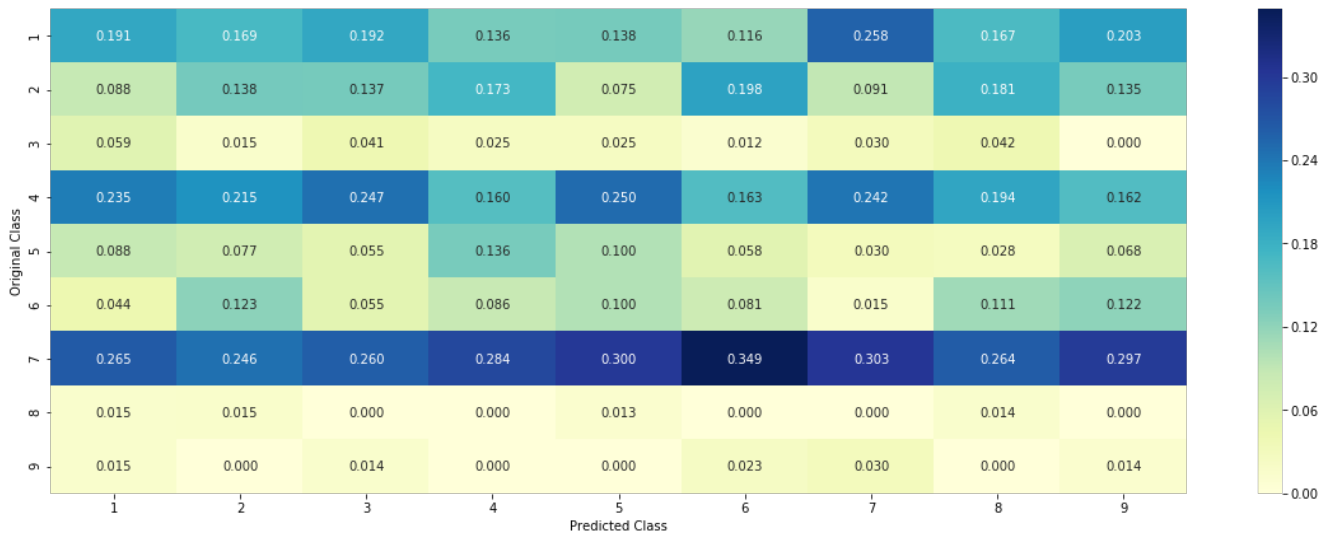
Log loss on Cross Validation Data using Random Model 2.520718570871371

Log loss on Test Data using Random Model 2.4951496087584557

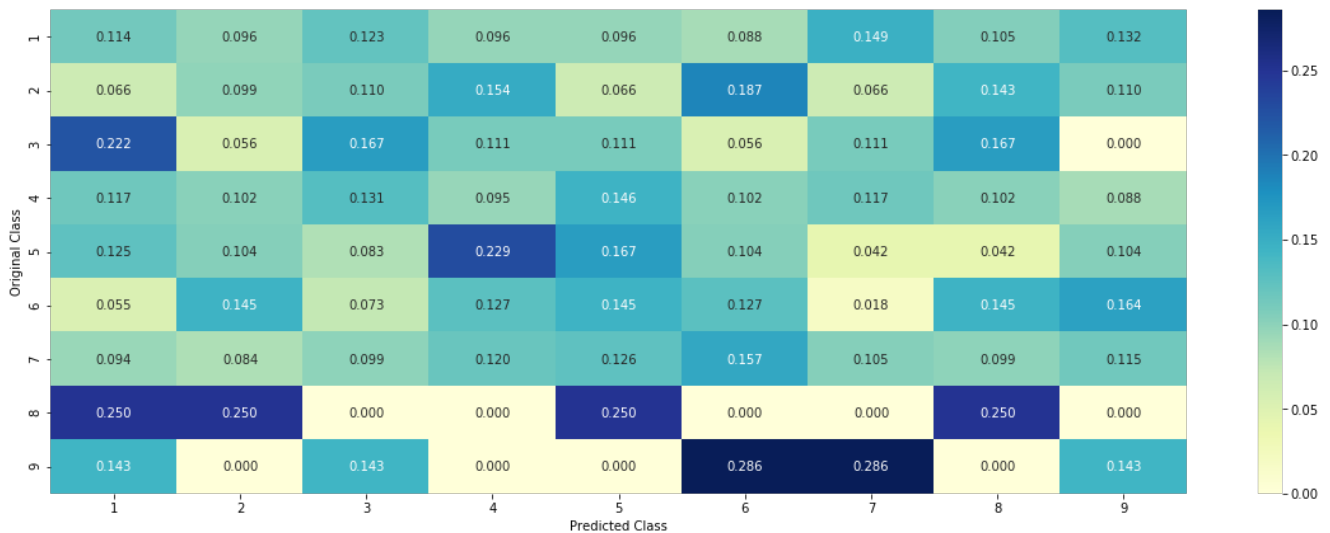
----- Confusion matrix -----



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



3.3 Univariate Analysis

In [15]:

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

# get_gv_fea_dict: Get Gene variation Feature Dict
def get_gv_fea_dict(alpha, feature, df):
    # value counts: it contains a dict like
```

```

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

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

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

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

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    #          {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.0681818181818177,
0.136363636363635, 0.25, 0.19318181818181818, 0.03787878787878788, 0.03787878787878788,
0.03787878787878788],
#           'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366,
0.27040816326530615, 0.061224489795918366, 0.066326530612244902, 0.051020408163265307, 0.051020408
163265307, 0.056122448979591837],
#           'EGFR': [0.0568181818181816, 0.21590909090909091, 0.0625, 0.0681818181818177,
0.06818181818177, 0.0625, 0.34659090909090912, 0.0625, 0.0568181818181816],
#           'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.0606060606060608,
0.0787878787878782, 0.1393939393939394, 0.34545454545454546, 0.0606060606060608,
0.0606060606060608, 0.0606060606060608],
#           'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917,
0.46540880503144655, 0.075471698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081
761006289, 0.062893081761006289],
#           'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295,
0.072847682119205295, 0.066225165562913912, 0.066225165562913912, 0.27152217890704702,
0.072847682119205295, 0.066225165562913912]

```

```

0.072047662119203293, 0.066225165562913912, 0.066225165562913912, 0.27132317680794702,
0.066225165562913912, 0.066225165562913912],
# 'BRAF': [0.06666666666666666, 0.17999999999999999, 0.07333333333333334,
0.07333333333333334, 0.09333333333333338, 0.08000000000000002, 0.29999999999999999,
0.06666666666666666, 0.06666666666666666],
# ...
# }
gv_dict = get_gv_fea_dict(alpha, feature, df)
# value_count is similar in get_gv_fea_dict
value_count = train_df[feature].value_counts()

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

```

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

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

3.2.1 Univariate Analysis on Gene Feature

Q1. Gene, What type of feature it is ?

Ans. Gene is a categorical variable

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

In [16]:

```

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

```

Number of Unique Genes : 233

```

BRCA1      171
TP53       102
EGFR        91
BRCA2       81
KIT         72
PTEN        72
BRAF        61
ALK         52
ERBB2       47
PDGFRA      38
Name: Gene, dtype: int64

```

In [17]:

```

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

```

Ans: There are 233 different categories of genes in the train data, and they are distributed as follows

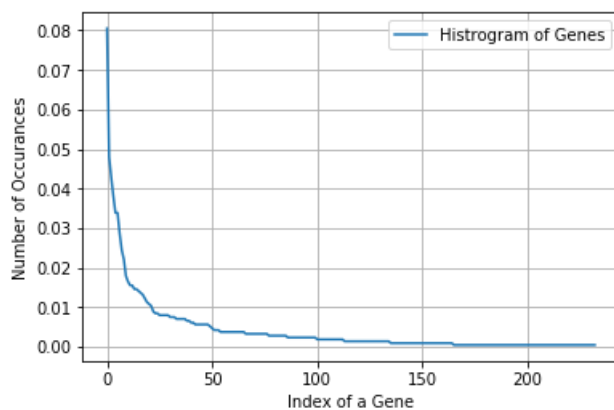
In [18]:

```

s = sum(unique_genes.values);
h = unique_genes.values/s;
plt.plot(h, label="Histogram of Genes")
plt.xlabel('Index of a Gene')

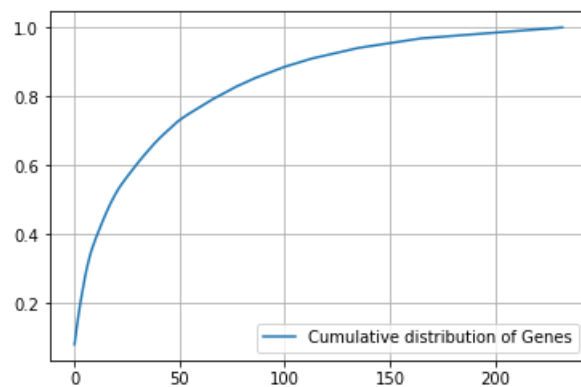
```

```
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```



In [19]:

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



Q3. How to featurize this Gene feature ?

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

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

1. One hot Encoding
2. Response coding

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

In [20]:

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

In [21]:

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

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

In [22]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer(max_features = 1000)
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

In [23]:

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

Out[23]:

```
1926      SMO
2847    BRCA2
3001      KIT
1913      SMO
1117    FANCC
Name: Gene, dtype: object
```

In [24]:

```
gene_vectorizer.get_feature_names()
```

Out[24]:

```
['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ar',
 'araf',
 'arid1b',
 'arid2',
 'arid5b',
 'atm',
 'atr',
 'atrx',
 'aurka',
 'axl',
 'b2m',
 'bap1',
 'bard1',
 'bcl10',
 'bcl2',
 'bcl2l11',
 'bcor',
 'braf',
 'brca1',
 'brca2',
 'brip1',
 'btk',
 'card11',
 'carm1',
 'casp8',
 'cbl',
 'ccnd1',
 'ccnd3',
 'ccne1',
 'cdh1',
 'cdk12',
 'cdk4',
```

'cdk6',
'cdk8',
'cdkn1a',
'cdkn1b',
'cdkn2a',
'cdkn2b',
'cdkn2c',
'cebpa',
'chek2',
'cic',
'crebbp',
'ctcf',
'ctla4',
'ctnnb1',
'ddr2',
'dicer1',
'dnmt3a',
'dnmt3b',
'egfr',
'elf3',
'ep300',
'epas1',
'epcam',
'erbb2',
'erbb3',
'erbb4',
'ercc2',
'ercc3',
'ercc4',
'erg',
'errfi1',
'esr1',
'etv1',
'etv6',
'ewsr1',
'ezh2',
'fam58a',
'fanca',
'fancc',
'fat1',
'fbxw7',
'fgf19',
'fgf3',
'fgf4',
'fgfr1',
'fgfr2',
'fgfr3',
'fgfr4',
'flt3',
'foxa1',
'foxl2',
'foxp1',
'gata3',
'gli1',
'gna11',
'gnas',
'h3f3a',
'hla',
'hnf1a',
'hras',
'idh1',
'idh2',
'igflr',
'ikbke',
'il7r',
'inpp4b',
'jak1',
'jak2',
'jun',
'kdm5a',
'kdm5c',
'kdm6a',
'kdr',
'keap1',
'kit',
'klf4',
'kmt2a',

'kmt2c',
'knstrn',
'kras',
'lats1',
'lats2',
'map2k1',
'map2k2',
'map2k4',
'map3k1',
'mapk1',
'mdm2',
'mdm4',
'med12',
'mef2b',
'met',
'mga',
'mlh1',
'msh2',
'msh6',
'mtor',
'myc',
'mycn',
'myd88',
'myod1',
'nfl',
'nf2',
'nfe2l2',
'nfkbia',
'nkx2',
'notch1',
'notch2',
'npm1',
'nras',
'nsd1',
'ntrk1',
'ntrk2',
'ntrk3',
'nup93',
'pbrm1',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pim1',
'pms1',
'pms2',
'pole',
'ppm1d',
'ppp2r1a',
'prdm1',
'ptch1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rac1',
'rad21',
'rad50',
'rad51c',
'rad51d',
'rad54l',
'raf1',
'rasal',
'rb1',
'rbm10',
'ret',
'rhoa',
'rictor',
'rit1',
'ros1',
'runx1',
'rxra',
'rybp',
'sdhb',


```
'sdhc',
'setd2',
'sf3b1',
'shoc2',
'shq1',
'smad2',
'smad3',
'smad4',
'smarca4',
'smarcb1',
'smo',
'sos1',
'sox9',
'spop',
'src',
'srsf2',
'stag2',
'stat3',
'stk11',
'tcf3',
'tcf7l2',
'tert',
'tet1',
'tet2',
'tgfbr1',
'tgfbr2',
'tmprss2',
'tp53',
'tp53bp1',
'tsc1',
'tsc2',
'u2af1',
'vhl',
'whsc1',
'whsc1l1',
'xpo1',
'xrcc2',
'yap1']
```

In [25]:

```
print("train_gene_feature_onehotCoding is converted feature using 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, 232)

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

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

In [26]:

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

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

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

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

```

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

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

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

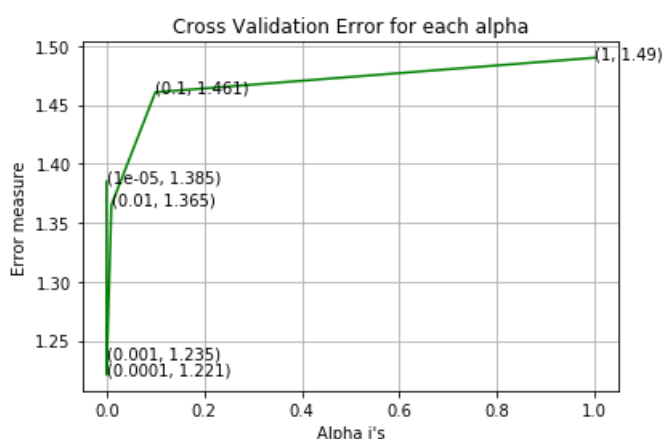
predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
print('For values of best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

```

```

For values of alpha = 1e-05 The log loss is: 1.384954953757045
For values of alpha = 0.0001 The log loss is: 1.220887545417012
For values of alpha = 0.001 The log loss is: 1.23527651920385
For values of alpha = 0.01 The log loss is: 1.364961439955826
For values of alpha = 0.1 The log loss is: 1.4606590908996107
For values of alpha = 1 The log loss is: 1.4899788637853473

```



```

For values of best alpha = 0.0001 The train log loss is: 1.0526060099705743
For values of best alpha = 0.0001 The cross validation log loss is: 1.220887545417012
For values of best alpha = 0.0001 The test log loss is: 1.1843654625498663

```

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.

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

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

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

Q6. How many data points in Test and CV datasets are covered by the 233 genes in train dataset?
 Ans
 1. In test data 643 out of 665 : 96.69172932330827
 2. In cross validation data 512 out of 532 : 96.2406015037594

3.2.2 Univariate Analysis on Variation Feature

Q7. Variation, What type of feature is it ?

Ans. Variation is a categorical variable

Q8. How many categories are there?

In [28]:

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

```
Number of Unique Variations : 1919
Truncating_Mutations      67
Deletion                   52
Amplification              46
Fusions                    26
M1R                        2
EWSR1-ETV1_Fusion         2
Q61H                      2
G12V                      2
Promoter_Hypermethylation  2
T167A                     2
Name: Variation, dtype: int64
```

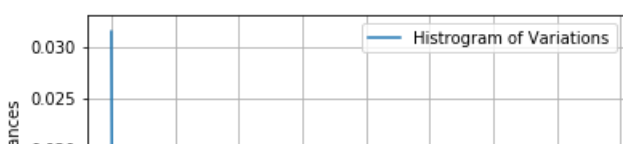
In [29]:

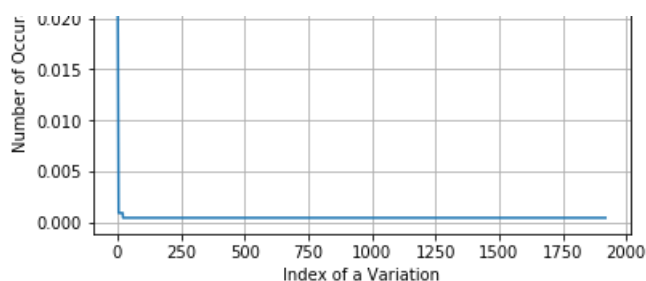
```
print("Ans: There are", unique_variations.shape[0], "different categories of variations in the train data, and they are distributed as follows",)
```

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

In [30]:

```
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()
```

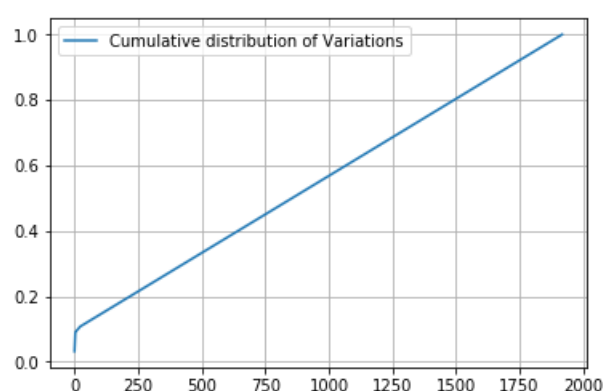




In [31]:

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

```
[0.03154426 0.05602637 0.07768362 ... 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.appliedaia.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

1. One hot Encoding
2. Response coding

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

In [32]:

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

In [33]:

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

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

In [34]:

```
# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer(max_features = 1000)
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(train_df['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(test_df['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(cv_df['Variation'])
```

In [35]:

```
print("train_variation_feature_onehotEncoded is converted feature using the 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, 1000)

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

Let's build a model just like the earlier!

In [36]:

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

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

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

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

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

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

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

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

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

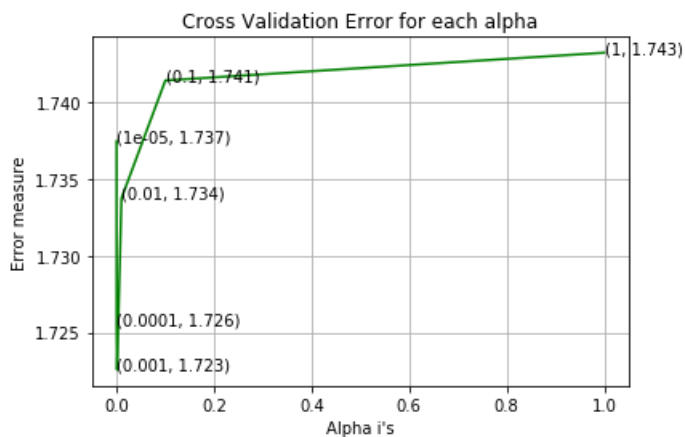
predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
```

```

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

```

For values of alpha = 1e-05 The log loss is: 1.737476973365291
 For values of alpha = 0.0001 The log loss is: 1.7255070794020986
 For values of alpha = 0.001 The log loss is: 1.7225571047902302
 For values of alpha = 0.01 The log loss is: 1.7337517466397754
 For values of alpha = 0.1 The log loss is: 1.741443301227446
 For values of alpha = 1 The log loss is: 1.7432509768377



For values of best alpha = 0.001 The train log loss is: 1.3611687786959856
 For values of best alpha = 0.001 The cross validation log loss is: 1.7225571047902302
 For values of best alpha = 0.001 The test log loss is: 1.74246364126082

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

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

In [37]:

```

print("Q12. How many data points are covered by total ", unique_variations.shape[0], " variation in test and cross validation data sets?")
test_coverage=test_df[test_df['Variation'].isin(list(set(train_df['Variation'])))]
cv_coverage=cv_df[cv_df['Variation'].isin(list(set(train_df['Variation'])))]
print('Ans\1. In test data',test_coverage, 'out of',test_df.shape[0], ":", (test_coverage/test_df.shape[0])*100)
print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[0],":", (cv_coverage/cv_df.shape[0])*100)

```

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

Ans

1. In test data 62 out of 665 : 9.323308270676693
2. In cross validation data 53 out of 532 : 9.962406015037594

3.2.3 Univariate Analysis on Text Feature

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

In [38]:

```

# cls_text is a data frame
# for every row in data fram consider the 'TEXT'

```

```

# split the words by space
# make a dict with those words
# increment its count whenever we see that word

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

```

In [39]:

```

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

```

In [40]:

```

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

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

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

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

```

Total number of unique words in train data : 1000

In [41]:

```

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

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

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

```

In [42]:

In [43]:

In [44]:

In [45]:

In [46]:

[illegible]

29.488239046426884: 1, 29.364605197463003: 1, 29.061709827602737: 1, 29.00707560911086: 1, 28.835610450201692: 1, 28.724436663271803: 1, 28.717545833205843: 1, 28.61807411194339: 1, 28.53210757213647: 1, 28.25463672412898: 1, 28.042138967240305: 1, 27.968867159172532: 1, 27.920780874664292: 1, 27.831394523584635: 1, 27.816808400978093: 1, 27.809680911038978: 1, 27.736742556325456: 1, 27.61706496562173: 1, 27.535601336155192: 1, 27.305596863029102: 1, 26.942905527706966: 1, 26.938781600860896: 1, 26.914947746360177: 1, 26.55094418897463: 1, 26.523662601040247: 1, 26.45340274229819: 1, 26.210484664259987: 1, 26.02018570993738: 1, 25.987965465774806: 1, 25.792172970433224: 1, 25.74056193129893: 1, 25.727823416301028: 1, 25.52668174875706: 1, 25.44840918012738: 1, 25.41685614209588: 1, 25.365200262669177: 1, 25.257488550447068: 1, 25.161211211708864: 1, 24.89939612475789: 1, 24.74726923847684: 1, 24.642023298959728: 1, 24.627997606263925: 1, 24.57211199305058: 1, 24.50445736903497: 1, 24.48481795047774: 1, 24.40316486384663: 1, 24.371344958716236: 1, 24.310413864561415: 1, 24.22710636822915: 1, 24.21764290998845: 1, 24.15279073879316: 1, 24.13784197029008: 1, 24.11570967445754: 1, 23.954659084968934: 1, 23.932102139549727: 1, 23.88624565384889: 1, 23.841057515907913: 1, 23.80410763079612: 1, 23.61924862062912: 1, 23.58510128061822: 1, 23.542682120415922: 1, 23.498071143900983: 1, 23.448447269956308: 1, 23.38891721808099: 1, 23.25021265850558: 1, 23.17795833426696: 1, 23.169920094010184: 1, 23.109652395025574: 1, 23.005242649001897: 1, 22.98639731868042: 1, 22.948981868486808: 1, 22.91961119831023: 1, 22.910368101738857: 1, 22.862233761897688: 1, 22.78613107709289: 1, 22.77614300515613: 1, 22.768119392545888: 1, 22.76231746619134: 1, 22.630461065267934: 1, 22.467389344689252: 1, 22.460570723804427: 1, 22.35959482666052: 1, 22.25145215016261: 1, 22.036888551322804: 1, 22.01275954842131: 1, 22.01082728407659: 1, 22.004844813026583: 1, 21.95982663567893: 1, 21.906874336970173: 1, 21.867609054392492: 1, 21.807785429640926: 1, 21.745882764373555: 1, 21.7430404200225: 1, 21.737050200748524: 1, 21.675226473723153: 1, 21.585336242899942: 1, 21.52658865153169: 1, 21.516973176020855: 1, 21.51149407045237: 1, 21.49966573743112: 1, 21.469003951153546: 1, 21.33966618174949: 1, 21.320478868722745: 1, 21.273629104156985: 1, 21.194469788962166: 1, 21.126553356724656: 1, 21.052513898001635: 1, 21.000134797236573: 1, 20.979525431771908: 1, 20.957180265519202: 1, 20.904853252726305: 1, 20.86974127547512: 1, 20.781542903897975: 1, 20.70283926361183: 1, 20.69634476725769: 1, 20.642654105129157: 1, 20.55939767765259: 1, 20.45500959713417: 1, 20.406665971782164: 1, 20.401785353758427: 1, 20.264490779406216: 1, 20.234494133479778: 1, 20.207585382010805: 1, 20.18324008374037: 1, 20.15572125915808: 1, 20.136543552085847: 1, 20.127025944532427: 1, 20.082772576784734: 1, 20.0733278760874: 1, 19.978618000524524: 1, 19.898023367778418: 1, 19.768616295835: 1, 19.75462387292715: 1, 19.671180222126477: 1, 19.6073135159309: 1, 19.5651434641379: 1, 19.560385617373036: 1, 19.550096010907716: 1, 19.51505650030504: 1, 19.505558749042898: 1, 19.457711998673258: 1, 19.434337506513835: 1, 19.30378313717418: 1, 19.25300797966021: 1, 19.225273522453: 1, 19.205898695612568: 1, 19.18846967875389: 1, 19.136986669153465: 1, 19.023324515689175: 1, 19.01169661866383: 1, 19.00757810121088: 1, 18.976711747896033: 1, 18.949765461271742: 1, 18.84363893373237: 1, 18.79288059072406: 1, 18.702317130277237: 1, 18.698372866529265: 1, 18.692516616234013: 1, 18.653204520745604: 1, 18.64801219149171: 1, 18.634416088354477: 1, 18.62900585518304: 1, 18.610992650325453: 1, 18.57419896727416: 1, 18.53284501448321: 1, 18.528709416606702: 1, 18.518165978988602: 1, 18.476833449920242: 1, 18.43207858986455: 1, 18.431602346543823: 1, 18.42427661883146: 1, 18.408665943510165: 1, 18.398448476649666: 1, 18.3683688368829: 1, 18.287972846747152: 1, 18.236276696027396: 1, 18.220794520085793: 1, 18.210929702681483: 1, 18.17770164742412: 1, 18.119870948368686: 1, 18.062693007851355: 1, 18.00395954921216: 1, 17.98936314429876: 1, 17.960440711651653: 1, 17.906914565705044: 1, 17.889148125153852: 1, 17.862364874178542: 1, 17.845445884092975: 1, 17.833663153711935: 1, 17.80812942713786: 1, 17.759780858795786: 1, 17.759205140368202: 1, 17.711498808767285: 1, 17.691838097462746: 1, 17.67992373749744: 1, 17.675251171808316: 1, 17.65824596314131: 1, 17.643431858515576: 1, 17.634567818948614: 1, 17.618732220903205: 1, 17.506810621322426: 1, 17.499641343499487: 1, 17.49473642479805: 1, 17.45449233956445: 1, 17.384529005721237: 1, 17.382300099647065: 1, 17.33592648946349: 1, 17.297996916844095: 1, 17.195836394583786: 1, 17.182551974322208: 1, 17.171193062247887: 1, 17.14513836238678: 1, 17.138139958975373: 1, 17.063983668262363: 1, 17.048398648171375: 1, 17.043909901926096: 1, 17.003967258186776: 1, 16.998298017486658: 1, 16.959174444793: 1, 16.9322883721242: 1, 16.929541034176207: 1, 16.909426273130347: 1, 16.90427751123065: 1, 16.85700516645319: 1, 16.833873314580746: 1, 16.82969274831408: 1, 16.814549010947353: 1, 16.80119541514669: 1, 16.78333586811674: 1, 16.781789617978788: 1, 16.76714270682867: 1, 16.749773723099796: 1, 16.72308031362737: 1, 16.691973450664623: 1, 16.691090288858852: 1, 16.643795437676644: 1, 16.632523658318306: 1, 16.586631699912182: 1, 16.55155902383042: 1, 16.551123779659303: 1, 16.519511622724735: 1, 16.460657511980923: 1, 16.426633774670567: 1, 16.402544208755547: 1, 16.392159935753: 1, 16.35870016721149: 1, 16.338027121338612: 1, 16.314324288880655: 1, 16.2291267106736: 1, 16.183977059713882: 1, 16.151844524536607: 1, 16.143042088812535: 1, 16.070433943197006: 1, 16.05560529237349: 1, 16.041975818182777: 1, 16.018040662233187: 1, 15.992038998984103: 1, 15.978748759280617: 1, 15.977370162060899: 1, 15.931631923690011: 1, 15.886151427259655: 1, 15.845068015871135: 1, 15.826721527987733: 1, 15.816935618259029: 1, 15.798851296612435: 1, 15.79246108676481: 1, 15.77658894851729: 1, 15.730007344660585: 1, 15.718282074154354: 1, 15.693798030068383: 1, 15.69267893335902: 1, 15.686705108632774: 1, 15.682549141359754: 1, 15.645283168168312: 1, 15.612755735026413: 1, 15.60747043798216: 1, 15.600062741411381: 1, 15.572290952661776: 1, 15.522929481995579: 1, 15.497170229475005: 1, 15.485846649437525: 1, 15.471169167154654: 1, 15.449630128675365: 1, 15.342407839211837: 1, 15.29326929106599: 1, 15.27512032349795: 1, 15.23075462697028: 1, 15.21549020910972: 1, 15.195494024390044: 1, 15.19263407709253: 1, 15.133123031789083: 1, 15.12481207357862: 1, 15.115283200411927: 1, 15.047658356987611: 1, 15.046130721235519: 1, 15.009811556461779: 1, 14.998148132205946: 1, 14.979354691535626: 1, 14.975111947459592: 1, 14.96850275375975: 1, 14.964551252629986: 1, 14.920086574950341: 1, 14.910074039051558: 1, 14.854658180072642: 1, 14.854478374233539: 1, 14.850960250528342: 1, 14.843684299522625: 1, 14.814196591707548: 1.

14.81035497535482: 1, 14.80991426707445: 1, 14.770266478060302: 1, 14.700941628960953: 1, 14.679495783255103: 1, 14.669261124715886: 1, 14.593443122980332: 1, 14.581164452675969: 1, 14.558043292457471: 1, 14.55328490083477: 1, 14.55121639628513: 1, 14.5333924368694: 1, 14.528916265552251: 1, 14.52221489988145: 1, 14.478651741357558: 1, 14.462366117930074: 1, 14.43617269106716: 1, 14.399417234709569: 1, 14.369654967814665: 1, 14.349565159799141: 1, 14.343753346220577: 1, 14.288675258388004: 1, 14.267005504771651: 1, 14.254917865759658: 1, 14.239543960150758: 1, 14.221806603834725: 1, 14.20976863371725: 1, 14.18317222026516: 1, 14.148210586578415: 1, 14.081669844561059: 1, 14.057204875378819: 1, 14.056065811644508: 1, 14.028820977183855: 1, 14.017160361129248: 1, 13.993331242376254: 1, 13.991401302411484: 1, 13.906482889785602: 1, 13.868181103605853: 1, 13.811459335567417: 1, 13.78414683596364: 1, 13.7785113871551: 1, 13.753236855463042: 1, 13.75024142423531: 1, 13.691530522199939: 1, 13.690571525536791: 1, 13.682401247227503: 1, 13.627840376240131: 1, 13.605770612141612: 1, 13.592179967942101: 1, 13.58385897136901: 1, 13.574560457791563: 1, 13.569309645607301: 1, 13.540143262595189: 1, 13.538306123977314: 1, 13.533829949661293: 1, 13.528509285886047: 1, 13.526579574843623: 1, 13.522918086966481: 1, 13.522047104850238: 1, 13.51538956039324: 1, 13.512055512943734: 1, 13.485917963558908: 1, 13.479636801911004: 1, 13.368294132591325: 1, 13.350325023613422: 1, 13.304902808337893: 1, 13.219794190246954: 1, 13.210648343311128: 1, 13.200239159217862: 1, 13.198723099077158: 1, 13.189207837134882: 1, 13.173971993082377: 1, 13.146998144653299: 1, 13.143330571862625: 1, 13.112125019702894: 1, 13.11152527846774: 1, 13.096801555799939: 1, 13.084112584043632: 1, 13.07146180564088: 1, 13.05148449501443: 1, 13.028571712451674: 1, 12.993691878746011: 1, 12.987872764045093: 1, 12.963480565262925: 1, 12.962059103335882: 1, 12.9612028177205: 1, 12.948515053697854: 1, 12.948466347760277: 1, 12.939835895343492: 1, 12.915414859571516: 1, 12.895424892937061: 1, 12.86460041276678: 1, 12.85799690049183: 1, 12.851618430763192: 1, 12.84388540296056: 1, 12.834933269278388: 1, 12.818661892896593: 1, 12.793781028334609: 1, 12.764992691943968: 1, 12.756030388884627: 1, 12.732830460453291: 1, 12.726255484859262: 1, 12.717945364535883: 1, 12.699052081643796: 1, 12.698643038451554: 1, 12.689311925883796: 1, 12.65083089287332: 1, 12.648002477642144: 1, 12.59905747631432: 1, 12.597553919394466: 1, 12.57277682486033: 1, 12.558212811641837: 1, 12.558042137409299: 1, 12.541200623146235: 1, 12.529349472714866: 1, 12.472936178230956: 1, 12.43745130714572: 1, 12.42474463714372: 1, 12.417636783647326: 1, 12.413061913837884: 1, 12.412107779919896: 1, 12.347830135601138: 1, 12.304473889890417: 1, 12.275308537558601: 1, 12.270381631616273: 1, 12.2581861649362: 1, 12.209342913815536: 1, 12.179881000253761: 1, 12.178674284311873: 1, 12.164176150899088: 1, 12.15038377129291: 1, 12.132425297722449: 1, 12.109718575072517: 1, 12.105252330277084: 1, 12.10457774469964: 1, 12.098564166092748: 1, 12.059162098969193: 1, 12.025697276214041: 1, 12.004173149917303: 1, 11.970454478898906: 1, 11.960002313852769: 1, 11.939236222752294: 1, 11.938788034195841: 1, 11.928356200124263: 1, 11.91601103631996: 1, 11.901873149955927: 1, 11.89251550516628: 1, 11.882362958381028: 1, 11.844453476778005: 1, 11.830872563177964: 1, 11.810123019190518: 1, 11.810037297228035: 1, 11.806180444529074: 1, 11.796158340978504: 1, 11.783606380308571: 1, 11.777673974856754: 1, 11.776087862168746: 1, 11.749225023891366: 1, 11.74516188748881: 1, 11.71045237772906: 1, 11.702460515463766: 1, 11.694191234756994: 1, 11.686699936165168: 1, 11.66137753263588: 1, 11.66010737927882: 1, 11.634767559130642: 1, 11.620040975957291: 1, 11.610637041408163: 1, 11.597878886126031: 1, 11.59587763005903: 1, 11.594378865079879: 1, 11.593601530425419: 1, 11.586051692108978: 1, 11.498045238810295: 1, 11.47254778526899: 1, 11.42775278312302: 1, 11.398962954769642: 1, 11.383343643969503: 1, 11.378897361440934: 1, 11.350207935403795: 1, 11.347724986494988: 1, 11.345056952802521: 1, 11.324786139623532: 1, 11.323039434824285: 1, 11.320517426898048: 1, 11.317511176616131: 1, 11.311615614322102: 1, 11.307852107523065: 1, 11.301783378171315: 1, 11.289341288447057: 1, 11.26931497698062: 1, 11.26850566941784: 1, 11.261778143672764: 1, 11.247617602001402: 1, 11.22725472895631: 1, 11.204816072030512: 1, 11.184573842729003: 1, 11.172202299947923: 1, 11.167242988054069: 1, 11.151639692479764: 1, 11.14043984469971: 1, 11.131091023782902: 1, 11.126078013882193: 1, 11.122724629238608: 1, 11.111935463604752: 1, 11.1002583854761: 1, 11.098406565238308: 1, 11.077797347127625: 1, 11.066535454774154: 1, 11.058328319568544: 1, 11.053251791920259: 1, 11.048471204410347: 1, 11.022983509230768: 1, 11.019648396813283: 1, 11.017973303593534: 1, 11.002637461596635: 1, 10.992752299577575: 1, 10.971756362508897: 1, 10.952771355446508: 1, 10.93929378123835: 1, 10.927468584144084: 1, 10.92305937596842: 1, 10.915226733948765: 1, 10.91190746482878: 1, 10.898899970225585: 1, 10.842020411116485: 1, 10.837007022653479: 1, 10.83494194302959: 1, 10.82723243925775: 1, 10.822250198459283: 1, 10.80660525744153: 1, 10.798719409884525: 1, 10.79608765980087: 1, 10.772265393023572: 1, 10.746495388607373: 1, 10.731973897391008: 1, 10.720036642681565: 1, 10.71734147540132: 1, 10.714520272244409: 1, 10.711022767943929: 1, 10.69776280696946: 1, 10.688292719460355: 1, 10.6816724286573: 1, 10.676910779573607: 1, 10.668215993047424: 1, 10.650640564253704: 1, 10.636655749634693: 1, 10.633362816235854: 1, 10.630777513987844: 1, 10.610896528011088: 1, 10.587363806494638: 1, 10.58481320585541: 1, 10.529053388741191: 1, 10.484748384875957: 1, 10.483967327972602: 1, 10.46562903934567: 1, 10.419679818577958: 1, 10.400880341611538: 1, 10.398465646962817: 1, 10.387686067296915: 1, 10.367070174983736: 1, 10.365869724683543: 1, 10.3625963222241: 1, 10.337926092377971: 1, 10.324436183583366: 1, 10.292828118666035: 1, 10.277545253654795: 1, 10.27642230208686: 1, 10.27143561219456: 1, 10.262627017756156: 1, 10.248108095723104: 1, 10.23901455722206: 1, 10.223984039443609: 1, 10.213761536945325: 1, 10.205081309655629: 1, 10.201875292948472: 1, 10.186880822943413: 1, 10.182549695674371: 1, 10.178335012145508: 1, 10.151207398461697: 1, 10.145785746762735: 1, 10.13179421133373: 1, 10.131046643710444: 1, 10.129165102123928: 1, 10.106479762605654: 1, 10.098419772961636: 1, 10.089376540255518: 1, 10.080275303042571: 1, 10.069708685512403: 1, 10.06911638952607: 1, 10.063232108092855: 1, 10.029035549028075: 1, 10.019172997847736: 1, 10.016033660864816: 1, 10.009188455403212: 1, 9.992051899261973: 1, 9.979132903988804: 1, 9.977260461058338: 1, 9.93830013150404: 1, 9.93496679136822: 1, 9.93045464926518: 1, 9.927056624702159: 1, 9.920413843166578: 1, 9.90804431367415: 1, 9.893223708760447: 1, 9.885044577784765: 1, 9.865368830177625: 1, 9.86221087051203: 1

9.857233256413005: 1, 9.854057169724117: 1, 9.83297807536788: 1, 9.829763391216083: 1, 9.811719045828994: 1, 9.801323547487495: 1, 9.769076066036488: 1, 9.756155242649355: 1, 9.753624293193118: 1, 9.752910650715545: 1, 9.745064590485649: 1, 9.716053079263714: 1, 9.700977648460498: 1, 9.676390434470399: 1, 9.65304549952451: 1, 9.630352484101033: 1, 9.625237089801272: 1, 9.622698800625864: 1, 9.614917073606032: 1, 9.603012200020135: 1, 9.58937121744456: 1, 9.585597178495975: 1, 9.568077005955768: 1, 9.564628782133962: 1, 9.553931183815129: 1, 9.55262561140146: 1, 9.541197076268979: 1, 9.534644397831208: 1, 9.527298016897305: 1, 9.52216346198049: 1, 9.516583590073129: 1, 9.478898663009689: 1, 9.44695215714148: 1, 9.445387745607096: 1, 9.420411658963982: 1, 9.41893010134917: 1, 9.392625278727559: 1, 9.391776105779469: 1, 9.377942796400518: 1, 9.369664739559047: 1, 9.365109443359131: 1, 9.358550342084051: 1, 9.351344113520241: 1, 9.31429245116283: 1, 9.309044182746662: 1, 9.299505102797324: 1, 9.299332461792705: 1, 9.295081849924603: 1, 9.290495084074937: 1, 9.284903734204981: 1, 9.277664976560324: 1, 9.274690186042541: 1, 9.263823879392993: 1, 9.257727391717216: 1, 9.254929770795766: 1, 9.25302132022473: 1, 9.246885192317013: 1, 9.24656926789055: 1, 9.246384471099342: 1, 9.234853106209357: 1, 9.226570518765953: 1, 9.22322268499077: 1, 9.221834106683927: 1, 9.213687024820219: 1, 9.212151076010237: 1, 9.21069861793393: 1, 9.192268764898294: 1, 9.188883840681802: 1, 9.160313233206457: 1, 9.152636438537565: 1, 9.149845480389265: 1, 9.14634031590986: 1, 9.107388290829112: 1, 9.08264902306266: 1, 9.082210029707985: 1, 9.070238866230866: 1, 9.063060946541981: 1, 9.051592325455035: 1, 9.047025971348441: 1, 9.045143692457312: 1, 9.044388652035252: 1, 9.043942076525905: 1, 9.030718022363507: 1, 9.027374548097642: 1, 9.01486386692439: 1, 9.010009992760493: 1, 9.003572291557195: 1, 9.001072138932571: 1, 8.997782217033421: 1, 8.99207915118547: 1, 8.990043870949616: 1, 8.969952678254042: 1, 8.967982552302216: 1, 8.95872207903843: 1, 8.938618505666533: 1, 8.918157337122084: 1, 8.881026576958876: 1, 8.878387990009061: 1, 8.878075581545357: 1, 8.859710688087135: 1, 8.859063516485634: 1, 8.828842032511828: 1, 8.823045684453577: 1, 8.821714601274186: 1, 8.810819577725287: 1, 8.804038632564442: 1, 8.793883123997196: 1, 8.782105892515547: 1, 8.778048920179819: 1, 8.776881018696573: 1, 8.756932594589449: 1, 8.75332719086131: 1, 8.723245824952578: 1, 8.719102821322045: 1, 8.717901423184813: 1, 8.688386727903962: 1, 8.669128607624017: 1, 8.663126722906956: 1, 8.650082947107162: 1, 8.64069229128806: 1, 8.636102203575344: 1, 8.632266288499471: 1, 8.623768060500325: 1, 8.62342345630528: 1, 8.621970483039874: 1, 8.612923162303321: 1, 8.605342172430964: 1, 8.591559323056323: 1, 8.585375998589193: 1, 8.562785466635933: 1, 8.560008056505538: 1, 8.557130551657197: 1, 8.542009448554058: 1, 8.536963877921178: 1, 8.504406434001627: 1, 8.501910319017997: 1, 8.500987504211372: 1, 8.448361934701392: 1, 8.432288947165043: 1, 8.400473579109335: 1, 8.385253138510702: 1, 8.37720180050142: 1, 8.371670892552087: 1, 8.362708177225175: 1, 8.351975791292618: 1, 8.348620240375752: 1, 8.341660057169738: 1, 8.338178467815082: 1, 8.324230955646872: 1, 8.317065445703888: 1, 8.309460224922267: 1, 8.297616990934717: 1, 8.296745878172457: 1, 8.293642196260274: 1, 8.262679930972425: 1, 8.252119481973635: 1, 8.243867778663118: 1, 8.242823023966512: 1, 8.236027002691339: 1, 8.217629716267378: 1, 8.21635856914741: 1, 8.213862323193915: 1, 8.213510354395915: 1, 8.20310477718944: 1, 8.189118213841905: 1, 8.171157077474117: 1, 8.167389899169601: 1, 8.161974081306555: 1, 8.15963110402869: 1, 8.1328850895338: 1, 8.131412531256768: 1, 8.130556957133532: 1, 8.120925595985893: 1, 8.10219290092524: 1, 8.102101301629734: 1, 8.08781064376597: 1, 8.084611653221216: 1, 8.068192175184498: 1, 8.052470856931896: 1, 8.026073679221838: 1, 8.02143997199466: 1, 8.02095629272802: 1, 8.0206056809269: 1, 8.017986244008378: 1, 8.008972540232191: 1, 7.9955706646393825: 1, 7.990284066943183: 1, 7.983968497714692: 1, 7.977341342294622: 1, 7.973057444603576: 1, 7.971663428800837: 1, 7.96680349667222: 1, 7.957358767415889: 1, 7.950752653010704: 1, 7.948659803182138: 1, 7.923784352248642: 1, 7.914708187886488: 1, 7.907190678786109: 1, 7.902972010443812: 1, 7.893236038788715: 1, 7.88574548500204: 1, 7.867343595187549: 1, 7.8529798589357265: 1, 7.82119001392506: 1, 7.816131173450084: 1, 7.767844065522989: 1, 7.753812382738819: 1, 7.727676952520121: 1, 7.723029275963992: 1, 7.72222375008493: 1, 7.70122642052673: 1, 7.695208944967322: 1, 7.6395439880342: 1, 7.636701867720308: 1, 7.59608190735593: 1, 7.569987150424848: 1, 7.56539428569615: 1, 7.5520594415855475: 1, 7.5427137826104484: 1, 7.523647242506462: 1, 7.479530166913577: 1, 7.4745343103127157: 1, 7.464687149857423: 1, 7.451954178378058: 1, 7.432365301009077: 1, 7.431826294263399: 1, 7.425196546789314: 1, 7.405407655867053: 1, 7.40282377479044: 1, 7.350143529153193: 1, 7.326459233614724: 1, 7.313528050144579: 1, 7.306586700033035: 1, 7.283813946263543: 1, 7.28098204831163: 1, 7.239929408870549: 1, 7.22817395100469: 1, 7.210503943085859: 1, 7.209324982796537: 1, 7.206555610181162: 1, 7.189544042202772: 1, 7.188708046778578: 1, 7.147572122035131: 1, 7.077279898442914: 1, 7.072702541741774: 1, 7.048016708688455: 1, 6.993603589254118: 1, 6.987776369120497: 1, 6.9760357790461365: 1, 6.95466247140416: 1, 6.897728410310045: 1, 6.87463547358919: 1, 6.870930497591952: 1, 6.859784847662268: 1, 6.82655466873579: 1, 6.82002350873249: 1, 6.7916210728313136: 1, 6.611351376543942: 1, 6.608292438217675: 1, 6.574525213895637: 1, 6.421310251836078: 1, 6.299693823083686: 1}})

In [47]:

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

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

```

# 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:
#-----

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

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

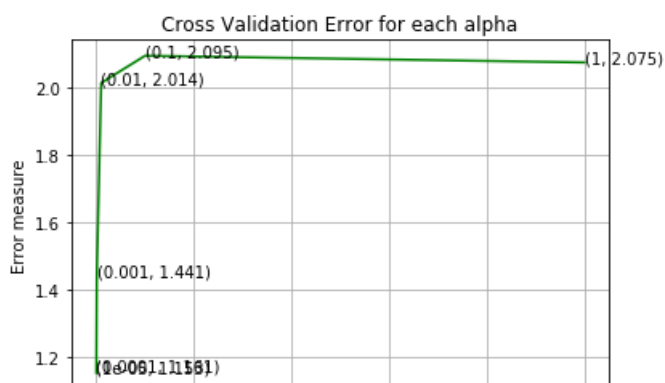
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best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

```

```

For values of alpha = 1e-05 The log loss is: 1.1532644383586574
For values of alpha = 0.0001 The log loss is: 1.1612679411635358
For values of alpha = 0.001 The log loss is: 1.4406696718788092
For values of alpha = 0.01 The log loss is: 2.014350554679497
For values of alpha = 0.1 The log loss is: 2.0950806538772286
For values of alpha = 1 The log loss is: 2.0751428702104078

```



0.0 0.2 0.4 0.6 0.8 1.0
Alpha i's

For values of best alpha = 1e-05 The train log loss is: 0.7412786504348623
For values of best alpha = 1e-05 The cross validation log loss is: 1.1532644383586574
For values of best alpha = 1e-05 The test log loss is: 1.1984396955548693

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

Ans. Yes, it seems like!

In [48]:

```
def get_intersec_text(df):
    df_text_vec = TfidfVectorizer(max_features = 1000, min_df = 3)
    df_text_fea = df_text_vec.fit_transform(df['TEXT'])
    df_text_features = df_text_vec.get_feature_names()

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

In [49]:

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

94.3 % of word of test data appeared in train data
93.6 % of word of Cross Validation appeared in train data

4. Machine Learning Models

In [63]:

```
#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 [64]:

```
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 [65]:

```
# 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_impreature_names(indices, text, gene, var, no_features):
    gene_count_vec = TfidfVectorizer(max_features = 1000)
    var_count_vec = TfidfVectorizer(max_features = 1000)
    text_count_vec = TfidfVectorizer(min_df=3,max_features = 1000)

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

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

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

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

```

Stacking the three types of features

In [66]:

```

# merging gene, variance and text features

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

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

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

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

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

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

```

```

np.hstack((cv_gene_feature_responseCoding,cv_variation_feature_responseCoding))

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

```

In [67]:

```

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, 2232)
(number of data points * number of features) in test data = (665, 2232)
(number of data points * number of features) in cross validation data = (532, 2232)

```

In [68]:

```

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 [69]:

```

# 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/
# -----

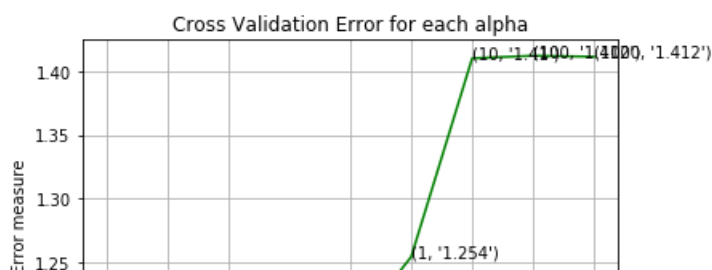
alpha = [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100, 1000]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = MultinomialNB(alpha=i)
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    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))

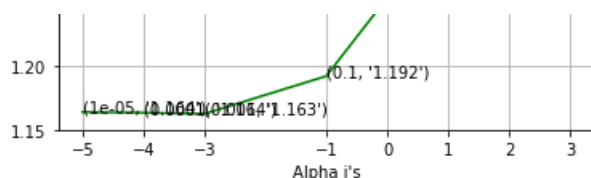
fig, ax = plt.subplots()
ax.plot(np.log10(alpha), cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], str(txt)), (np.log10(alpha[i]), cv_log_error_array[i]))
plt.grid()
plt.xticks(np.log10(alpha))
plt.title("Cross Validation Error for each alpha")
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_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, p
redict_y, labels=clf.classes_, eps=1e-15))
```

```
for alpha = 1e-05
Log Loss : 1.164111647557351
for alpha = 0.0001
Log Loss : 1.1636446127113975
for alpha = 0.001
Log Loss : 1.1627400999467232
for alpha = 0.1
Log Loss : 1.1924304270718766
for alpha = 1
Log Loss : 1.2541366464909083
for alpha = 10
Log Loss : 1.4103758317317077
for alpha = 100
Log Loss : 1.4124142354020242
for alpha = 1000
Log Loss : 1.4116634312711283
```





For values of best alpha = 0.001 The train log loss is: 0.7389550367162582
 For values of best alpha = 0.001 The cross validation log loss is: 1.1627400999467232
 For values of best alpha = 0.001 The test log loss is: 1.1886205094193842

4.1.1.2. Testing the model with best hyper paramters

In [70]:

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

# some of methods of MultinomialNB()
# fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
# predict(X) Perform classification on an array of test vectors X.
# predict_log_proba(X) Return log-probability estimates for the test vector X.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/naive-bayes-algorithm-1/
# -----

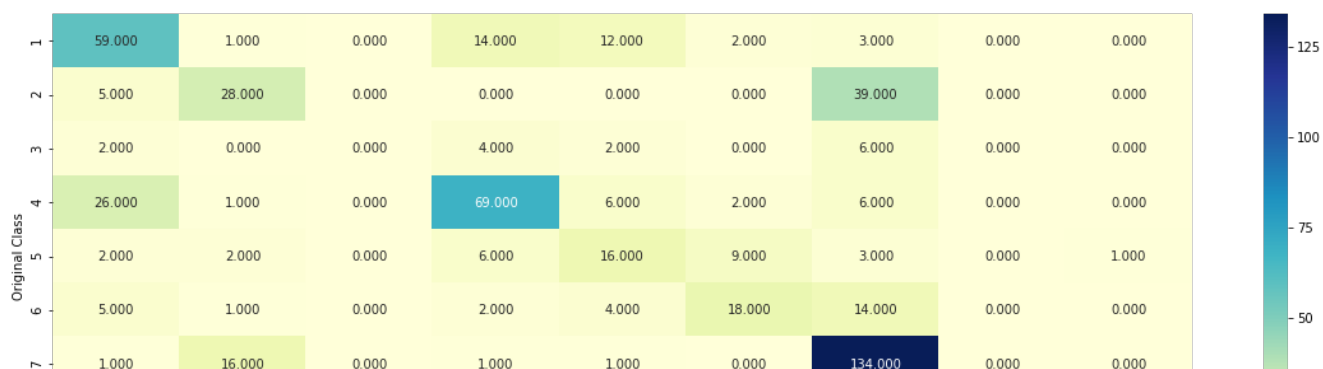
# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default paramters
# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
# some of the methods of CalibratedClassifierCV()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

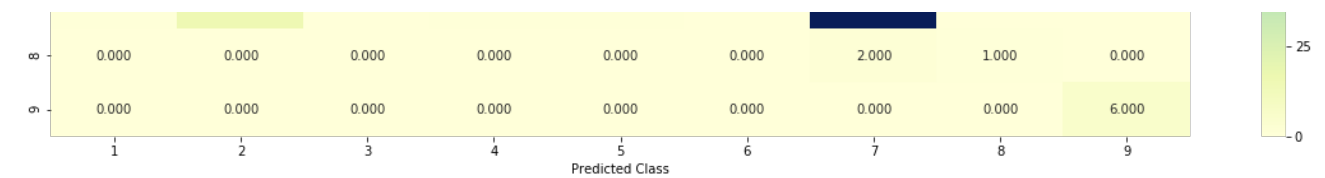
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
# to avoid rounding error while multiplying probabilties we use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

Log Loss : 1.1627400999467232

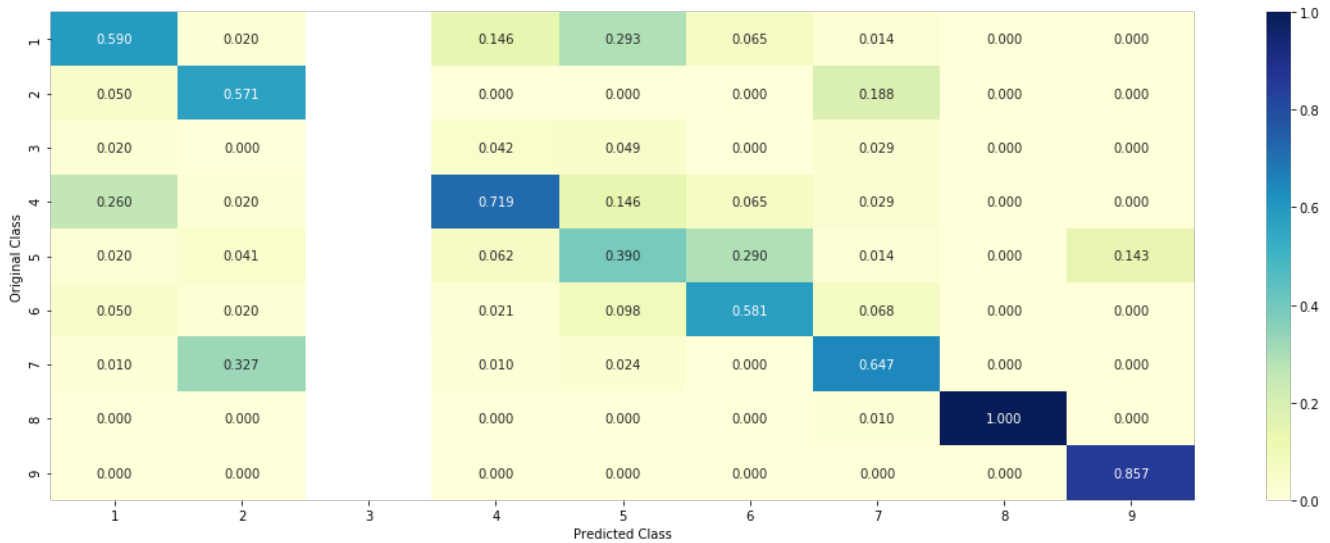
Number of missclassified point : 0.37781954887218044

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

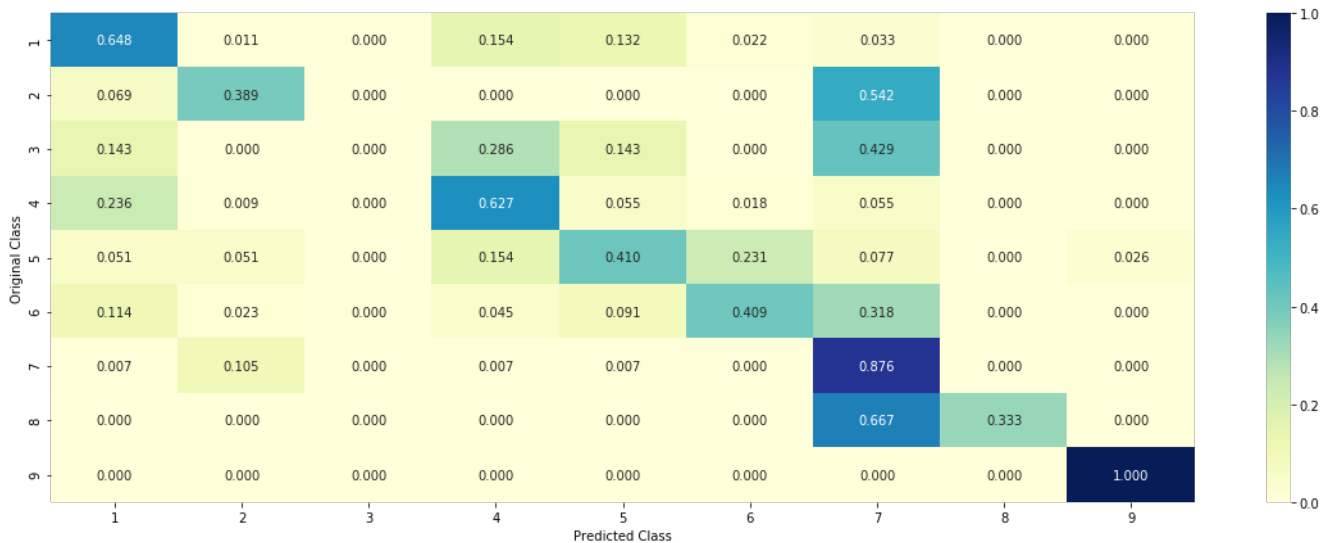




----- Precision matrix (Column Sum=1) -----



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



4.1.1.3. Feature Importance, Correctly classified point

In [71]:

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

Predicted Class : 4

Predicted Class Probabilities: [0.0576 0.0472 0.0122 0.7485 0.0216 0.0213 0.0644 0.0024 0.0027 1]

Predicted Class Probabilities: [[0.0376 0.0472 0.0122 0.7493 0.0516 0.0513 0.0644 0.0024 0.0037]]
Actual Class : 4

10 Text feature [activity] present in test data point [True]
11 Text feature [protein] present in test data point [True]
13 Text feature [proteins] present in test data point [True]
14 Text feature [function] present in test data point [True]
15 Text feature [missense] present in test data point [True]
16 Text feature [results] present in test data point [True]
17 Text feature [whereas] present in test data point [True]
18 Text feature [experiments] present in test data point [True]
19 Text feature [acid] present in test data point [True]
20 Text feature [amino] present in test data point [True]
21 Text feature [suppressor] present in test data point [True]
24 Text feature [type] present in test data point [True]
25 Text feature [also] present in test data point [True]
26 Text feature [functional] present in test data point [True]
27 Text feature [shown] present in test data point [True]
28 Text feature [whether] present in test data point [True]
29 Text feature [described] present in test data point [True]
31 Text feature [wild] present in test data point [True]
32 Text feature [two] present in test data point [True]
33 Text feature [important] present in test data point [True]
34 Text feature [related] present in test data point [True]
35 Text feature [reduced] present in test data point [True]
36 Text feature [mutations] present in test data point [True]
37 Text feature [determined] present in test data point [True]
39 Text feature [either] present in test data point [True]
40 Text feature [may] present in test data point [True]
42 Text feature [ability] present in test data point [True]
43 Text feature [although] present in test data point [True]
44 Text feature [three] present in test data point [True]
45 Text feature [previously] present in test data point [True]
46 Text feature [determine] present in test data point [True]
47 Text feature [loss] present in test data point [True]
48 Text feature [containing] present in test data point [True]
49 Text feature [discussion] present in test data point [True]
50 Text feature [associated] present in test data point [True]
51 Text feature [therefore] present in test data point [True]
52 Text feature [analysis] present in test data point [True]
53 Text feature [levels] present in test data point [True]
54 Text feature [thus] present in test data point [True]
57 Text feature [mammalian] present in test data point [True]
58 Text feature [vitro] present in test data point [True]
63 Text feature [one] present in test data point [True]
65 Text feature [effect] present in test data point [True]
66 Text feature [lower] present in test data point [True]
67 Text feature [contribute] present in test data point [True]
68 Text feature [show] present in test data point [True]
69 Text feature [site] present in test data point [True]
70 Text feature [however] present in test data point [True]
71 Text feature [30] present in test data point [True]
72 Text feature [several] present in test data point [True]
74 Text feature [effects] present in test data point [True]
75 Text feature [indicated] present in test data point [True]
77 Text feature [similar] present in test data point [True]
78 Text feature [could] present in test data point [True]
80 Text feature [binding] present in test data point [True]
81 Text feature [introduction] present in test data point [True]
82 Text feature [possible] present in test data point [True]
83 Text feature [10] present in test data point [True]
85 Text feature [mutant] present in test data point [True]
86 Text feature [general] present in test data point [True]
87 Text feature [cells] present in test data point [True]
88 Text feature [addition] present in test data point [True]
89 Text feature [dna] present in test data point [True]
90 Text feature [bind] present in test data point [True]
91 Text feature [mutation] present in test data point [True]
92 Text feature [changes] present in test data point [True]
93 Text feature [found] present in test data point [True]
94 Text feature [see] present in test data point [True]
95 Text feature [substitutions] present in test data point [True]
96 Text feature [50] present in test data point [True]
97 Text feature [stability] present in test data point [True]
98 Text feature [mutants] present in test data point [True]
99 Text feature [terminal] present in test data point [True]
Out of the top 100 features 73 are present in query point

4.1.1.4. Feature Importance, Incorrectly classified point

In [72]:

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

Predicted Class : 1

Predicted Class Probabilities: [[0.7247 0.0507 0.0131 0.0682 0.0339 0.0335 0.0694 0.0026 0.004]]

Actual Class : 1

```
-----
11 Text feature [one] present in test data point [True]
13 Text feature [results] present in test data point [True]
14 Text feature [therefore] present in test data point [True]
15 Text feature [protein] present in test data point [True]
16 Text feature [function] present in test data point [True]
17 Text feature [also] present in test data point [True]
18 Text feature [two] present in test data point [True]
19 Text feature [specific] present in test data point [True]
20 Text feature [however] present in test data point [True]
21 Text feature [role] present in test data point [True]
22 Text feature [type] present in test data point [True]
23 Text feature [region] present in test data point [True]
24 Text feature [table] present in test data point [True]
25 Text feature [whether] present in test data point [True]
26 Text feature [control] present in test data point [True]
27 Text feature [using] present in test data point [True]
28 Text feature [human] present in test data point [True]
29 Text feature [functions] present in test data point [True]
31 Text feature [may] present in test data point [True]
32 Text feature [either] present in test data point [True]
33 Text feature [shown] present in test data point [True]
34 Text feature [well] present in test data point [True]
35 Text feature [least] present in test data point [True]
36 Text feature [dna] present in test data point [True]
37 Text feature [essential] present in test data point [True]
38 Text feature [possible] present in test data point [True]
39 Text feature [discussion] present in test data point [True]
40 Text feature [deletion] present in test data point [True]
42 Text feature [wild] present in test data point [True]
43 Text feature [performed] present in test data point [True]
44 Text feature [effect] present in test data point [True]
45 Text feature [result] present in test data point [True]
46 Text feature [suggest] present in test data point [True]
48 Text feature [25] present in test data point [True]
49 Text feature [large] present in test data point [True]
50 Text feature [although] present in test data point [True]
52 Text feature [gene] present in test data point [True]
53 Text feature [proteins] present in test data point [True]
54 Text feature [expression] present in test data point [True]
55 Text feature [determined] present in test data point [True]
56 Text feature [containing] present in test data point [True]
57 Text feature [present] present in test data point [True]
58 Text feature [including] present in test data point [True]
59 Text feature [within] present in test data point [True]
60 Text feature [analysis] present in test data point [True]
61 Text feature [three] present in test data point [True]
62 Text feature [fig] present in test data point [True]
63 Text feature [similar] present in test data point [True]
64 Text feature [important] present in test data point [True]
65 Text feature [reduced] present in test data point [True]
66 Text feature [used] present in test data point [True]
67 Text feature [cell] present in test data point [True]
68 Text feature [mediated] present in test data point [True]
```

```

68 Text feature [indicated] present in test data point [True]
69 Text feature [another] present in test data point [True]
70 Text feature [data] present in test data point [True]
71 Text feature [different] present in test data point [True]
72 Text feature [together] present in test data point [True]
73 Text feature [previously] present in test data point [True]
75 Text feature [several] present in test data point [True]
76 Text feature [addition] present in test data point [True]
77 Text feature [indicated] present in test data point [True]
78 Text feature [example] present in test data point [True]
79 Text feature [additional] present in test data point [True]
80 Text feature [likely] present in test data point [True]
81 Text feature [described] present in test data point [True]
82 Text feature [30] present in test data point [True]
83 Text feature [whereas] present in test data point [True]
84 Text feature [15] present in test data point [True]
85 Text feature [significant] present in test data point [True]
86 Text feature [corresponding] present in test data point [True]
87 Text feature [observed] present in test data point [True]
88 Text feature [mutations] present in test data point [True]
89 Text feature [cancer] present in test data point [True]
90 Text feature [associated] present in test data point [True]
91 Text feature [previous] present in test data point [True]
92 Text feature [required] present in test data point [True]
93 Text feature [yet] present in test data point [True]
94 Text feature [furthermore] present in test data point [True]
95 Text feature [directly] present in test data point [True]
96 Text feature [indicating] present in test data point [True]
97 Text feature [first] present in test data point [True]
98 Text feature [respectively] present in test data point [True]
99 Text feature [ability] present in test data point [True]
Out of the top 100 features 83 are present in query point

```

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

In [73]:

```

# 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-ne
ighbors-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 = []

```

```

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

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

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

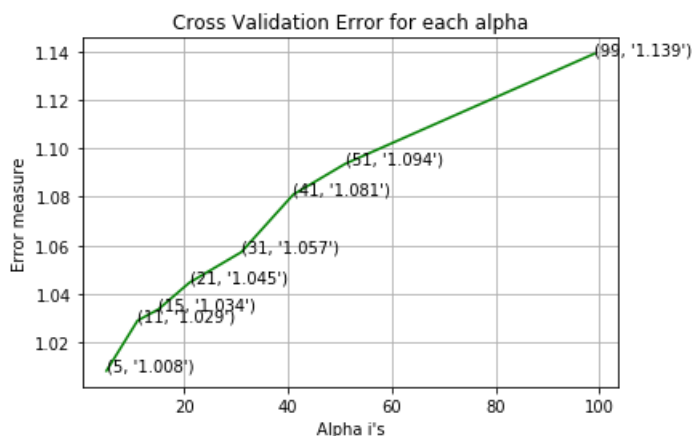
predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 5
Log Loss : 1.0082901840724636
for alpha = 11
Log Loss : 1.0290099490795332
for alpha = 15
Log Loss : 1.0336422613933147
for alpha = 21
Log Loss : 1.044671240525676
for alpha = 31
Log Loss : 1.0572406169994706
for alpha = 41
Log Loss : 1.080967997199732
for alpha = 51
Log Loss : 1.0936172427068644
for alpha = 99
Log Loss : 1.1389937619186208

```



```

For values of best alpha = 5 The train log loss is: 0.4755178649543858
For values of best alpha = 5 The cross validation log loss is: 1.0082901840724636
For values of best alpha = 5 The test log loss is: 1.1184776899071618

```

For values of best_alpha = 0 the loss log loss is: 1.1101760000/1010

4.2.2. Testing the model with best hyper paramters

In [74]:

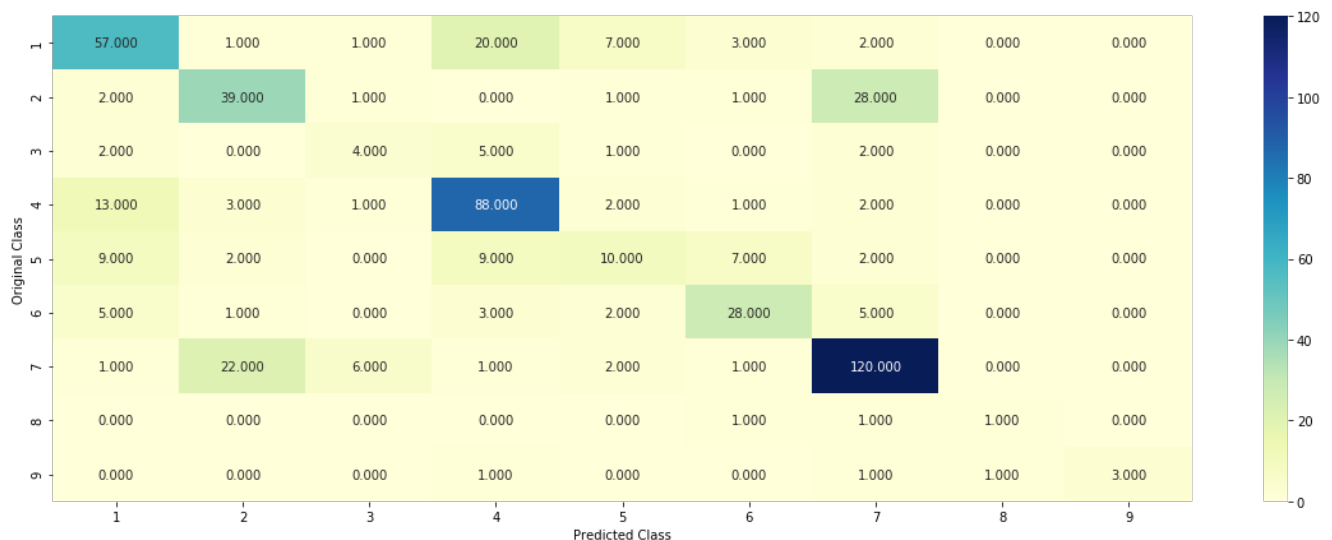
```
# find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None,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

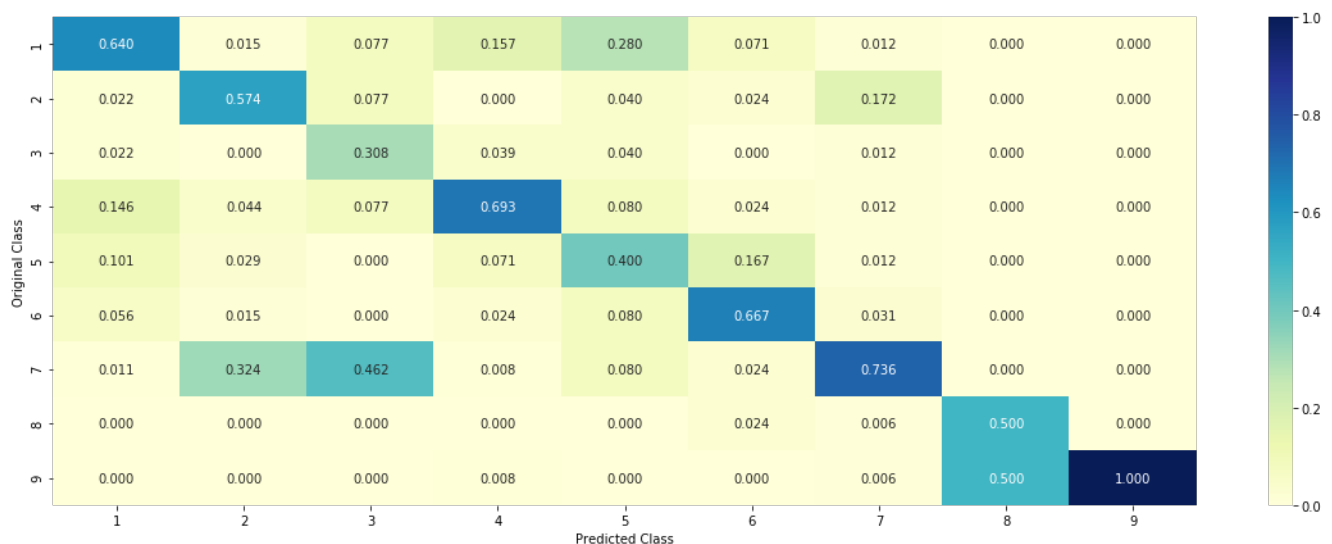
Log loss : 1.0082901840724636

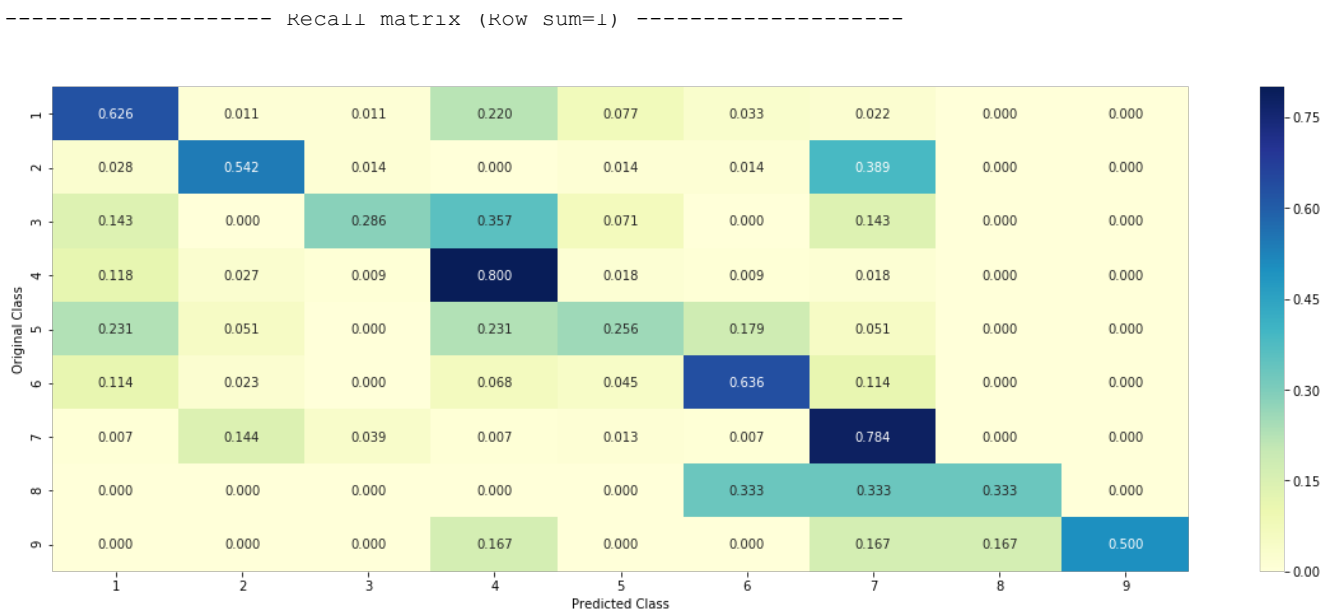
Number of mis-classified points : 0.34210526315789475

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



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





4.2.3. Sample Query point -1

In [75]:

```
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 : 4

Actual Class : 4

The 5 nearest neighbours of the test points belongs to classes [4 4 4 4 4]

Fequency of nearest points : Counter({4: 5})

4.2.4. Sample Query Point-2

In [76]:

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

test_point_index = 100

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

Predicted Class : 1

Actual Class : 1

the k value for knn is 5 and the nearest neighbours of the test points belongs to classes [1 1 1 1 1]

Frequency of nearest points : Counter({1: 5})

4.3. Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper paramter tuning

In [77]:

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

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

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

# 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 probabilitites 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)

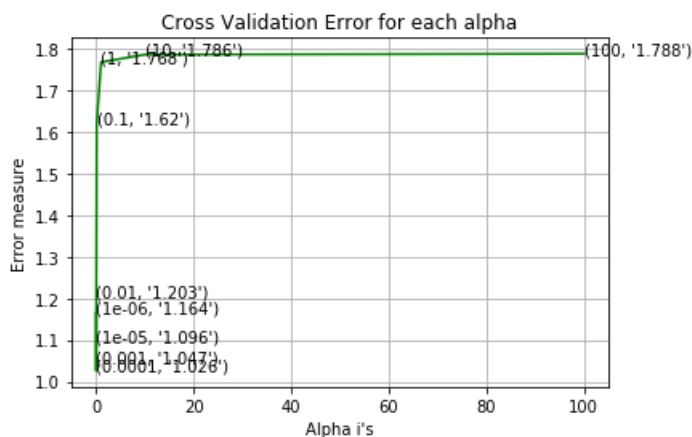
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 1e-06
Log Loss : 1.1642713953663995
for alpha = 1e-05
Log Loss : 1.0958173677051277
for alpha = 0.0001
Log Loss : 1.0255817916912926
for alpha = 0.001
Log Loss : 1.0468027218518803
for alpha = 0.01
Log Loss : 1.2031750288141916
for alpha = 0.1
Log Loss : 1.6202351308931224
for alpha = 1
Log Loss : 1.7679919763730352
for alpha = 10
Log Loss : 1.7862074549485683
for alpha = 100
Log Loss : 1.7882999037580773

```



```

For values of best alpha = 0.0001 The train log loss is: 0.5861623344161264
For values of best alpha = 0.0001 The cross validation log loss is: 1.0255817916912926
For values of best alpha = 0.0001 The test log loss is: 1.044162743983839

```

4.3.1.2. Testing the model with best hyper paramters

In [78]:

```

# 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/
#-----
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 : 1.0255817916912926

Number of mis-classified points : 0.35150375939849626

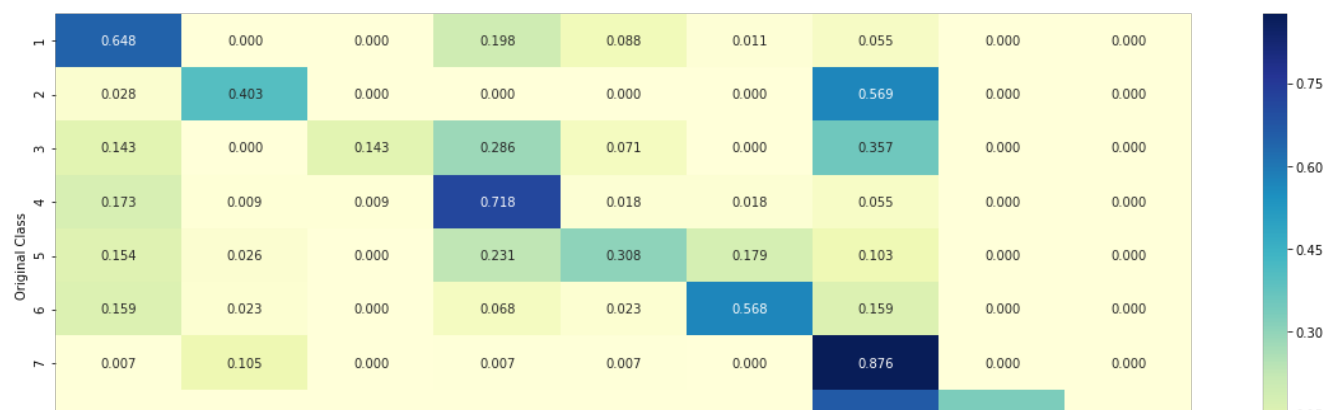
----- Confusion matrix -----

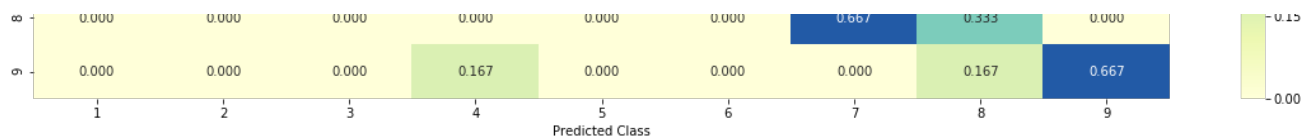


----- Precision matrix (Column Sum=1) -----



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





4.3.1.3. Feature Importance

In [79]:

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

4.3.1.3.1. Correctly Classified point

In [80]:

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

Predicted Class : 4

Predicted Class Probabilities: [[0.0324 0.023 0.0067 0.8869 0.0195 0.0131 0.0148 0.0018 0.0017]]

Actual Class : 4

```
-----
3 Text feature [suppressor] present in test data point [True]
17 Text feature [missense] present in test data point [True]
38 Text feature [families] present in test data point [True]
73 Text feature [mm] present in test data point [True]
80 Text feature [inactivation] present in test data point [True]
89 Text feature [mammalian] present in test data point [True]
92 Text feature [see] present in test data point [True]
99 Text feature [protein] present in test data point [True]
105 Text feature [yeast] present in test data point [True]
112 Text feature [null] present in test data point [True]
122 Text feature [proportion] present in test data point [True]
125 Text feature [homozygous] present in test data point [True]
133 Text feature [reduced] present in test data point [True]
147 Text feature [dna] present in test data point [True]
153 Text feature [consequences] present in test data point [True]
154 Text feature [family] present in test data point [True]
```

165 Text feature [localization] present in test data point [True]
168 Text feature [changes] present in test data point [True]
173 Text feature [germline] present in test data point [True]
189 Text feature [functional] present in test data point [True]
193 Text feature [appears] present in test data point [True]
194 Text feature [loss] present in test data point [True]
195 Text feature [require] present in test data point [True]
205 Text feature [represent] present in test data point [True]
219 Text feature [bind] present in test data point [True]
222 Text feature [cycle] present in test data point [True]
224 Text feature [kinases] present in test data point [True]
225 Text feature [show] present in test data point [True]
226 Text feature [38] present in test data point [True]
228 Text feature [specifically] present in test data point [True]
230 Text feature [plates] present in test data point [True]
232 Text feature [determine] present in test data point [True]
233 Text feature [1998] present in test data point [True]
234 Text feature [contribute] present in test data point [True]
237 Text feature [stability] present in test data point [True]
249 Text feature [plasmid] present in test data point [True]
255 Text feature [stained] present in test data point [True]
263 Text feature [2010] present in test data point [True]
266 Text feature [high] present in test data point [True]
271 Text feature [investigated] present in test data point [True]
279 Text feature [displayed] present in test data point [True]
280 Text feature [affected] present in test data point [True]
282 Text feature [defective] present in test data point [True]
291 Text feature [bound] present in test data point [True]
292 Text feature [risk] present in test data point [True]
293 Text feature [26] present in test data point [True]
294 Text feature [cannot] present in test data point [True]
299 Text feature [particular] present in test data point [True]
302 Text feature [several] present in test data point [True]
304 Text feature [western] present in test data point [True]
305 Text feature [cases] present in test data point [True]
309 Text feature [function] present in test data point [True]
314 Text feature [transfected] present in test data point [True]
315 Text feature [terminal] present in test data point [True]
321 Text feature [cellular] present in test data point [True]
324 Text feature [site] present in test data point [True]
333 Text feature [view] present in test data point [True]
334 Text feature [splice] present in test data point [True]
338 Text feature [early] present in test data point [True]
339 Text feature [observed] present in test data point [True]
340 Text feature [amino] present in test data point [True]
342 Text feature [relative] present in test data point [True]
343 Text feature [lower] present in test data point [True]
345 Text feature [comparison] present in test data point [True]
346 Text feature [red] present in test data point [True]
353 Text feature [37] present in test data point [True]
356 Text feature [associated] present in test data point [True]
362 Text feature [figure] present in test data point [True]
363 Text feature [mutants] present in test data point [True]
368 Text feature [described] present in test data point [True]
371 Text feature [altered] present in test data point [True]
375 Text feature [sds] present in test data point [True]
377 Text feature [involved] present in test data point [True]
380 Text feature [nuclear] present in test data point [True]
385 Text feature [inhibitory] present in test data point [True]
386 Text feature [groups] present in test data point [True]
392 Text feature [mouse] present in test data point [True]
394 Text feature [2a] present in test data point [True]
400 Text feature [low] present in test data point [True]
402 Text feature [transfection] present in test data point [True]
404 Text feature [cdk4] present in test data point [True]
405 Text feature [linked] present in test data point [True]
408 Text feature [antibodies] present in test data point [True]
409 Text feature [3b] present in test data point [True]
412 Text feature [similarly] present in test data point [True]
413 Text feature [times] present in test data point [True]
415 Text feature [isolated] present in test data point [True]
416 Text feature [2000] present in test data point [True]
418 Text feature [45] present in test data point [True]
421 Text feature [containing] present in test data point [True]
429 Text feature [genomic] present in test data point [True]
430 Text feature [buffer] present in test data point [True]
435 Text feature [indeed] present in test data point [True]

```

436 Text feature [40] present in test data point [True]
437 Text feature [pathogenic] present in test data point [True]
441 Text feature [key] present in test data point [True]
443 Text feature [phosphorylation] present in test data point [True]
444 Text feature [alterations] present in test data point [True]
445 Text feature [larger] present in test data point [True]
449 Text feature [decrease] present in test data point [True]
453 Text feature [partial] present in test data point [True]
454 Text feature [ca] present in test data point [True]
455 Text feature [flag] present in test data point [True]
457 Text feature [induced] present in test data point [True]
459 Text feature [genetic] present in test data point [True]
460 Text feature [29] present in test data point [True]
462 Text feature [via] present in test data point [True]
463 Text feature [gst] present in test data point [True]
465 Text feature [locus] present in test data point [True]
466 Text feature [individuals] present in test data point [True]
469 Text feature [proteins] present in test data point [True]
471 Text feature [distribution] present in test data point [True]
472 Text feature [reporter] present in test data point [True]
473 Text feature [assay] present in test data point [True]
474 Text feature [average] present in test data point [True]
481 Text feature [deletion] present in test data point [True]
483 Text feature [variant] present in test data point [True]
485 Text feature [whether] present in test data point [True]
488 Text feature [consistent] present in test data point [True]
490 Text feature [green] present in test data point [True]
496 Text feature [predicted] present in test data point [True]
499 Text feature [assays] present in test data point [True]
Out of the top 500 features 122 are present in query point

```

4.3.1.3.2. Incorrectly Classified point

In [81]:

```

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

```

```

Predicted Class : 1
Predicted Class Probabilities: [[9.648e-01 1.880e-02 1.000e-04 1.380e-02 1.300e-03 2.000e-04 9.000
e-04

```

```

1.000e-04 0.000e+00]]

```

```

Actual Class : 1

```

```

-----
10 Text feature [panel] present in test data point [True]
52 Text feature [nuclear] present in test data point [True]
62 Text feature [deletion] present in test data point [True]
65 Text feature [repeats] present in test data point [True]
68 Text feature [rather] present in test data point [True]
99 Text feature [ability] present in test data point [True]
109 Text feature [therefore] present in test data point [True]
112 Text feature [defined] present in test data point [True]
117 Text feature [box] present in test data point [True]
121 Text feature [essential] present in test data point [True]
141 Text feature [functions] present in test data point [True]
144 Text feature [effect] present in test data point [True]
146 Text feature [21] present in test data point [True]
150 Text feature [deficient] present in test data point [True]
152 Text feature [corresponding] present in test data point [True]
154 Text feature [region] present in test data point [True]
155 Text feature [strong] present in test data point [True]
156 Text feature [analyses] present in test data point [True]
163 Text feature [mutational] present in test data point [True]
165 Text feature [17] present in test data point [True]
172 Text feature [panel] present in test data point [True]

```

173 Text feature [page] present in test data point [True]
179 Text feature [calculated] present in test data point [True]
181 Text feature [indicated] present in test data point [True]
182 Text feature [notably] present in test data point [True]
187 Text feature [whole] present in test data point [True]
188 Text feature [gel] present in test data point [True]
191 Text feature [subjected] present in test data point [True]
192 Text feature [function] present in test data point [True]
194 Text feature [position] present in test data point [True]
195 Text feature [value] present in test data point [True]
196 Text feature [sequenced] present in test data point [True]
197 Text feature [change] present in test data point [True]
201 Text feature [impaired] present in test data point [True]
203 Text feature [identify] present in test data point [True]
204 Text feature [located] present in test data point [True]
206 Text feature [encoding] present in test data point [True]
218 Text feature [previous] present in test data point [True]
221 Text feature [percentage] present in test data point [True]
222 Text feature [showing] present in test data point [True]
223 Text feature [blood] present in test data point [True]
226 Text feature [splice] present in test data point [True]
227 Text feature [interact] present in test data point [True]
228 Text feature [splicing] present in test data point [True]
231 Text feature [set] present in test data point [True]
232 Text feature [another] present in test data point [True]
235 Text feature [next] present in test data point [True]
237 Text feature [sequences] present in test data point [True]
240 Text feature [least] present in test data point [True]
243 Text feature [localization] present in test data point [True]
244 Text feature [residues] present in test data point [True]
251 Text feature [performed] present in test data point [True]
252 Text feature [possibility] present in test data point [True]
254 Text feature [upon] present in test data point [True]
257 Text feature [induce] present in test data point [True]
258 Text feature [supplemental] present in test data point [True]
259 Text feature [signals] present in test data point [True]
260 Text feature [values] present in test data point [True]
262 Text feature [tested] present in test data point [True]
263 Text feature [conserved] present in test data point [True]
264 Text feature [primers] present in test data point [True]
266 Text feature [obtained] present in test data point [True]
268 Text feature [genome] present in test data point [True]
269 Text feature [strand] present in test data point [True]
270 Text feature [via] present in test data point [True]
271 Text feature [length] present in test data point [True]
274 Text feature [treated] present in test data point [True]
276 Text feature [signal] present in test data point [True]
277 Text feature [specific] present in test data point [True]
281 Text feature [possible] present in test data point [True]
282 Text feature [sample] present in test data point [True]
283 Text feature [defects] present in test data point [True]
285 Text feature [insertion] present in test data point [True]
286 Text feature [16] present in test data point [True]
287 Text feature [de] present in test data point [True]
289 Text feature [carrying] present in test data point [True]
290 Text feature [development] present in test data point [True]
291 Text feature [sequencing] present in test data point [True]
294 Text feature [medium] present in test data point [True]
295 Text feature [coding] present in test data point [True]
299 Text feature [size] present in test data point [True]
300 Text feature [remaining] present in test data point [True]
301 Text feature [taken] present in test data point [True]
303 Text feature [large] present in test data point [True]
304 Text feature [screening] present in test data point [True]
305 Text feature [chain] present in test data point [True]
306 Text feature [exon] present in test data point [True]
309 Text feature [wild] present in test data point [True]
310 Text feature [within] present in test data point [True]
311 Text feature [vitro] present in test data point [True]
317 Text feature [following] present in test data point [True]
321 Text feature [type] present in test data point [True]
322 Text feature [rt] present in test data point [True]
323 Text feature [frame] present in test data point [True]
327 Text feature [mean] present in test data point [True]
328 Text feature [rna] present in test data point [True]
329 Text feature [incubated] present in test data point [True]
332 Text feature [nucleotide] present in test data point [True]
333 Text feature [catalyze] present in test data point [True]

333 Text feature [del] present in test data point [True]
334 Text feature [dependent] present in test data point [True]
336 Text feature [table] present in test data point [True]
339 Text feature [impact] present in test data point [True]
341 Text feature [selection] present in test data point [True]
342 Text feature [one] present in test data point [True]
347 Text feature [somatic] present in test data point [True]
348 Text feature [fold] present in test data point [True]
349 Text feature [heterozygous] present in test data point [True]
350 Text feature [cell] present in test data point [True]
351 Text feature [manufacturer] present in test data point [True]
352 Text feature [additional] present in test data point [True]
353 Text feature [33] present in test data point [True]
357 Text feature [49] present in test data point [True]
360 Text feature [among] present in test data point [True]
361 Text feature [control] present in test data point [True]
362 Text feature [none] present in test data point [True]
363 Text feature [context] present in test data point [True]
371 Text feature [24] present in test data point [True]
373 Text feature [absence] present in test data point [True]
374 Text feature [exhibited] present in test data point [True]
375 Text feature [complex] present in test data point [True]
376 Text feature [established] present in test data point [True]
377 Text feature [deletions] present in test data point [True]
379 Text feature [stable] present in test data point [True]
380 Text feature [2001] present in test data point [True]
381 Text feature [sporadic] present in test data point [True]
385 Text feature [pcr] present in test data point [True]
387 Text feature [reverse] present in test data point [True]
388 Text feature [individual] present in test data point [True]
389 Text feature [domains] present in test data point [True]
396 Text feature [derived] present in test data point [True]
398 Text feature [reduced] present in test data point [True]
400 Text feature [analyzed] present in test data point [True]
401 Text feature [05] present in test data point [True]
402 Text feature [role] present in test data point [True]
403 Text feature [54] present in test data point [True]
404 Text feature [right] present in test data point [True]
405 Text feature [non] present in test data point [True]
406 Text feature [example] present in test data point [True]
407 Text feature [32] present in test data point [True]
409 Text feature [cause] present in test data point [True]
412 Text feature [allele] present in test data point [True]
413 Text feature [early] present in test data point [True]
414 Text feature [recently] present in test data point [True]
416 Text feature [evidence] present in test data point [True]
417 Text feature [complete] present in test data point [True]
419 Text feature [22] present in test data point [True]
420 Text feature [55] present in test data point [True]
421 Text feature [http] present in test data point [True]
422 Text feature [indeed] present in test data point [True]
423 Text feature [population] present in test data point [True]
425 Text feature [1997] present in test data point [True]
430 Text feature [assay] present in test data point [True]
432 Text feature [distinct] present in test data point [True]
435 Text feature [might] present in test data point [True]
436 Text feature [key] present in test data point [True]
438 Text feature [could] present in test data point [True]
439 Text feature [protein] present in test data point [True]
440 Text feature [results] present in test data point [True]
443 Text feature [revealed] present in test data point [True]
445 Text feature [residue] present in test data point [True]
446 Text feature [1b] present in test data point [True]
450 Text feature [25] present in test data point [True]
451 Text feature [lines] present in test data point [True]
452 Text feature [exons] present in test data point [True]
453 Text feature [generation] present in test data point [True]
454 Text feature [19] present in test data point [True]
455 Text feature [transcription] present in test data point [True]
458 Text feature [presence] present in test data point [True]
459 Text feature [furthermore] present in test data point [True]
466 Text feature [larger] present in test data point [True]
467 Text feature [finally] present in test data point [True]
468 Text feature [effects] present in test data point [True]
470 Text feature [limited] present in test data point [True]
471 Text feature [predicted] present in test data point [True]
473 Text feature [reported] present in test data point [True]
475 Text feature [] present in test data point [True]


```

475 Text feature [nm] present in test data point [True]
476 Text feature [human] present in test data point [True]
478 Text feature [proteins] present in test data point [True]
480 Text feature [respectively] present in test data point [True]
481 Text feature [criteria] present in test data point [True]
484 Text feature [using] present in test data point [True]
485 Text feature [based] present in test data point [True]
488 Text feature [43] present in test data point [True]
490 Text feature [30] present in test data point [True]
491 Text feature [51] present in test data point [True]
492 Text feature [sensitivity] present in test data point [True]
496 Text feature [cohort] present in test data point [True]
497 Text feature [42] present in test data point [True]
499 Text feature [common] present in test data point [True]
Out of the top 500 features 188 are present in query point

```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

In [82]:

```

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

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

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

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

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

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

```

```

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

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

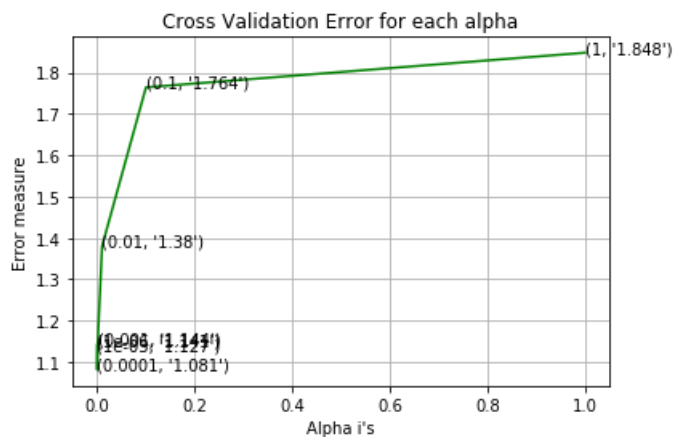
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 1e-06
Log Loss : 1.1412991489533721
for alpha = 1e-05
Log Loss : 1.127222351318073
for alpha = 0.0001
Log Loss : 1.0809042565760936
for alpha = 0.001
Log Loss : 1.1442140220862396
for alpha = 0.01
Log Loss : 1.3801612857584817
for alpha = 0.1
Log Loss : 1.7642984247364517
for alpha = 1
Log Loss : 1.8480675582772192

```



```

For values of best alpha = 0.0001 The train log loss is: 0.5784459460818665
For values of best alpha = 0.0001 The cross validation log loss is: 1.0809042565760936
For values of best alpha = 0.0001 The test log loss is: 1.088411515127868

```

4.3.2.2. Testing model with best hyper parameters

In [83]:

```

# 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.
```

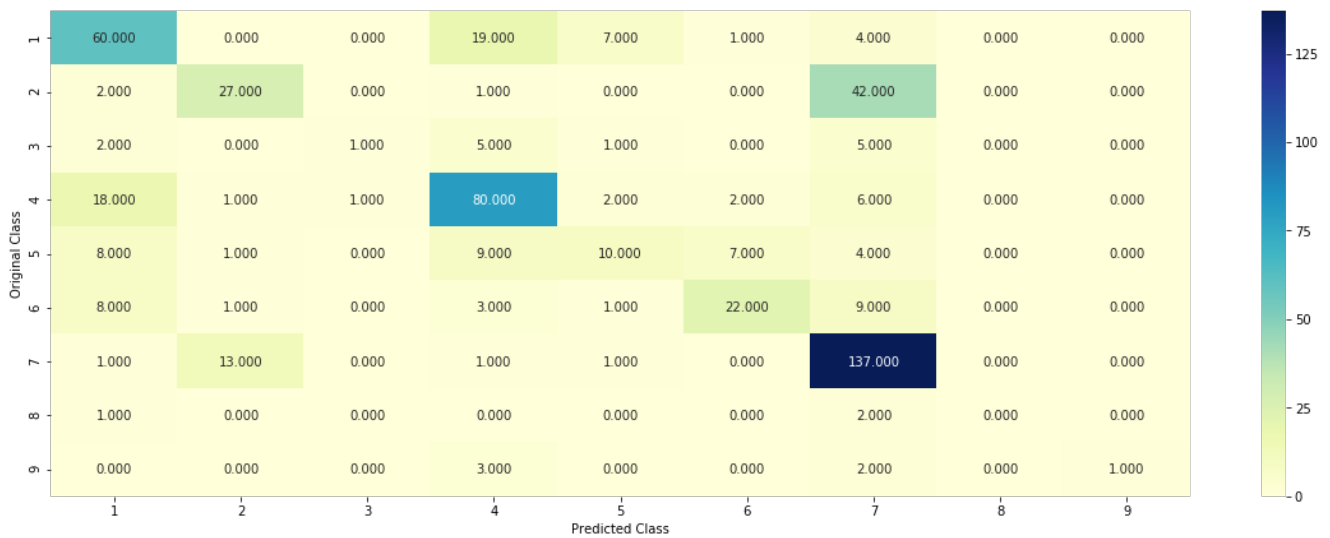
```
#-----
# video link:
#-----
```

```
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)
```

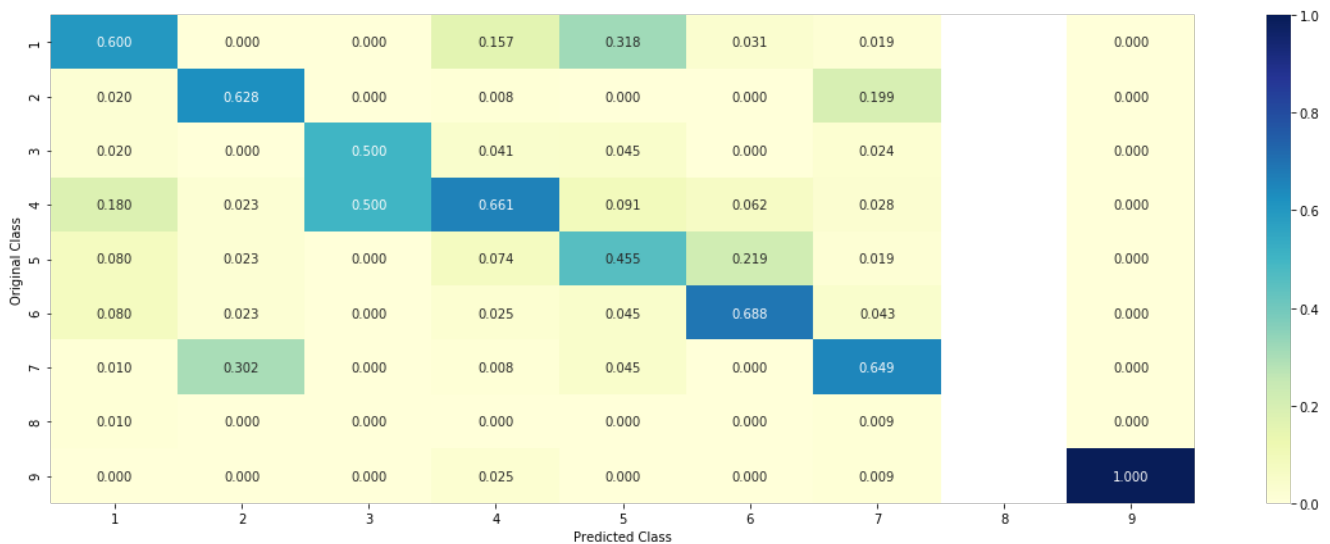
Log loss : 1.0809042565760936

Number of mis-classified points : 0.36466165413533835

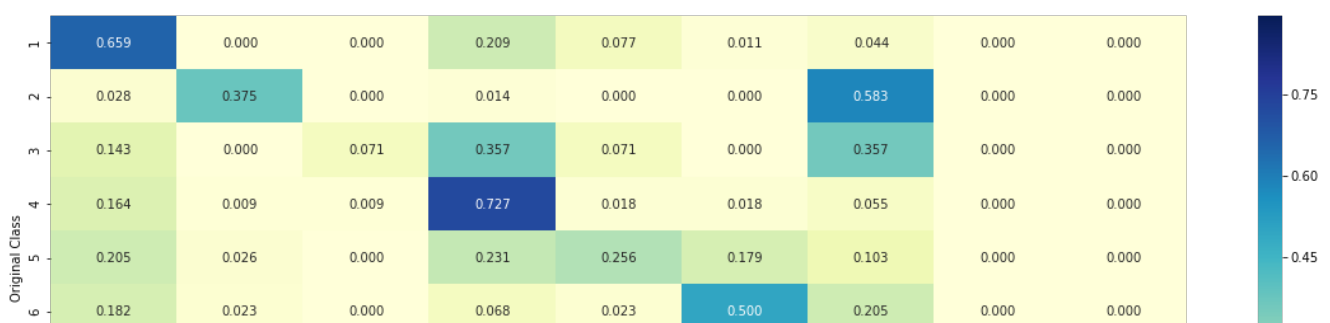
----- Confusion matrix -----

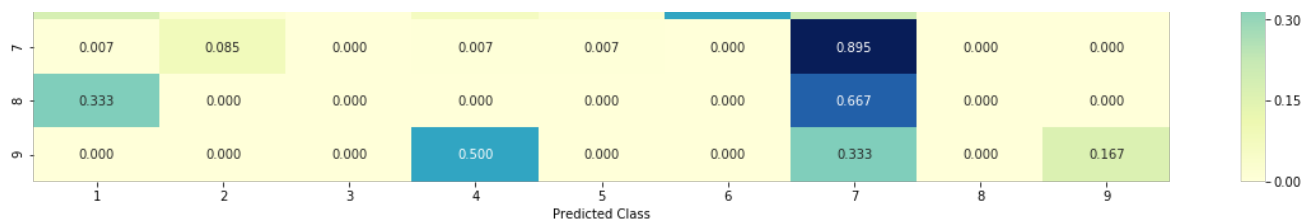


----- Precision matrix (Column Sum=1) -----



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





4.3.2.3. Feature Importance, Correctly Classified point

In [84]:

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

Predicted Class : 4

Predicted Class Probabilities: [[0.032 0.0241 0.0053 0.8883 0.0182 0.0124 0.017 0.0012 0.0016]]

Actual Class : 4

```
-----
0 Text feature [suppressor] present in test data point [True]
27 Text feature [missense] present in test data point [True]
62 Text feature [families] present in test data point [True]
83 Text feature [protein] present in test data point [True]
89 Text feature [mammalian] present in test data point [True]
93 Text feature [inactivation] present in test data point [True]
99 Text feature [yeast] present in test data point [True]
100 Text feature [proportion] present in test data point [True]
104 Text feature [see] present in test data point [True]
108 Text feature [mm] present in test data point [True]
127 Text feature [bind] present in test data point [True]
128 Text feature [null] present in test data point [True]
150 Text feature [homozygous] present in test data point [True]
155 Text feature [reduced] present in test data point [True]
156 Text feature [family] present in test data point [True]
157 Text feature [specifically] present in test data point [True]
162 Text feature [dna] present in test data point [True]
163 Text feature [consequences] present in test data point [True]
164 Text feature [functional] present in test data point [True]
169 Text feature [require] present in test data point [True]
182 Text feature [represent] present in test data point [True]
186 Text feature [show] present in test data point [True]
191 Text feature [2010] present in test data point [True]
192 Text feature [germline] present in test data point [True]
194 Text feature [changes] present in test data point [True]
195 Text feature [determine] present in test data point [True]
196 Text feature [loss] present in test data point [True]
200 Text feature [plates] present in test data point [True]
201 Text feature [kinases] present in test data point [True]
210 Text feature [38] present in test data point [True]
221 Text feature [cannot] present in test data point [True]
231 Text feature [localization] present in test data point [True]
233 Text feature [defective] present in test data point [True]
235 Text feature [high] present in test data point [True]
238 Text feature [cycle] present in test data point [True]
243 Text feature [appears] present in test data point [True]
246 Text feature [contribute] present in test data point [True]
247 Text feature [plasmid] present in test data point [True]
253 Text feature [bound] present in test data point [True]
256 Text feature [1998] present in test data point [True]
262 Text feature [stained] present in test data point [True]
271 Text feature [displayed] present in test data point [True]
```

273 Text feature [several] present in test data point [True]
274 Text feature [risk] present in test data point [True]
275 Text feature [stability] present in test data point [True]
277 Text feature [function] present in test data point [True]
287 Text feature [26] present in test data point [True]
291 Text feature [cellular] present in test data point [True]
292 Text feature [investigated] present in test data point [True]
298 Text feature [western] present in test data point [True]
299 Text feature [early] present in test data point [True]
303 Text feature [site] present in test data point [True]
306 Text feature [cases] present in test data point [True]
310 Text feature [37] present in test data point [True]
312 Text feature [splice] present in test data point [True]
313 Text feature [observed] present in test data point [True]
316 Text feature [described] present in test data point [True]
317 Text feature [particular] present in test data point [True]
318 Text feature [affected] present in test data point [True]
320 Text feature [mouse] present in test data point [True]
321 Text feature [45] present in test data point [True]
325 Text feature [amino] present in test data point [True]
326 Text feature [view] present in test data point [True]
330 Text feature [relative] present in test data point [True]
332 Text feature [red] present in test data point [True]
339 Text feature [associated] present in test data point [True]
341 Text feature [transfected] present in test data point [True]
343 Text feature [inhibitory] present in test data point [True]
347 Text feature [lower] present in test data point [True]
350 Text feature [terminal] present in test data point [True]
351 Text feature [figure] present in test data point [True]
361 Text feature [3b] present in test data point [True]
368 Text feature [antibodies] present in test data point [True]
369 Text feature [partial] present in test data point [True]
370 Text feature [2a] present in test data point [True]
371 Text feature [40] present in test data point [True]
372 Text feature [mutants] present in test data point [True]
373 Text feature [induced] present in test data point [True]
380 Text feature [low] present in test data point [True]
384 Text feature [average] present in test data point [True]
392 Text feature [transfection] present in test data point [True]
393 Text feature [similarly] present in test data point [True]
395 Text feature [phosphorylation] present in test data point [True]
397 Text feature [medium] present in test data point [True]
399 Text feature [genomic] present in test data point [True]
402 Text feature [comparison] present in test data point [True]
403 Text feature [groups] present in test data point [True]
405 Text feature [times] present in test data point [True]
406 Text feature [key] present in test data point [True]
407 Text feature [larger] present in test data point [True]
410 Text feature [pathogenic] present in test data point [True]
418 Text feature [indeed] present in test data point [True]
422 Text feature [sds] present in test data point [True]
423 Text feature [29] present in test data point [True]
424 Text feature [assay] present in test data point [True]
425 Text feature [containing] present in test data point [True]
426 Text feature [variant] present in test data point [True]
427 Text feature [first] present in test data point [True]
439 Text feature [involved] present in test data point [True]
441 Text feature [breast] present in test data point [True]
445 Text feature [whether] present in test data point [True]
446 Text feature [correlated] present in test data point [True]
447 Text feature [isolated] present in test data point [True]
448 Text feature [2000] present in test data point [True]
449 Text feature [genetic] present in test data point [True]
451 Text feature [combined] present in test data point [True]
452 Text feature [via] present in test data point [True]
456 Text feature [predicted] present in test data point [True]
459 Text feature [assessment] present in test data point [True]
461 Text feature [decrease] present in test data point [True]
464 Text feature [cdk4] present in test data point [True]
466 Text feature [1999] present in test data point [True]
471 Text feature [ca] present in test data point [True]
472 Text feature [altered] present in test data point [True]
473 Text feature [nuclear] present in test data point [True]
477 Text feature [previously] present in test data point [True]
480 Text feature [assays] present in test data point [True]
481 Text feature [locus] present in test data point [True]
482 Text feature [proteins] present in test data point [True]
...

```

487 Text feature [considered] present in test data point [True]
489 Text feature [70] present in test data point [True]
494 Text feature [linked] present in test data point [True]
496 Text feature [bp] present in test data point [True]
497 Text feature [various] present in test data point [True]
498 Text feature [critical] present in test data point [True]
499 Text feature [eight] present in test data point [True]
Out of the top 500 features 126 are present in query point

```

4.3.2.4. Feature Importance, Inorrectly Classified point

In [85]:

```

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

```

Predicted Class : 1

Predicted Class Probabilities: [[9.60e-01 2.38e-02 1.00e-04 1.35e-02 8.00e-04 2.00e-04 1.60e-03 0.00e+00
0.00e+00]]

Actual Class : 1

```

-----
12 Text feature [panel] present in test data point [True]
53 Text feature [nuclear] present in test data point [True]
60 Text feature [rather] present in test data point [True]
70 Text feature [repeats] present in test data point [True]
72 Text feature [ability] present in test data point [True]
73 Text feature [deletion] present in test data point [True]
88 Text feature [therefore] present in test data point [True]
97 Text feature [box] present in test data point [True]
112 Text feature [21] present in test data point [True]
121 Text feature [defined] present in test data point [True]
123 Text feature [effect] present in test data point [True]
130 Text feature [analyses] present in test data point [True]
138 Text feature [strong] present in test data point [True]
139 Text feature [essential] present in test data point [True]
141 Text feature [functions] present in test data point [True]
143 Text feature [corresponding] present in test data point [True]
153 Text feature [deficient] present in test data point [True]
155 Text feature [17] present in test data point [True]
158 Text feature [mutational] present in test data point [True]
162 Text feature [region] present in test data point [True]
163 Text feature [page] present in test data point [True]
173 Text feature [indicated] present in test data point [True]
176 Text feature [notably] present in test data point [True]
181 Text feature [calculated] present in test data point [True]
182 Text feature [position] present in test data point [True]
184 Text feature [whole] present in test data point [True]
189 Text feature [change] present in test data point [True]
193 Text feature [previous] present in test data point [True]
196 Text feature [another] present in test data point [True]
197 Text feature [function] present in test data point [True]
199 Text feature [gel] present in test data point [True]
200 Text feature [identify] present in test data point [True]
201 Text feature [impaired] present in test data point [True]
206 Text feature [encoding] present in test data point [True]
208 Text feature [value] present in test data point [True]
211 Text feature [subjected] present in test data point [True]
212 Text feature [blood] present in test data point [True]
213 Text feature [sequenced] present in test data point [True]
219 Text feature [supplemental] present in test data point [True]
220 Text feature [showing] present in test data point [True]
223 Text feature [induce] present in test data point [True]
229 Text feature [located] present in test data point [True]
230 Text feature [north] present in test data point [True]

```

230 Text feature [next] present in test data point [True]
232 Text feature [residues] present in test data point [True]
233 Text feature [least] present in test data point [True]
234 Text feature [set] present in test data point [True]
235 Text feature [signals] present in test data point [True]
236 Text feature [splicing] present in test data point [True]
238 Text feature [localization] present in test data point [True]
239 Text feature [sequences] present in test data point [True]
240 Text feature [possibility] present in test data point [True]
242 Text feature [values] present in test data point [True]
245 Text feature [splice] present in test data point [True]
246 Text feature [defects] present in test data point [True]
247 Text feature [interact] present in test data point [True]
250 Text feature [tested] present in test data point [True]
252 Text feature [via] present in test data point [True]
256 Text feature [selection] present in test data point [True]
257 Text feature [strand] present in test data point [True]
262 Text feature [screening] present in test data point [True]
264 Text feature [performed] present in test data point [True]
265 Text feature [insertion] present in test data point [True]
267 Text feature [percentage] present in test data point [True]
269 Text feature [obtained] present in test data point [True]
270 Text feature [upon] present in test data point [True]
272 Text feature [conserved] present in test data point [True]
273 Text feature [specific] present in test data point [True]
274 Text feature [medium] present in test data point [True]
278 Text feature [genome] present in test data point [True]
279 Text feature [signal] present in test data point [True]
281 Text feature [chain] present in test data point [True]
284 Text feature [sequencing] present in test data point [True]
286 Text feature [sample] present in test data point [True]
287 Text feature [length] present in test data point [True]
288 Text feature [16] present in test data point [True]
291 Text feature [within] present in test data point [True]
292 Text feature [wild] present in test data point [True]
295 Text feature [development] present in test data point [True]
296 Text feature [coding] present in test data point [True]
298 Text feature [treated] present in test data point [True]
300 Text feature [vitro] present in test data point [True]
301 Text feature [size] present in test data point [True]
304 Text feature [carrying] present in test data point [True]
305 Text feature [impact] present in test data point [True]
307 Text feature [taken] present in test data point [True]
309 Text feature [incubated] present in test data point [True]
311 Text feature [following] present in test data point [True]
313 Text feature [primers] present in test data point [True]
316 Text feature [del] present in test data point [True]
317 Text feature [possible] present in test data point [True]
318 Text feature [large] present in test data point [True]
319 Text feature [type] present in test data point [True]
322 Text feature [absence] present in test data point [True]
327 Text feature [fold] present in test data point [True]
328 Text feature [de] present in test data point [True]
329 Text feature [nucleotide] present in test data point [True]
331 Text feature [somatic] present in test data point [True]
336 Text feature [exon] present in test data point [True]
338 Text feature [dependent] present in test data point [True]
339 Text feature [rt] present in test data point [True]
340 Text feature [one] present in test data point [True]
343 Text feature [cell] present in test data point [True]
344 Text feature [remaining] present in test data point [True]
345 Text feature [none] present in test data point [True]
346 Text feature [table] present in test data point [True]
348 Text feature [mean] present in test data point [True]
350 Text feature [context] present in test data point [True]
353 Text feature [exhibited] present in test data point [True]
354 Text feature [additional] present in test data point [True]
357 Text feature [frame] present in test data point [True]
358 Text feature [among] present in test data point [True]
359 Text feature [33] present in test data point [True]
363 Text feature [rna] present in test data point [True]
365 Text feature [heterozygous] present in test data point [True]
369 Text feature [manufacturer] present in test data point [True]
372 Text feature [reduced] present in test data point [True]
374 Text feature [indeed] present in test data point [True]
375 Text feature [control] present in test data point [True]
377 Text feature [sporadic] present in test data point [True]
378 Text feature [evidence] present in test data point [True]

```

378 Text feature [evidence] present in test data point [True]
380 Text feature [analyzed] present in test data point [True]
381 Text feature [individual] present in test data point [True]
383 Text feature [05] present in test data point [True]
384 Text feature [reverse] present in test data point [True]
387 Text feature [cause] present in test data point [True]
389 Text feature [22] present in test data point [True]
390 Text feature [recently] present in test data point [True]
391 Text feature [established] present in test data point [True]
392 Text feature [stable] present in test data point [True]
394 Text feature [49] present in test data point [True]
396 Text feature [24] present in test data point [True]
401 Text feature [1b] present in test data point [True]
402 Text feature [allele] present in test data point [True]
403 Text feature [2001] present in test data point [True]
405 Text feature [example] present in test data point [True]
406 Text feature [pcr] present in test data point [True]
408 Text feature [could] present in test data point [True]
409 Text feature [lines] present in test data point [True]
410 Text feature [role] present in test data point [True]
412 Text feature [assay] present in test data point [True]
413 Text feature [non] present in test data point [True]
414 Text feature [1997] present in test data point [True]
415 Text feature [distinct] present in test data point [True]
417 Text feature [might] present in test data point [True]
418 Text feature [54] present in test data point [True]
420 Text feature [presence] present in test data point [True]
422 Text feature [32] present in test data point [True]
427 Text feature [deletions] present in test data point [True]
428 Text feature [http] present in test data point [True]
430 Text feature [sensitivity] present in test data point [True]
431 Text feature [population] present in test data point [True]
432 Text feature [complex] present in test data point [True]
434 Text feature [right] present in test data point [True]
435 Text feature [domains] present in test data point [True]
437 Text feature [43] present in test data point [True]
442 Text feature [key] present in test data point [True]
444 Text feature [derived] present in test data point [True]
445 Text feature [55] present in test data point [True]
446 Text feature [19] present in test data point [True]
448 Text feature [early] present in test data point [True]
449 Text feature [effects] present in test data point [True]
450 Text feature [results] present in test data point [True]
451 Text feature [limited] present in test data point [True]
452 Text feature [reported] present in test data point [True]
454 Text feature [larger] present in test data point [True]
455 Text feature [residue] present in test data point [True]
457 Text feature [predicted] present in test data point [True]
459 Text feature [protein] present in test data point [True]
463 Text feature [complete] present in test data point [True]
464 Text feature [generation] present in test data point [True]
467 Text feature [cohort] present in test data point [True]
470 Text feature [42] present in test data point [True]
473 Text feature [revealed] present in test data point [True]
474 Text feature [nm] present in test data point [True]
476 Text feature [furthermore] present in test data point [True]
479 Text feature [transcription] present in test data point [True]
480 Text feature [proteins] present in test data point [True]
484 Text feature [different] present in test data point [True]
486 Text feature [51] present in test data point [True]
487 Text feature [common] present in test data point [True]
488 Text feature [finally] present in test data point [True]
489 Text feature [based] present in test data point [True]
490 Text feature [exons] present in test data point [True]
492 Text feature [using] present in test data point [True]
493 Text feature [respectively] present in test data point [True]
494 Text feature [25] present in test data point [True]
Out of the top 500 features 186 are present in query point

```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

In [86]:

```
# 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/
# -----

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

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

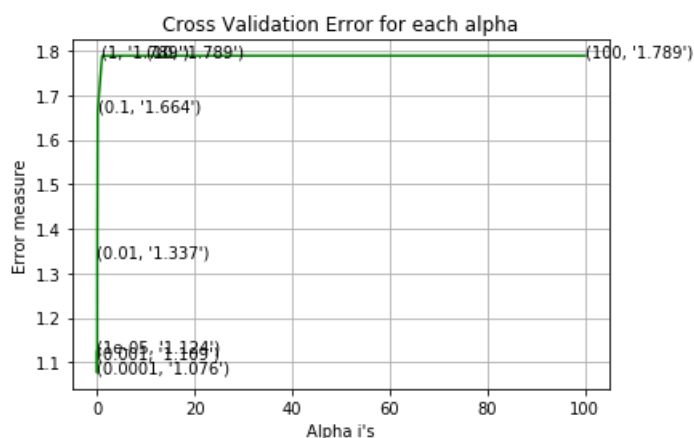
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], str(txt)), (alpha[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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

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

```
for C = 1e-05
Log Loss : 1.1239671615818072
for C = 0.0001
Log Loss : 1.0760324189972004
for C = 0.001
Log Loss : 1.108906137026918
for C = 0.01
Log Loss : 1.337427198724904
for C = 0.1
Log Loss : 1.6642874744580503
for C = 1
Log Loss : 1.7887823762757156
for C = 10
Log Loss : 1.7887823316450397
for C = 100
Log Loss : 1.788782267348181
```



For values of best alpha = 0.0001 The train log loss is: 0.6797907549629029
 For values of best alpha = 0.0001 The cross validation log loss is: 1.0760324189972004
 For values of best alpha = 0.0001 The test log loss is: 1.110967409077816

4.4.2. Testing model with best hyper parameters

In [87]:

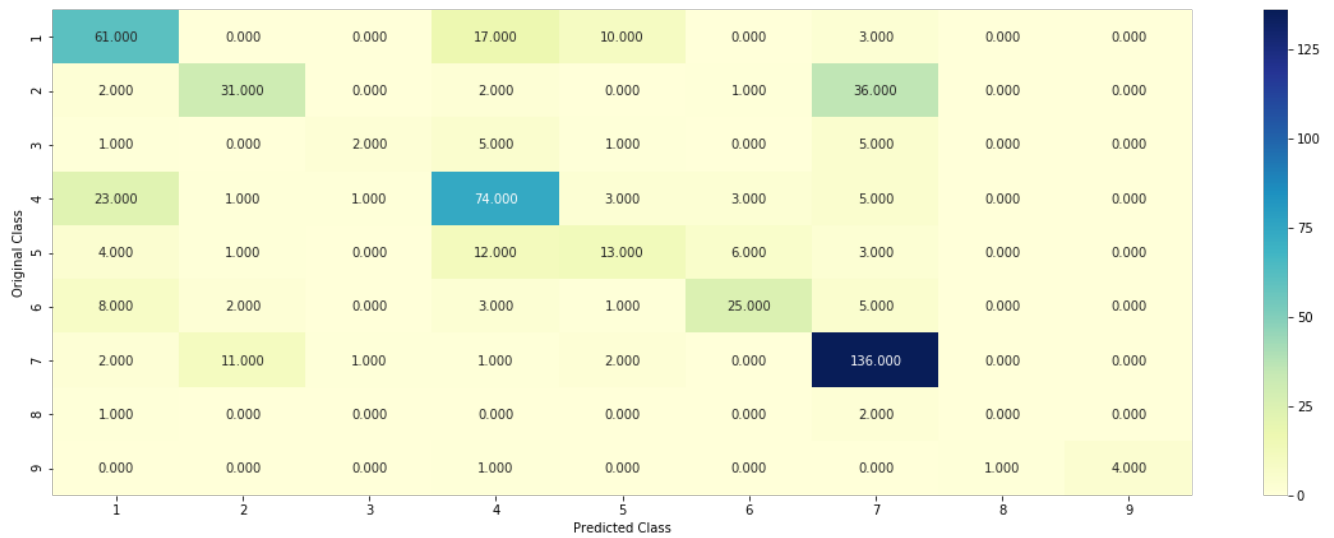
```
# read more about support vector machines with linear kernels here http://scikit-
learn.org/stable/modules/generated/sklearn.svm.SVC.html

# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, t
ol=0.001,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', ra
ndom_state=None)

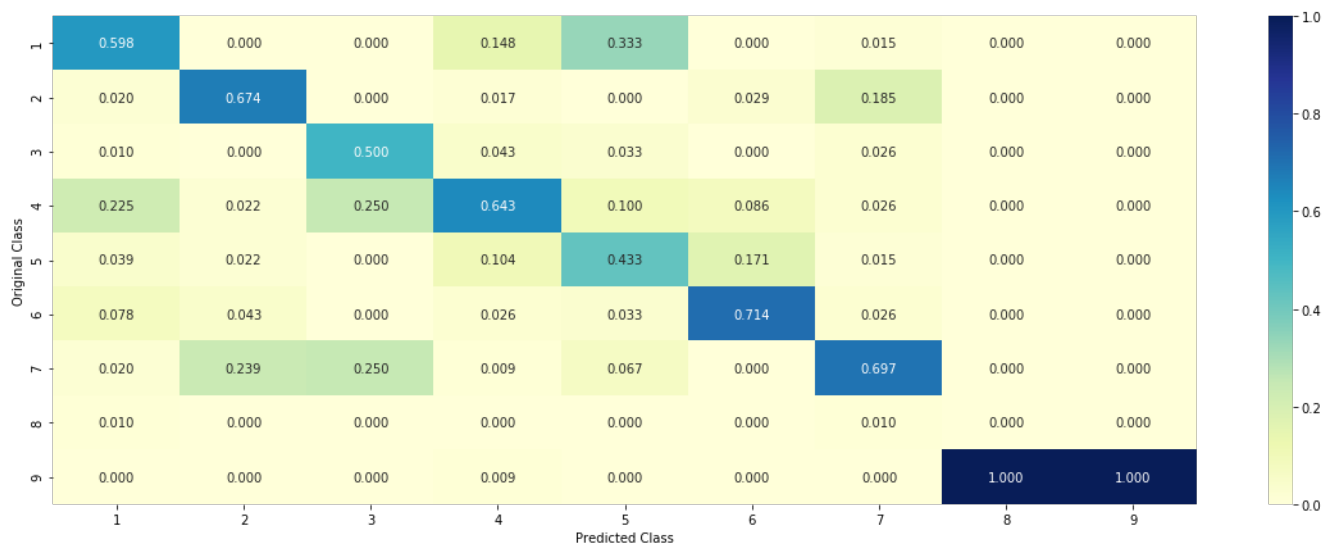
# Some of methods of SVM()
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
# predict(X) Perform classification on samples in X.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-
online/lessons/mathematical-derivation-copy-8/
# -----

# 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)
```

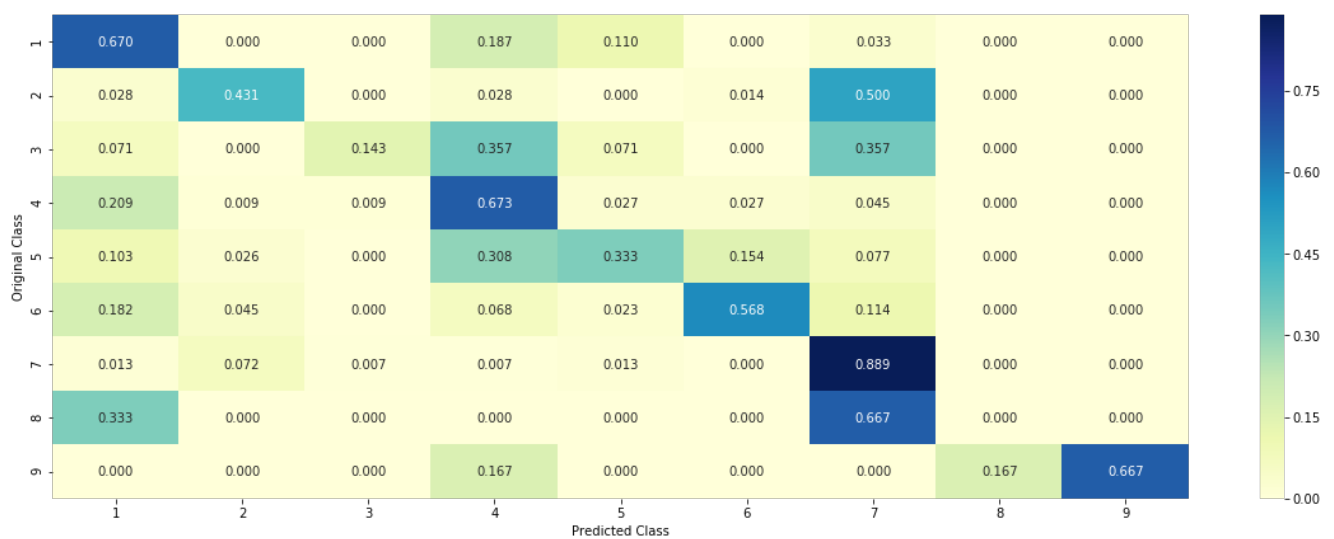
```
Log loss : 1.0760324189972004
Number of mis-classified points : 0.34962406015037595
----- Confusion matrix -----
```



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



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



4.3.3. Feature Importance

4.3.3.1. For Correctly classified point

In [88]:

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

Predicted Class : 4

Predicted Class Probabilities: [[0.04 0.0622 0.0122 0.7697 0.0288 0.0121 0.0702 0.0018 0.0029]]

Actual Class : 4

0 Text feature [suppressor] present in test data point [True]
59 Text feature [missense] present in test data point [True]
61 Text feature [families] present in test data point [True]
64 Text feature [proportion] present in test data point [True]
142 Text feature [mammalian] present in test data point [True]
143 Text feature [bind] present in test data point [True]
145 Text feature [mm] present in test data point [True]
147 Text feature [yeast] present in test data point [True]
150 Text feature [family] present in test data point [True]
151 Text feature [defective] present in test data point [True]
160 Text feature [specifically] present in test data point [True]
162 Text feature [1998] present in test data point [True]
163 Text feature [plates] present in test data point [True]
164 Text feature [see] present in test data point [True]
167 Text feature [require] present in test data point [True]
168 Text feature [inactivation] present in test data point [True]
169 Text feature [protein] present in test data point [True]
170 Text feature [bound] present in test data point [True]
171 Text feature [homozygous] present in test data point [True]
173 Text feature [dna] present in test data point [True]
174 Text feature [determine] present in test data point [True]
175 Text feature [38] present in test data point [True]
179 Text feature [high] present in test data point [True]
182 Text feature [consequences] present in test data point [True]
185 Text feature [changes] present in test data point [True]
187 Text feature [functional] present in test data point [True]
188 Text feature [show] present in test data point [True]
191 Text feature [site] present in test data point [True]
197 Text feature [2010] present in test data point [True]
198 Text feature [induced] present in test data point [True]
199 Text feature [medium] present in test data point [True]
200 Text feature [average] present in test data point [True]
201 Text feature [reduced] present in test data point [True]
202 Text feature [several] present in test data point [True]
204 Text feature [represent] present in test data point [True]
205 Text feature [plasmid] present in test data point [True]
208 Text feature [cannot] present in test data point [True]
209 Text feature [partial] present in test data point [True]
210 Text feature [loss] present in test data point [True]
211 Text feature [37] present in test data point [True]
214 Text feature [figure] present in test data point [True]
215 Text feature [null] present in test data point [True]
216 Text feature [mutants] present in test data point [True]
217 Text feature [phosphorylation] present in test data point [True]
219 Text feature [view] present in test data point [True]
221 Text feature [inhibitory] present in test data point [True]
222 Text feature [correlated] present in test data point [True]
224 Text feature [function] present in test data point [True]
225 Text feature [cases] present in test data point [True]
227 Text feature [29] present in test data point [True]
229 Text feature [described] present in test data point [True]
230 Text feature [similarly] present in test data point [True]
231 Text feature [observed] present in test data point [True]
235 Text feature [induced] present in test data point [True]

235 Text feature [kinases] present in test data point [True]
236 Text feature [line] present in test data point [True]
237 Text feature [localization] present in test data point [True]
358 Text feature [critical] present in test data point [True]
360 Text feature [sds] present in test data point [True]
361 Text feature [contribute] present in test data point [True]
362 Text feature [2000] present in test data point [True]
365 Text feature [germline] present in test data point [True]
367 Text feature [cellular] present in test data point [True]
368 Text feature [early] present in test data point [True]
369 Text feature [affected] present in test data point [True]
371 Text feature [comparison] present in test data point [True]
373 Text feature [displayed] present in test data point [True]
375 Text feature [mouse] present in test data point [True]
376 Text feature [amino] present in test data point [True]
378 Text feature [1999] present in test data point [True]
379 Text feature [stained] present in test data point [True]
381 Text feature [growth] present in test data point [True]
382 Text feature [first] present in test data point [True]
383 Text feature [times] present in test data point [True]
384 Text feature [investigated] present in test data point [True]
385 Text feature [groups] present in test data point [True]
386 Text feature [26] present in test data point [True]
387 Text feature [associated] present in test data point [True]
388 Text feature [absence] present in test data point [True]
389 Text feature [particular] present in test data point [True]
390 Text feature [affinity] present in test data point [True]
391 Text feature [isolated] present in test data point [True]
393 Text feature [primary] present in test data point [True]
396 Text feature [western] present in test data point [True]
397 Text feature [suggested] present in test data point [True]
398 Text feature [70] present in test data point [True]
399 Text feature [relative] present in test data point [True]
402 Text feature [splice] present in test data point [True]
403 Text feature [decrease] present in test data point [True]
405 Text feature [containing] present in test data point [True]
406 Text feature [considered] present in test data point [True]
408 Text feature [materials] present in test data point [True]
411 Text feature [low] present in test data point [True]
415 Text feature [45] present in test data point [True]
417 Text feature [transfected] present in test data point [True]
420 Text feature [effective] present in test data point [True]
421 Text feature [impact] present in test data point [True]
422 Text feature [23] present in test data point [True]
423 Text feature [combined] present in test data point [True]
425 Text feature [3b] present in test data point [True]
426 Text feature [following] present in test data point [True]
428 Text feature [reports] present in test data point [True]
431 Text feature [indeed] present in test data point [True]
433 Text feature [40] present in test data point [True]
434 Text feature [locus] present in test data point [True]
435 Text feature [significantly] present in test data point [True]
436 Text feature [various] present in test data point [True]
437 Text feature [cycle] present in test data point [True]
438 Text feature [2a] present in test data point [True]
439 Text feature [appears] present in test data point [True]
440 Text feature [via] present in test data point [True]
442 Text feature [variant] present in test data point [True]
444 Text feature [university] present in test data point [True]
447 Text feature [antibodies] present in test data point [True]
452 Text feature [risk] present in test data point [True]
455 Text feature [breast] present in test data point [True]
456 Text feature [ca] present in test data point [True]
457 Text feature [48] present in test data point [True]
460 Text feature [lines] present in test data point [True]
463 Text feature [http] present in test data point [True]
464 Text feature [larger] present in test data point [True]
466 Text feature [pathogenic] present in test data point [True]
467 Text feature [stability] present in test data point [True]
468 Text feature [likely] present in test data point [True]
469 Text feature [consistent] present in test data point [True]
470 Text feature [kinase] present in test data point [True]
477 Text feature [genomic] present in test data point [True]
478 Text feature [number] present in test data point [True]
480 Text feature [followed] present in test data point [True]
485 Text feature [large] present in test data point [True]
486 Text feature [gst] present in test data point [True]
487 Text feature [1999] present in test data point [True]

```

48/ Text feature [based] present in test data point [True]
491 Text feature [terminal] present in test data point [True]
494 Text feature [identified] present in test data point [True]
495 Text feature [buffer] present in test data point [True]
496 Text feature [key] present in test data point [True]
498 Text feature [history] present in test data point [True]
499 Text feature [genetic] present in test data point [True]
Out of the top 500 features 137 are present in query point

```

4.3.3.2. For Incorrectly classified point

In [89]:

```

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

```

```

Predicted Class : 1
Predicted Class Probabilities: [[9.003e-01 7.040e-02 3.000e-04 1.790e-02 4.300e-03 4.000e-04 6.000
e-03
2.000e-04 2.000e-04]]
Actual Class : 1

```

```

-----
10 Text feature [panel] present in test data point [True]
55 Text feature [21] present in test data point [True]
56 Text feature [ability] present in test data point [True]
57 Text feature [rather] present in test data point [True]
58 Text feature [analyses] present in test data point [True]
59 Text feature [repeats] present in test data point [True]
127 Text feature [box] present in test data point [True]
129 Text feature [deletion] present in test data point [True]
131 Text feature [nuclear] present in test data point [True]
132 Text feature [therefore] present in test data point [True]
134 Text feature [previous] present in test data point [True]
136 Text feature [selection] present in test data point [True]
138 Text feature [identify] present in test data point [True]
139 Text feature [supplemental] present in test data point [True]
140 Text feature [strong] present in test data point [True]
142 Text feature [page] present in test data point [True]
143 Text feature [whole] present in test data point [True]
144 Text feature [another] present in test data point [True]
147 Text feature [mutational] present in test data point [True]
149 Text feature [showing] present in test data point [True]
151 Text feature [defined] present in test data point [True]
153 Text feature [impaired] present in test data point [True]
154 Text feature [values] present in test data point [True]
156 Text feature [effect] present in test data point [True]
158 Text feature [17] present in test data point [True]
159 Text feature [essential] present in test data point [True]
161 Text feature [sequenced] present in test data point [True]
162 Text feature [calculated] present in test data point [True]
164 Text feature [gel] present in test data point [True]
166 Text feature [signals] present in test data point [True]
167 Text feature [corresponding] present in test data point [True]
168 Text feature [taken] present in test data point [True]
170 Text feature [indicated] present in test data point [True]
172 Text feature [position] present in test data point [True]
175 Text feature [screening] present in test data point [True]
176 Text feature [signal] present in test data point [True]
178 Text feature [blood] present in test data point [True]
181 Text feature [wild] present in test data point [True]
182 Text feature [functions] present in test data point [True]
183 Text feature [deficient] present in test data point [True]
184 Text feature [induce] present in test data point [True]
185 Text feature [located] present in test data point [True]
186 Text feature [value] present in test data point [True]

```

186 Text feature [value] present in test data point [True]
187 Text feature [subjected] present in test data point [True]
190 Text feature [fold] present in test data point [True]
191 Text feature [least] present in test data point [True]
192 Text feature [notably] present in test data point [True]
193 Text feature [residues] present in test data point [True]
194 Text feature [16] present in test data point [True]
196 Text feature [next] present in test data point [True]
197 Text feature [percentage] present in test data point [True]
199 Text feature [development] present in test data point [True]
201 Text feature [somatic] present in test data point [True]
204 Text feature [encoding] present in test data point [True]
206 Text feature [type] present in test data point [True]
208 Text feature [sequencing] present in test data point [True]
209 Text feature [medium] present in test data point [True]
210 Text feature [none] present in test data point [True]
211 Text feature [obtained] present in test data point [True]
212 Text feature [performed] present in test data point [True]
214 Text feature [localization] present in test data point [True]
215 Text feature [possibility] present in test data point [True]
216 Text feature [defects] present in test data point [True]
217 Text feature [sample] present in test data point [True]
296 Text feature [sporadic] present in test data point [True]
299 Text feature [carrying] present in test data point [True]
300 Text feature [incubated] present in test data point [True]
301 Text feature [sequences] present in test data point [True]
302 Text feature [conditions] present in test data point [True]
303 Text feature [heterozygous] present in test data point [True]
304 Text feature [22] present in test data point [True]
306 Text feature [change] present in test data point [True]
308 Text feature [exon] present in test data point [True]
309 Text feature [specific] present in test data point [True]
312 Text feature [upon] present in test data point [True]
313 Text feature [exhibited] present in test data point [True]
314 Text feature [genome] present in test data point [True]
315 Text feature [del] present in test data point [True]
316 Text feature [evidence] present in test data point [True]
317 Text feature [rt] present in test data point [True]
318 Text feature [mean] present in test data point [True]
319 Text feature [stable] present in test data point [True]
320 Text feature [splice] present in test data point [True]
324 Text feature [indeed] present in test data point [True]
325 Text feature [analyzed] present in test data point [True]
326 Text feature [region] present in test data point [True]
328 Text feature [05] present in test data point [True]
331 Text feature [55] present in test data point [True]
332 Text feature [table] present in test data point [True]
333 Text feature [chain] present in test data point [True]
334 Text feature [coding] present in test data point [True]
337 Text feature [set] present in test data point [True]
339 Text feature [cell] present in test data point [True]
340 Text feature [1997] present in test data point [True]
341 Text feature [tested] present in test data point [True]
342 Text feature [length] present in test data point [True]
345 Text feature [within] present in test data point [True]
346 Text feature [right] present in test data point [True]
347 Text feature [54] present in test data point [True]
348 Text feature [presence] present in test data point [True]
349 Text feature [population] present in test data point [True]
350 Text feature [absence] present in test data point [True]
351 Text feature [signaling] present in test data point [True]
352 Text feature [2001] present in test data point [True]
354 Text feature [reduced] present in test data point [True]
355 Text feature [reported] present in test data point [True]
356 Text feature [cause] present in test data point [True]
357 Text feature [sensitivity] present in test data point [True]
359 Text feature [allele] present in test data point [True]
362 Text feature [limited] present in test data point [True]
367 Text feature [primers] present in test data point [True]
368 Text feature [51] present in test data point [True]
369 Text feature [splicing] present in test data point [True]
370 Text feature [could] present in test data point [True]
371 Text feature [one] present in test data point [True]
373 Text feature [distinct] present in test data point [True]
374 Text feature [impact] present in test data point [True]
375 Text feature [insertion] present in test data point [True]
376 Text feature [treated] present in test data point [True]
377 Text feature [might] present in test data point [True]

```
377 Text feature [might] present in test data point [True]
378 Text feature [via] present in test data point [True]
380 Text feature [recently] present in test data point [True]
381 Text feature [key] present in test data point [True]
382 Text feature [function] present in test data point [True]
384 Text feature [context] present in test data point [True]
385 Text feature [cultured] present in test data point [True]
386 Text feature [insertions] present in test data point [True]
387 Text feature [possible] present in test data point [True]
389 Text feature [non] present in test data point [True]
390 Text feature [lines] present in test data point [True]
391 Text feature [pcr] present in test data point [True]
392 Text feature [interact] present in test data point [True]
394 Text feature [studied] present in test data point [True]
395 Text feature [exons] present in test data point [True]
396 Text feature [among] present in test data point [True]
397 Text feature [predicted] present in test data point [True]
399 Text feature [remaining] present in test data point [True]
401 Text feature [lb] present in test data point [True]
402 Text feature [nucleotide] present in test data point [True]
404 Text feature [conserved] present in test data point [True]
406 Text feature [even] present in test data point [True]
407 Text feature [additional] present in test data point [True]
408 Text feature [particular] present in test data point [True]
409 Text feature [43] present in test data point [True]
412 Text feature [patients] present in test data point [True]
413 Text feature [plasmid] present in test data point [True]
414 Text feature [de] present in test data point [True]
415 Text feature [common] present in test data point [True]
416 Text feature [related] present in test data point [True]
418 Text feature [strand] present in test data point [True]
419 Text feature [30] present in test data point [True]
420 Text feature [provided] present in test data point [True]
421 Text feature [different] present in test data point [True]
423 Text feature [33] present in test data point [True]
424 Text feature [52] present in test data point [True]
425 Text feature [reverse] present in test data point [True]
427 Text feature [selected] present in test data point [True]
429 Text feature [gene] present in test data point [True]
430 Text feature [combination] present in test data point [True]
431 Text feature [following] present in test data point [True]
433 Text feature [cohort] present in test data point [True]
434 Text feature [rna] present in test data point [True]
436 Text feature [min] present in test data point [True]
437 Text feature [finally] present in test data point [True]
438 Text feature [24] present in test data point [True]
439 Text feature [frame] present in test data point [True]
440 Text feature [mutagenesis] present in test data point [True]
441 Text feature [inhibitors] present in test data point [True]
443 Text feature [part] present in test data point [True]
444 Text feature [control] present in test data point [True]
445 Text feature [method] present in test data point [True]
447 Text feature [frequently] present in test data point [True]
448 Text feature [independent] present in test data point [True]
449 Text feature [phenotype] present in test data point [True]
450 Text feature [vitro] present in test data point [True]
451 Text feature [large] present in test data point [True]
454 Text feature [mutant] present in test data point [True]
455 Text feature [100] present in test data point [True]
457 Text feature [effects] present in test data point [True]
458 Text feature [42] present in test data point [True]
459 Text feature [assay] present in test data point [True]
460 Text feature [current] present in test data point [True]
462 Text feature [involving] present in test data point [True]
468 Text feature [diagnosis] present in test data point [True]
471 Text feature [association] present in test data point [True]
472 Text feature [2005] present in test data point [True]
473 Text feature [using] present in test data point [True]
474 Text feature [homozygous] present in test data point [True]
476 Text feature [study] present in test data point [True]
478 Text feature [first] present in test data point [True]
479 Text feature [manufacturer] present in test data point [True]
481 Text feature [established] present in test data point [True]
482 Text feature [ng] present in test data point [True]
483 Text feature [role] present in test data point [True]
484 Text feature [provide] present in test data point [True]
486 Text feature [example] present in test data point [True]
487 Text feature [moderately] present in test data point [True]
```



```

487 Text feature [majority] present in test data point [True]
488 Text feature [dependent] present in test data point [True]
489 Text feature [13] present in test data point [True]
491 Text feature [49] present in test data point [True]
492 Text feature [compared] present in test data point [True]
493 Text feature [codon] present in test data point [True]
494 Text feature [blot] present in test data point [True]
495 Text feature [41] present in test data point [True]
497 Text feature [1a] present in test data point [True]
498 Text feature [2004] present in test data point [True]
Out of the top 500 features 206 are present in query point

```

4.5 Random Forest Classifier

4.5.1. Hyper paramter tuning (With One hot Encoding)

In [90]:

```

# -----
# 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 = [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()

```

```

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 : 1.2104650476104701
 for n_estimators = 100 and max depth = 10
 Log Loss : 1.229839981617805
 for n_estimators = 200 and max depth = 5
 Log Loss : 1.194056714239944
 for n_estimators = 200 and max depth = 10
 Log Loss : 1.2169872464100477
 for n_estimators = 500 and max depth = 5
 Log Loss : 1.181379091988994
 for n_estimators = 500 and max depth = 10
 Log Loss : 1.2134859279603696
 for n_estimators = 1000 and max depth = 5
 Log Loss : 1.182933407723163
 for n_estimators = 1000 and max depth = 10
 Log Loss : 1.2089507869204237
 for n_estimators = 2000 and max depth = 5
 Log Loss : 1.1830596828378452
 for n_estimators = 2000 and max depth = 10
 Log Loss : 1.2095645950553566
 For values of best estimator = 500 The train log loss is: 0.8409143472123352
 For values of best estimator = 500 The cross validation log loss is: 1.181379091988994
 For values of best estimator = 500 The test log loss is: 1.2216369378538312

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

In [91]:

```

# -----
# 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]

```

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

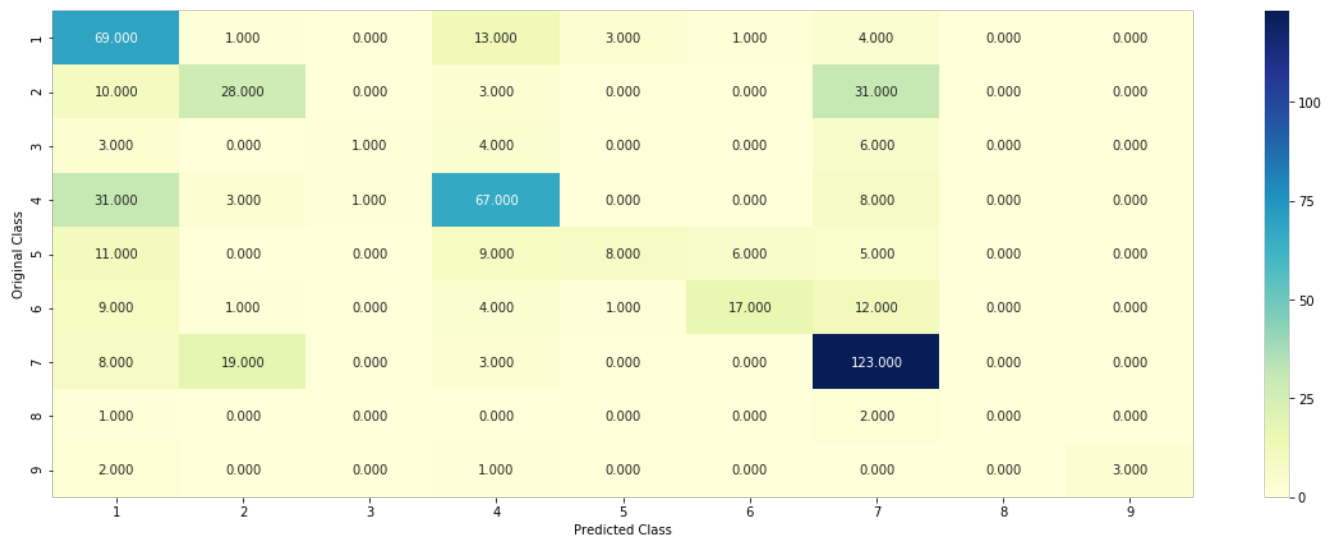
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

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

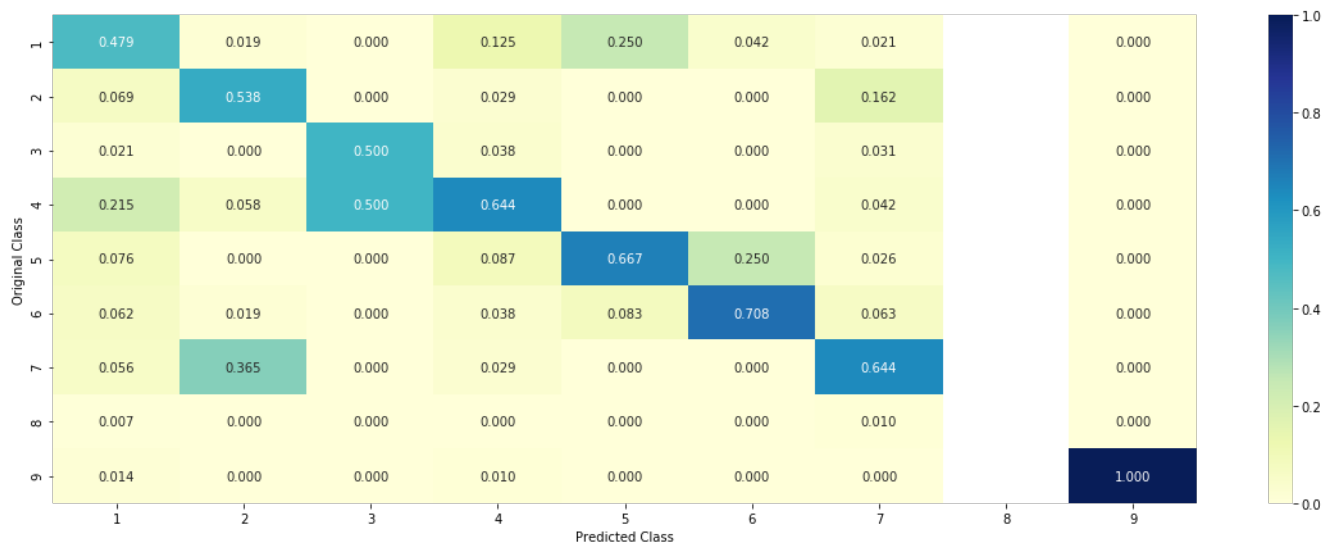
Log loss : 1.181379091988994

Number of mis-classified points : 0.40601503759398494

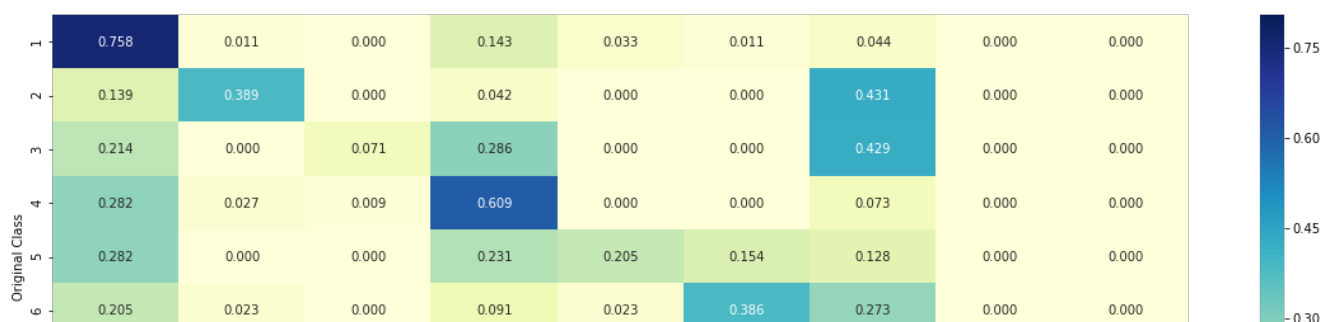
----- Confusion matrix -----

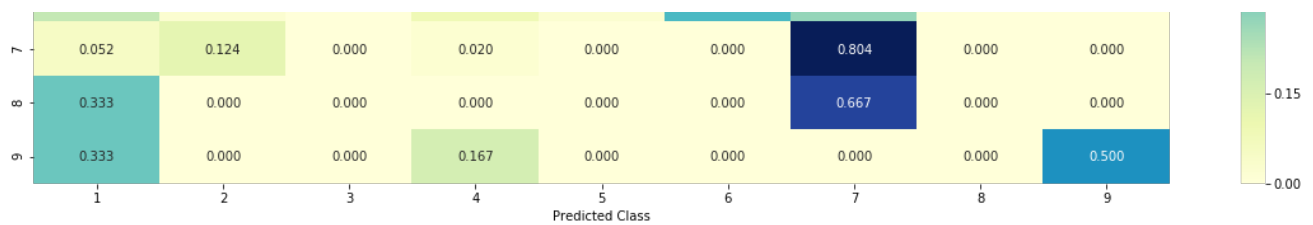


----- Precision matrix (Column Sum=1) -----



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





4.5.3. Feature Importance

4.5.3.1. Correctly Classified point

In [92]:

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

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

Predicted Class : 4

Predicted Class Probabilities: [[0.1574 0.0133 0.0113 0.6882 0.0552 0.0456 0.0225 0.0032 0.0033]]

Actual Class : 4

```
-----
0 Text feature [kinase] present in test data point [True]
2 Text feature [suppressor] present in test data point [True]
4 Text feature [inhibitors] present in test data point [True]
5 Text feature [function] present in test data point [True]
6 Text feature [missense] present in test data point [True]
9 Text feature [phosphorylation] present in test data point [True]
10 Text feature [loss] present in test data point [True]
12 Text feature [deleterious] present in test data point [True]
13 Text feature [protein] present in test data point [True]
15 Text feature [stability] present in test data point [True]
19 Text feature [pathogenic] present in test data point [True]
20 Text feature [neutral] present in test data point [True]
25 Text feature [functional] present in test data point [True]
27 Text feature [cell] present in test data point [True]
33 Text feature [classified] present in test data point [True]
34 Text feature [proteins] present in test data point [True]
36 Text feature [kinases] present in test data point [True]
40 Text feature [downstream] present in test data point [True]
41 Text feature [cells] present in test data point [True]
46 Text feature [expected] present in test data point [True]
48 Text feature [growth] present in test data point [True]
52 Text feature [predicted] present in test data point [True]
55 Text feature [functions] present in test data point [True]
56 Text feature [defective] present in test data point [True]
61 Text feature [variants] present in test data point [True]
63 Text feature [splice] present in test data point [True]
65 Text feature [months] present in test data point [True]
66 Text feature [57] present in test data point [True]
69 Text feature [proliferation] present in test data point [True]
70 Text feature [expression] present in test data point [True]
71 Text feature [yeast] present in test data point [True]
74 Text feature [nuclear] present in test data point [True]
75 Text feature [inactivation] present in test data point [True]
76 Text feature [phosphorylated] present in test data point [True]
80 Text feature [inhibition] present in test data point [True]
```

```

81 Text feature [potential] present in test data point [True]
82 Text feature [breast] present in test data point [True]
84 Text feature [use] present in test data point [True]
85 Text feature [expressing] present in test data point [True]
86 Text feature [sensitivity] present in test data point [True]
88 Text feature [variant] present in test data point [True]
89 Text feature [null] present in test data point [True]
92 Text feature [ability] present in test data point [True]
94 Text feature [clinical] present in test data point [True]
98 Text feature [type] present in test data point [True]
99 Text feature [21] present in test data point [True]
Out of the top 100 features 46 are present in query point

```

4.5.3.2. Incorrectly Classified point

In [93]:

```

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

```

```

Predicted Class : 1
Predicted Class Probabilities: [[0.6281 0.0386 0.0059 0.1477 0.0504 0.0502 0.0324 0.017 0.0298]]
Actual Class : 1

```

```

-----
4 Text feature [inhibitors] present in test data point [True]
5 Text feature [function] present in test data point [True]
6 Text feature [missense] present in test data point [True]
7 Text feature [activation] present in test data point [True]
9 Text feature [phosphorylation] present in test data point [True]
13 Text feature [protein] present in test data point [True]
15 Text feature [stability] present in test data point [True]
17 Text feature [treatment] present in test data point [True]
18 Text feature [brca1] present in test data point [True]
19 Text feature [pathogenic] present in test data point [True]
25 Text feature [functional] present in test data point [True]
26 Text feature [therapeutic] present in test data point [True]
27 Text feature [cell] present in test data point [True]
28 Text feature [activate] present in test data point [True]
30 Text feature [signaling] present in test data point [True]
31 Text feature [therapy] present in test data point [True]
32 Text feature [brca2] present in test data point [True]
33 Text feature [classified] present in test data point [True]
34 Text feature [proteins] present in test data point [True]
35 Text feature [treated] present in test data point [True]
40 Text feature [downstream] present in test data point [True]
41 Text feature [cells] present in test data point [True]
46 Text feature [expected] present in test data point [True]
50 Text feature [survival] present in test data point [True]
52 Text feature [predicted] present in test data point [True]
53 Text feature [patients] present in test data point [True]
55 Text feature [functions] present in test data point [True]
60 Text feature [brca] present in test data point [True]
61 Text feature [variants] present in test data point [True]
63 Text feature [splice] present in test data point [True]
64 Text feature [response] present in test data point [True]
68 Text feature [repair] present in test data point [True]
70 Text feature [expression] present in test data point [True]
74 Text feature [nuclear] present in test data point [True]
81 Text feature [potential] present in test data point [True]
82 Text feature [breast] present in test data point [True]
84 Text feature [use] present in test data point [True]
85 Text feature [expressing] present in test data point [True]
86 Text feature [sensitivity] present in test data point [True]
87 Text feature [pathway] present in test data point [True]
88 Text feature [variant] present in test data point [True]

```

```

89 Text feature [null] present in test data point [True]
92 Text feature [ability] present in test data point [True]
94 Text feature [clinical] present in test data point [True]
98 Text feature [type] present in test data point [True]
99 Text feature [21] present in test data point [True]
Out of the top 100 features 46 are present in query point

```

4.5.3. Hyper paramter tuning (With Response Coding)

In [94]:

```

# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_
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 : 2.201308963715666
for n_estimators = 10 and max depth = 3
Log Loss : 1.715576402352545
for n_estimators = 10 and max depth = 5
Log Loss : 1.4611230896289988
for n_estimators = 10 and max depth = 10
Log Loss : 1.7468155715605362
for n_estimators = 50 and max depth = 2
Log Loss : 1.7190133197710857
for n_estimators = 50 and max depth = 3
Log Loss : 1.4837863479799054
for n_estimators = 50 and max depth = 5
Log Loss : 1.413886352151792
for n_estimators = 50 and max depth = 10
Log Loss : 1.5606902745576023
for n_estimators = 100 and max depth = 2
Log Loss : 1.557195317852617
for n_estimators = 100 and max depth = 3
Log Loss : 1.4832220897946087
for n_estimators = 100 and max depth = 5
Log Loss : 1.2537573027728035
for n_estimators = 100 and max depth = 10
Log Loss : 1.5128031426075377
for n_estimators = 200 and max depth = 2
Log Loss : 1.624233611986707
for n_estimators = 200 and max depth = 3
Log Loss : 1.5255797829657924
for n_estimators = 200 and max depth = 5
Log Loss : 1.2654943694462788
for n_estimators = 200 and max depth = 10
Log Loss : 1.5317255553062172
for n_estimators = 500 and max depth = 2
Log Loss : 1.6791549372006445
for n_estimators = 500 and max depth = 3
Log Loss : 1.565164538412316
for n_estimators = 500 and max depth = 5
Log Loss : 1.293726277866204
for n_estimators = 500 and max depth = 10
Log Loss : 1.5660587192387527
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6482182778866807
for n_estimators = 1000 and max depth = 3
Log Loss : 1.5569256071386133
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2778922308136347
for n_estimators = 1000 and max depth = 10
Log Loss : 1.5484922596041149
For values of best alpha = 100 The train log loss is: 0.054526937802532705
For values of best alpha = 100 The cross validation log loss is: 1.2537573027728035
For values of best alpha = 100 The test log loss is: 1.324091335741171

```

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

In [95]:

```
# -----
# 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).

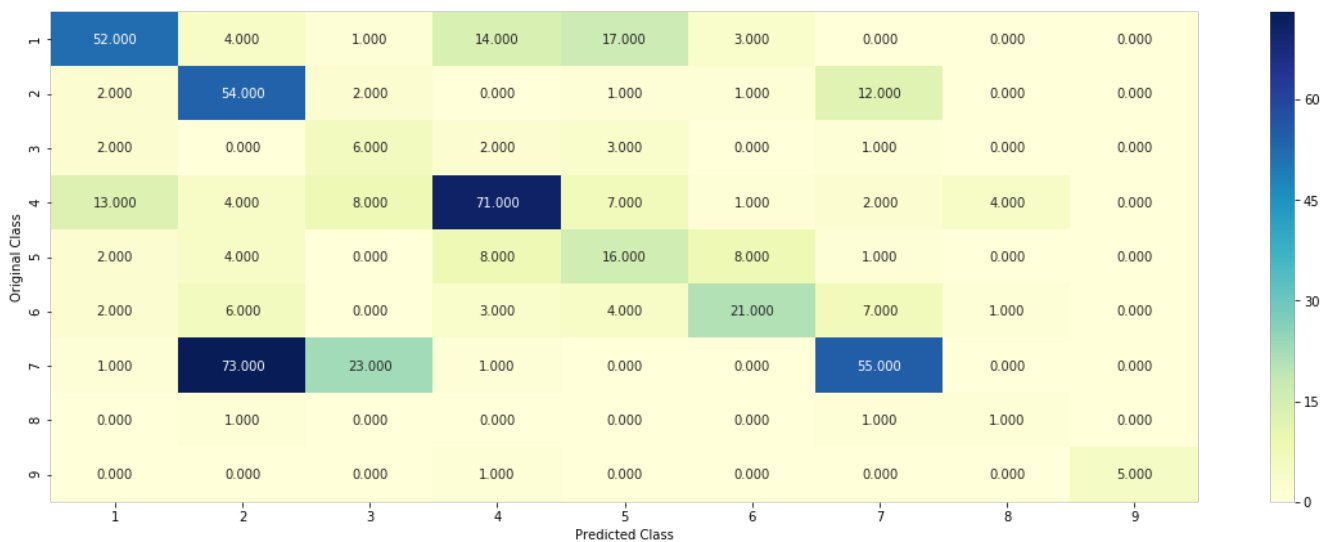
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

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)
```

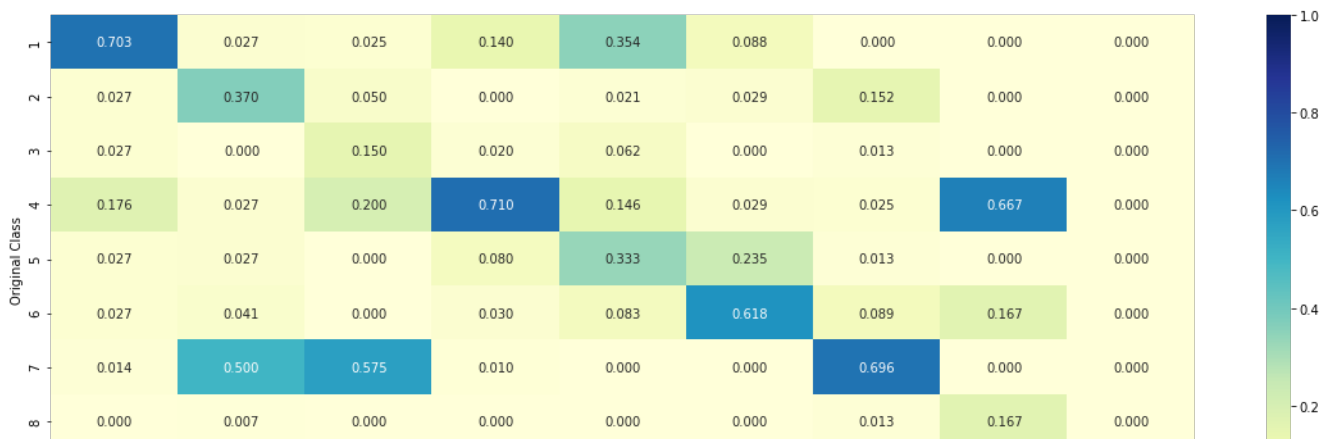
Log loss : 1.2537573027728035

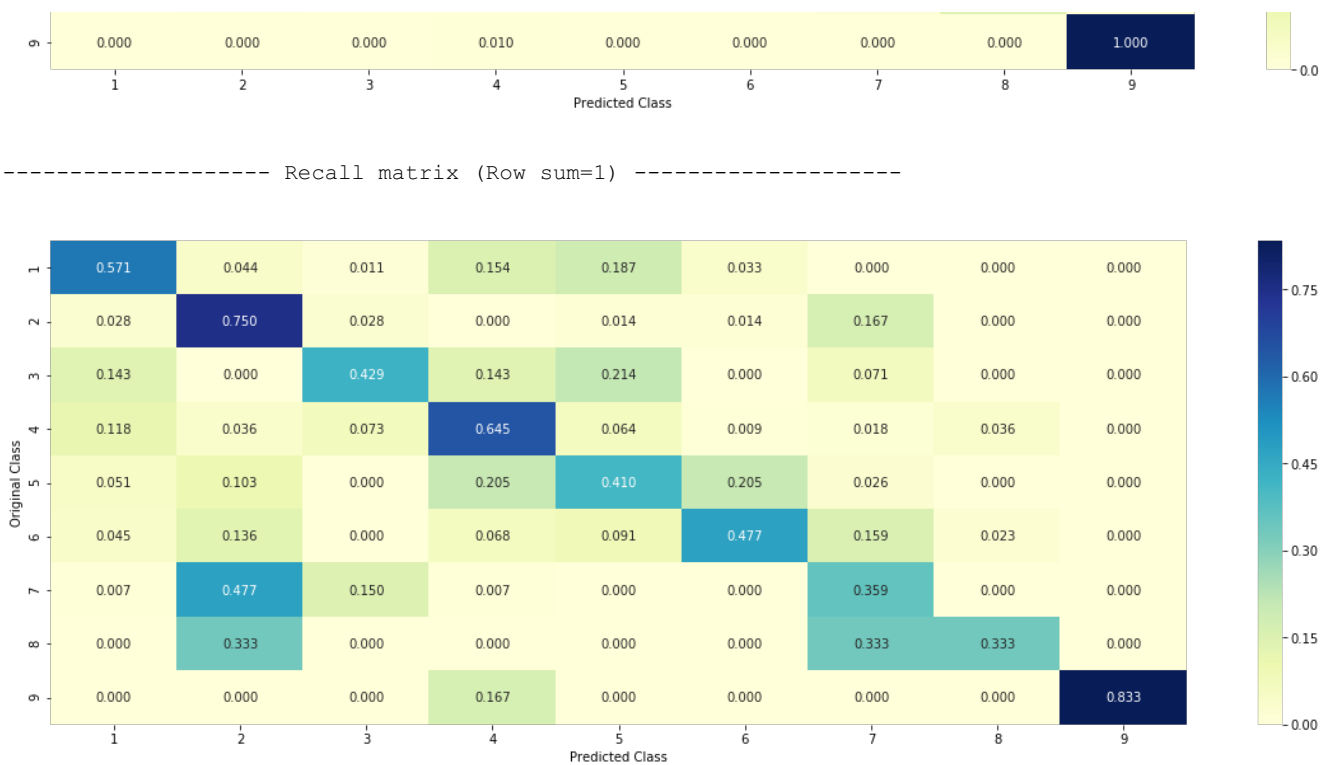
Number of mis-classified points : 0.4718045112781955

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



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





4.5.5. Feature Importance

4.5.5.1. Correctly Classified point

In [96]:

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

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.1094 0.0174 0.2285 0.5241 0.0206 0.0293 0.0074 0.0259 0.0374]]
Actual Class : 4
```

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

```

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

```

4.5.5.2. Incorrectly Classified point

In [97]:

```

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

```

```

Predicted Class : 1
Predicted Class Probabilities: [[0.9829 0.0014 0.0018 0.0046 0.001 0.0023 0.0015 0.0018 0.0027]]
Actual Class : 1

```

```

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

```

4.7 Stack the models

4.7.1 testing with hyper parameter tuning

In [98]:

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

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

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

# read more about support vector machines with linear kernels here http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html
# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_state=None)

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

# read more about support vector machines with linear kernels here http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html
# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
# class_weight=None)

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

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

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

clf1 = SGDClassifier(alpha=0.001, penalty='l2', loss='log', class_weight='balanced', random_state=0)
clf1.fit(train_x_onehotCoding, train_y)
sig_clf1 = CalibratedClassifierCV(clf1, method="sigmoid")

clf2 = SGDClassifier(alpha=1, penalty='l2', loss='hinge', class_weight='balanced', random_state=0)
```

```

clf2.fit(train_x_onehotCoding, train_y)
sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")

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

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

```

```

Logistic Regression : Log Loss: 1.04
Support vector machines : Log Loss: 1.79
Naive Bayes : Log Loss: 1.16
-----
Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.178
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.033
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.505
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.124
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.134
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.253

```

4.7.2 testing the model with the best hyper parameters

In [99]:

```

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))

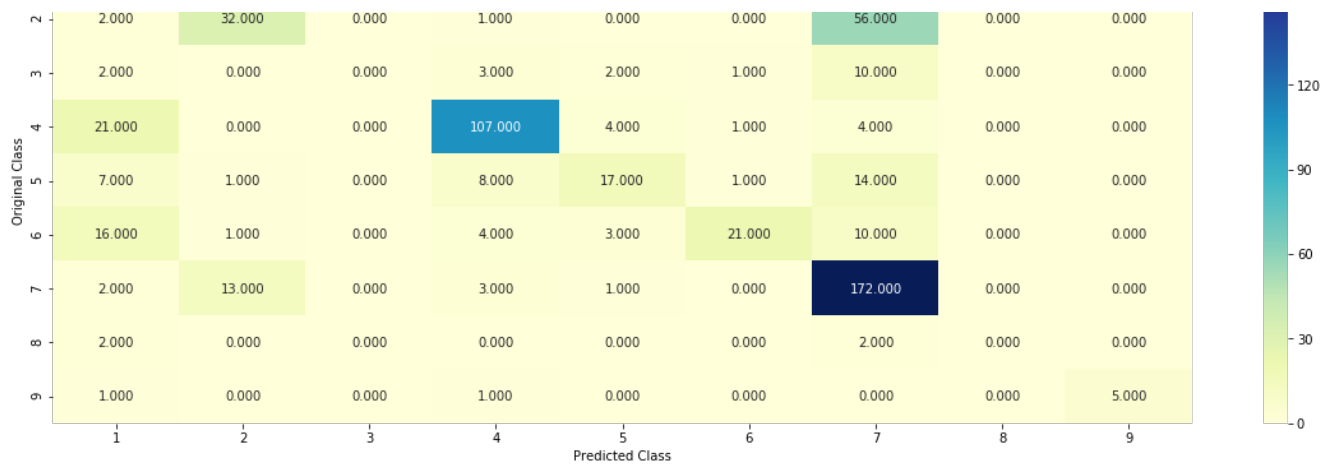
```

```

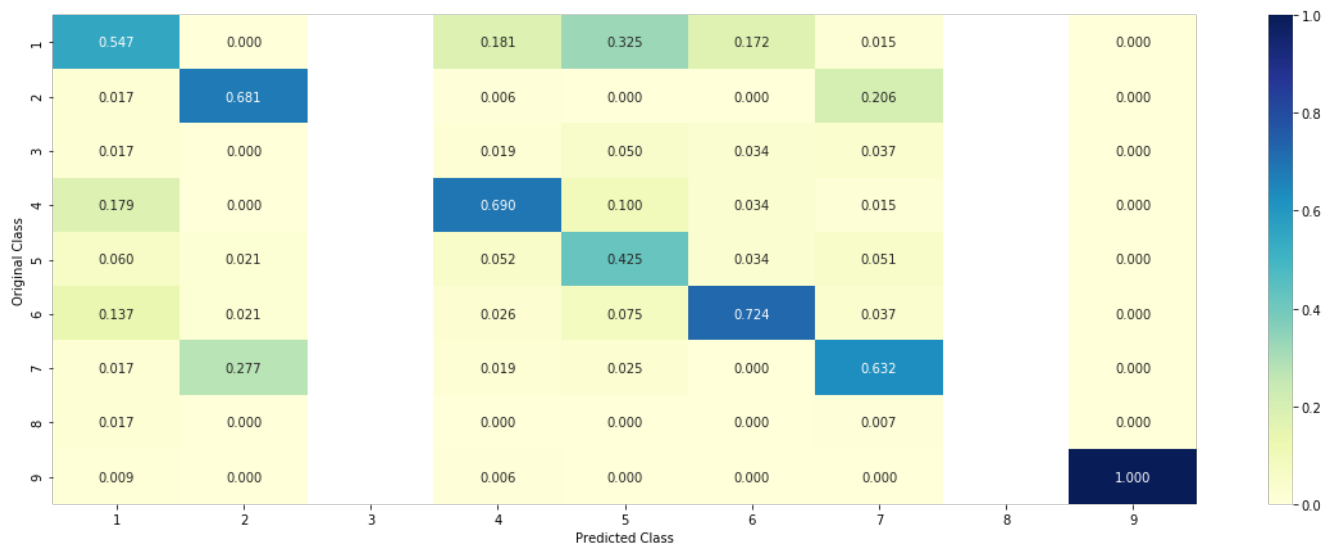
Log loss (train) on the stacking classifier : 0.7988623633997219
Log loss (CV) on the stacking classifier : 1.1236622548128592
Log loss (test) on the stacking classifier : 1.1440981600121238
Number of missclassified point : 0.37142857142857144
----- Confusion matrix -----

```

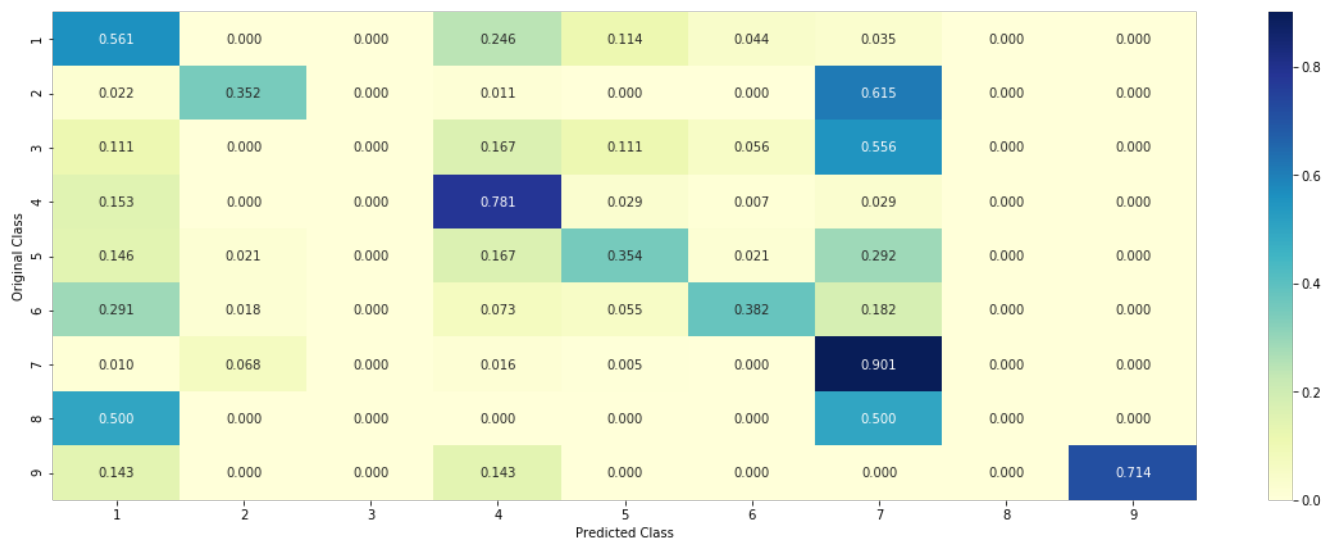
64.000	0.000	0.000	28.000	13.000	5.000	4.000	0.000	0.000
--------	-------	-------	--------	--------	-------	-------	-------	-------



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



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



4.7.3 Maximum Voting classifier

In [100]:

```
#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')
```

```

vcclf.fit(train_x_onehotCoding, train_y)
print("Log loss (train) on the VotingClassifier :", log_loss(train_y,
vcclf.predict_proba(train_x_onehotCoding)))
print("Log loss (CV) on the VotingClassifier :", log_loss(cv_y,
vcclf.predict_proba(cv_x_onehotCoding)))
print("Log loss (test) on the VotingClassifier :", log_loss(test_y,
vcclf.predict_proba(test_x_onehotCoding)))
print("Number of missclassified point :", np.count_nonzero((vcclf.predict(test_x_onehotCoding)-
test_y))/test_y.shape[0])
plot_confusion_matrix(test_y=test_y, predict_y=vcclf.predict(test_x_onehotCoding))

```

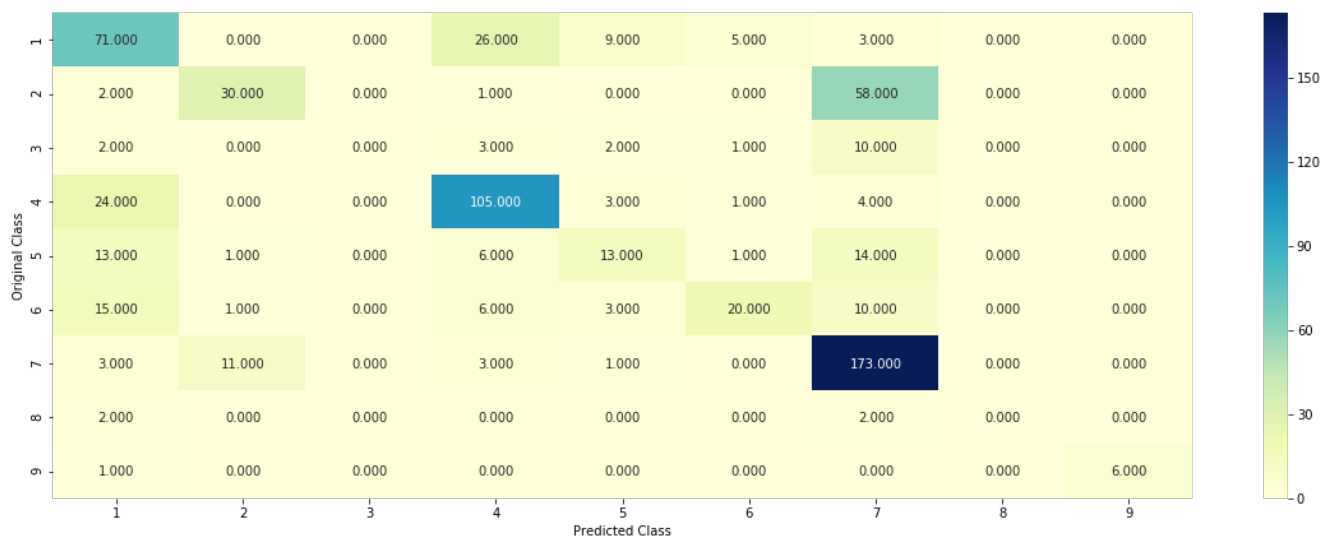
Log loss (train) on the VotingClassifier : 0.9414372681430161

Log loss (CV) on the VotingClassifier : 1.1671893322321885

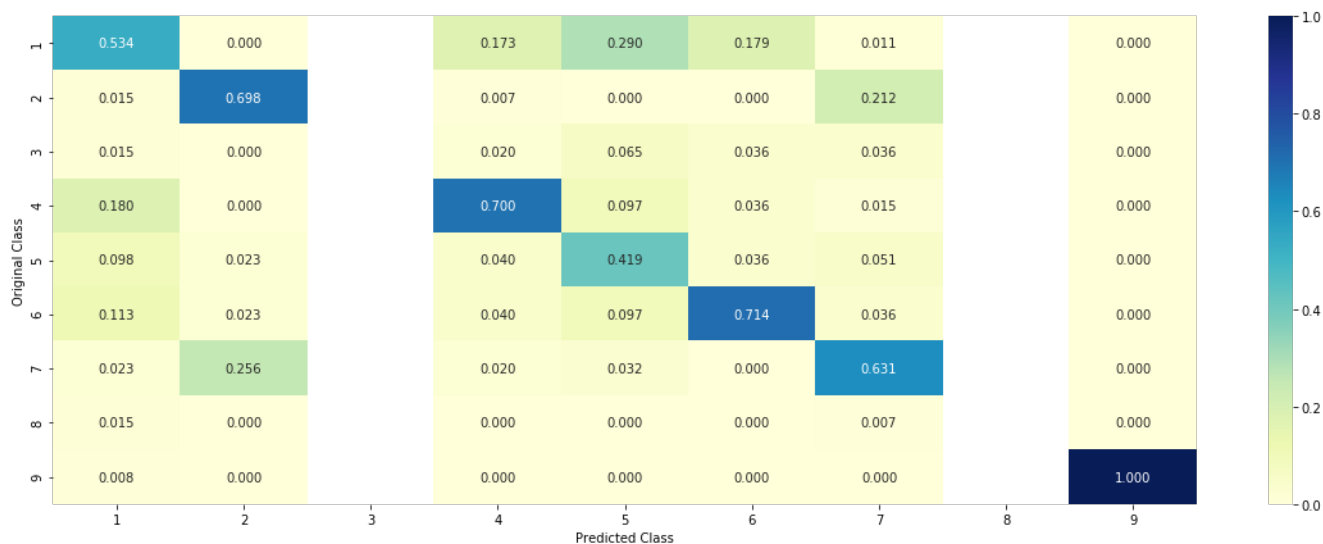
Log loss (test) on the VotingClassifier : 1.194325116744259

Number of missclassified point : 0.37142857142857144

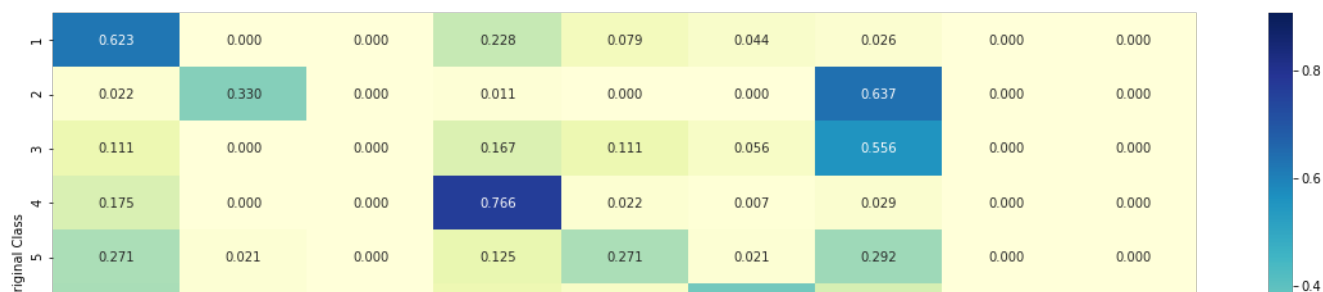
----- Confusion matrix -----

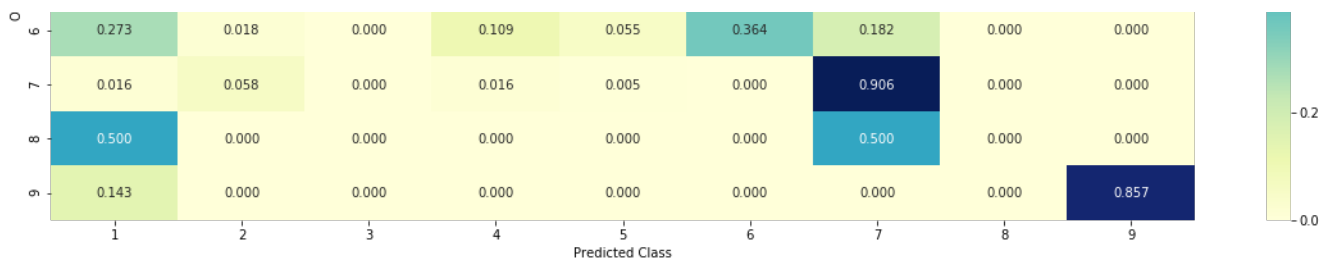


----- Precision matrix (Column Sum=1) -----



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





Conclusion

- According to our problem statement: We need to predict the probability of each data-point belonging to each of the nine classes. i.e Classify the given genetic variations/mutations based on evidence from text-based clinical literature
- Lets Start ->
 1. As we know we have dataset which contains ID, Gene, Variation, Test, Class and we have given class labels i.e 1-9 and genetic mutation has been classified on this.
 2. So let's start with Exploratory Data Analysis but before that we will do some Preprocessing of text so that we will be able to rectify that there is any row which has null value and after doing this we will split our whole data into Train, Cv and test and then we will do some EDA on these we will be able to visualize the distribution of data or the class in the train test and cv.
 3. And after that let's apply some random models to that we will be able to rectify the worst performance of the models, and then we will start working on the models
 4. Now after getting performance of the random models, let's do some featurization on our data set as we know we have dataset which contains ID, Gene, Variation, Test, Class from which we have Gene, variance and test as features. So let's start with univariate Analysis on each feature one by one and try to get the performance of each. Here we will work with two types of featurization i.e
 - Response coding
 - One hot coding

Note : We will choose the appropriate featurization based on the ML model we use. For this problem of multi-class classification with categorical features, one-hot encoding is better for Logistic regression while response coding is better for Random Forests. And at the same time we will try to get how good is this feature in predicting y_i , and is the feature stable across all the data sets (Test, Train, Cross validation), the distribution of the features and lot more.

5. After doing all above now we will start with our Machine learning models but before that we will combine our all featurized features into one i.e Stacking the three types of features and then will apply different machine learning models on it and try to get the best out of it by applying hyperparameter tuning on it each of the models.
6. The most imp this is to visualize the confusion matrix, precision and recall using heat map which help us to visualize the performance of our models so that we will be able to pick best out of these.

Note: In this we will Apply All the models with tf-idf features (Replace CountVectorizer with tfidfVectorizer) and Instead of using all the words in the dataset, use only the top 1000 words based on tf-idf values

Performance

In [19]:

```
from prettytable import PrettyTable
from termcolor import colored
print(colored('Performance Table', 'green'))
x = PrettyTable()
x.field_names = ["Models", "Train", "CV", "Test", "Misclassified(%)"]

x.add_row(["Naive Bayes (One hot coding)", 0.73, 1.16, 1.18, 0.377])
x.add_row(["KNN (Response)", 0.47, 1.008, 1.11, 0.3421])
x.add_row(["LR(Class balanced) one hot coding", 0.58, 1.02, 1.044, 0.3515])
x.add_row(["LR(Class unbalanced) one hot coding", 0.57, 1.08, 1.088, 0.3646])
x.add_row(["Lr SVM one hot encoding", 0.67, 1.07, 1.11, 0.34])
x.add_row(["Random Forest one hot coding", 0.84, 1.18, 1.22, 0.40])
x.add_row(["Random Forest Response coding", 0.05, 1.25, 1.32, 0.47])
x.add_row(["Stacking classifier", 0.79, 1.12, 1.14, 0.37])
x.add_row(["Maximum Voting Classifier", 0.94, 1.16, 1.19, 0.37])
```

```
print(x)
```

Performance Table

Models	Train	CV	Test	Misclassified(%)
Naive Bayes (One hot coding)	0.73	1.16	1.18	0.377
KNN (Response)	0.47	1.008	1.11	0.3421
LR(Class balanced) one hot coding	0.58	1.02	1.044	0.3515
LR(Class unbalanced) one hot coding	0.57	1.08	1.088	0.3646
Lr SVM one hot encoding	0.67	1.07	1.11	0.34
Random Forest one hot coding	0.84	1.18	1.22	0.4
Random Forest Response coding	0.05	1.25	1.32	0.47
Stacking classifier	0.79	1.12	1.14	0.37
Maximum Voting Classifier	0.94	1.16	1.19	0.37