

Personalized cancer diagnosis

1. Business Problem

1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

Context:

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

Problem statement :

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

1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

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

1.3. Real-world/Business objectives and constraints.

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

2. Machine Learning Problem Formulation

2.1. Data

2.1.1. Data Overview

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

2.1.2. Example Data Point

training_variants

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

training_text

ID, Text

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

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

2.2.1. Type of Machine Learning Problem

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

2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

2.3. Train, CV and Test Datasets

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

3. Exploratory Data Analysis

In [1]:

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

from mlxtend.classifier import StackingClassifier

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

Using TensorFlow backend.

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]:

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 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.process_time()
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.process_time() - 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 : 45.296875 seconds
```

In [6]:

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

Out[6]:

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

In [7]:

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

Out[7]:

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

In [8]:

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

In [9]:

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

Out[9]:

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

3.1.4. Test, Train and Cross Validation Split

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

In [10]:

```
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().sort_index()
test_class_distribution = test_df['Class'].value_counts().sort_index()
cv_class_distribution = cv_df['Class'].value_counts().sort_index()

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

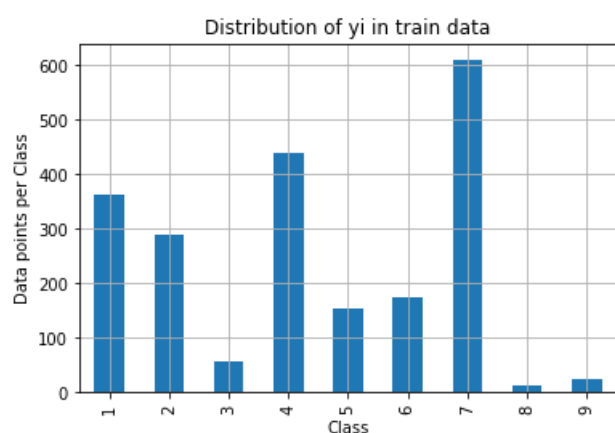
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

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

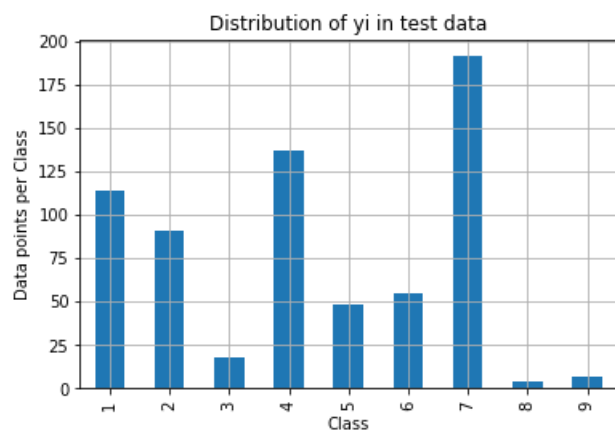
```
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1,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

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

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



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

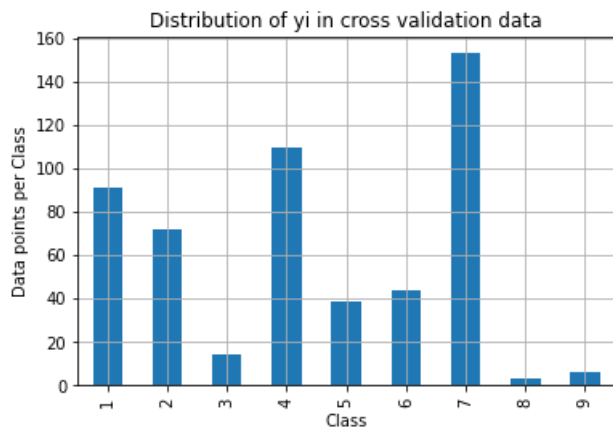


Number of data points in class 7 : 191 (28.722 %)
Number of data points in class 4 : 137 (20.602 %)
Number of data points in class 1 : 114 (17.143 %)
Number of data points in class 2 : 91 (13.684 %)
Number of data points in class 6 : 55 (8.271 %)
Number of data points in class 5 : 48 (7.219 %)

```

Number of data points in class 3 : 46 ( 7.216 %)
Number of data points in class 3 : 18 ( 2.707 %)
Number of data points in class 9 : 7 ( 1.053 %)
Number of data points in class 8 : 4 ( 0.602 %)
-----

```



```

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

```

3.2 Prediction using a 'Random' Model

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

In [13]:

```

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

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

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

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

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

    labels = [1,2,3,4,5,6,7,8,9]
    # representing A in heatmap format

```



```

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

print("-"*20, "Precision matrix (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))

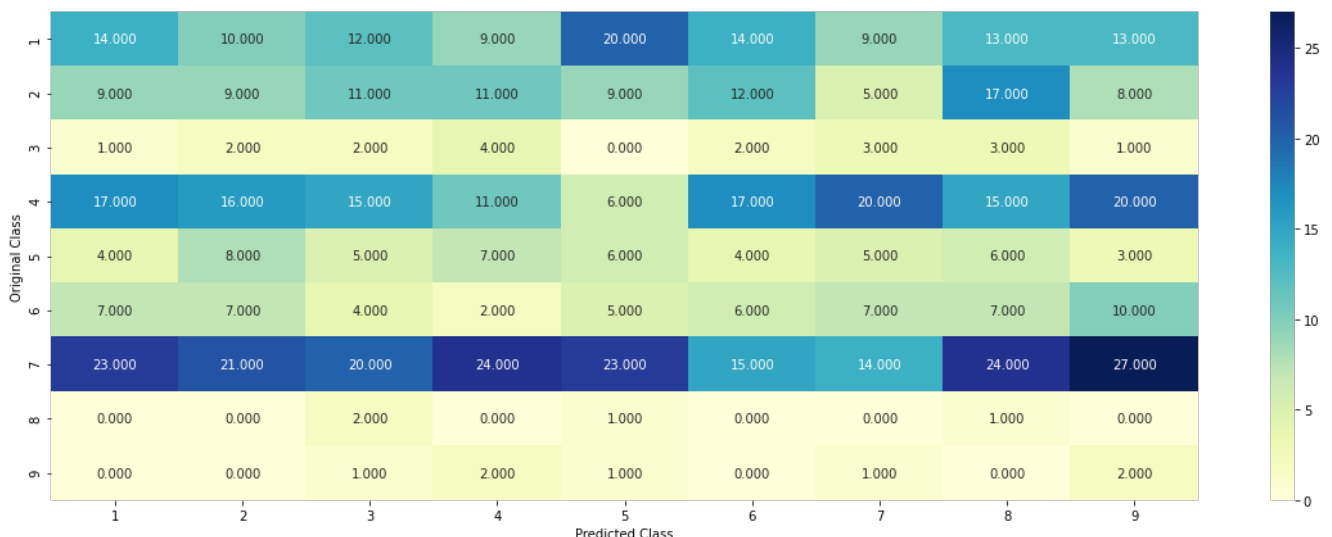
predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

```

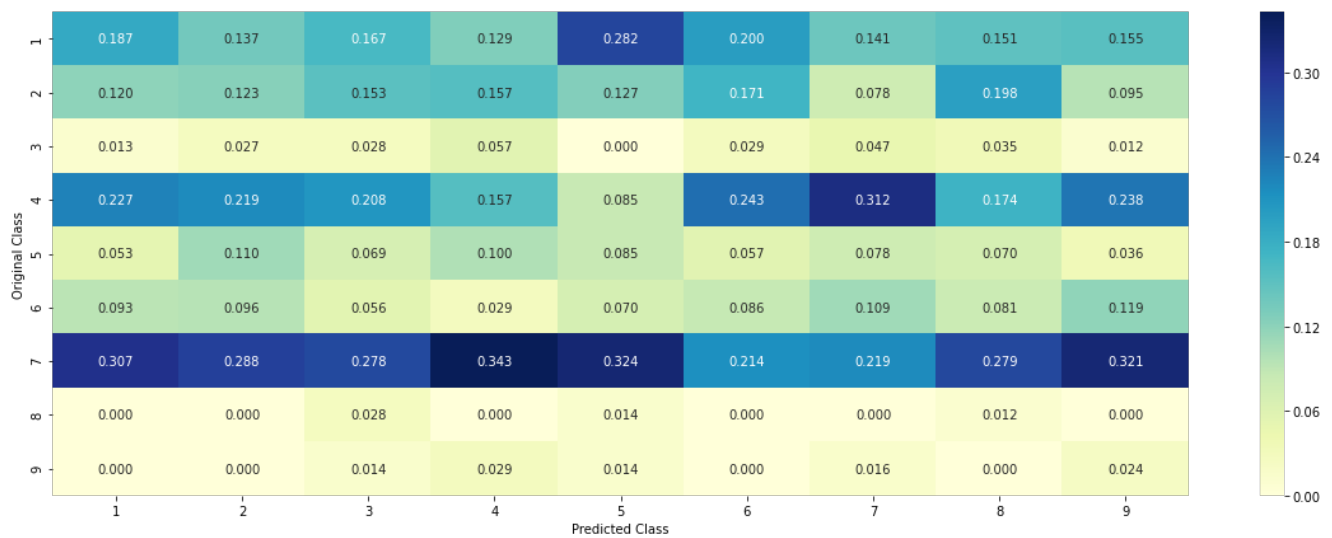
Log loss on Cross Validation Data using Random Model 2.44365319425409

Log loss on Test Data using Random Model 2.4772076655106283

----- 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 store a (1*9) representation of it
# for a value of feature in df:
# if it is in train data:
# we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
# if it is not there is train:
# we add [1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
# return 'gv_fea'
# -----

# get_gv_fea_dict: Get Gene variation Feature Dict
def get_gv_fea_dict(alpha, feature, df):
    # value_count: it contains a dict like
    # print(train_df[feature].value_counts())
```

```

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

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

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

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

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    #      {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.0681818181818177,
0.13636363636363635, 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': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.0681818181818177,
0.0681818181818177, 0.0625, 0.34659090909090912, 0.0625, 0.0568181818181816],
#      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608,
0.0787878787878782, 0.1393939393939394, 0.34545454545454546, 0.060606060606060608,
0.060606060606060608, 0.060606060606060608],
#      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917,
0.46540880503144655, 0.075471698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081
761006289, 0.062893081761006289],
#      'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295,
0.072847682119205295, 0.066225165562913912, 0.066225165562913912, 0.27152317880794702,
0.066225165562913912, 0.066225165562913912]}

```

```

0.066223163362913912, 0.066223163362913912],
# 'BRAF': [0.06666666666666666, 0.17999999999999999, 0.07333333333333334,
0.07333333333333334, 0.09333333333333338, 0.08000000000000002, 0.29999999999999999,
0.06666666666666666, 0.06666666666666666],
# ...
# }
gv_dict = get_gv_fea_dict(alpha, feature, df)
# value_count is similar in get_gv_fea_dict
value_count = train_df[feature].value_counts()

# gv_fea: Gene_variation feature, it will contain the feature for each feature value in the da
ta
gv_fea = []
# for every feature values in the given data frame we will check if it is there in the train
data then we will add the feature to gv_fea
# if not we will add [1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9] 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 : 235
BRCA1      169
TP53       103
EGFR        87
PTEN        80
BRCA2       75
KIT         62
BRAF        56
ERBB2       45
ALK         44
PIK3CA      41
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 235 different categories of genes in the train data, and they are distributed as follows

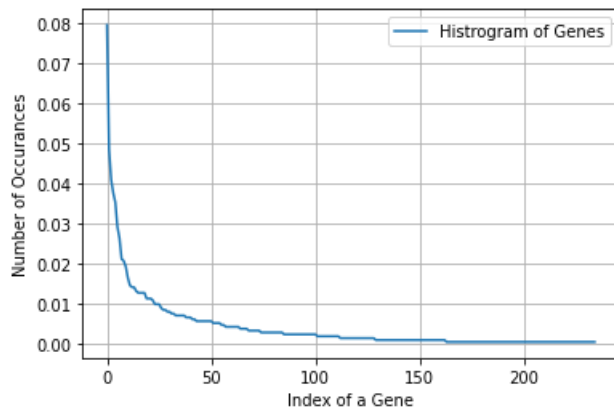
In [18]:

```

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')

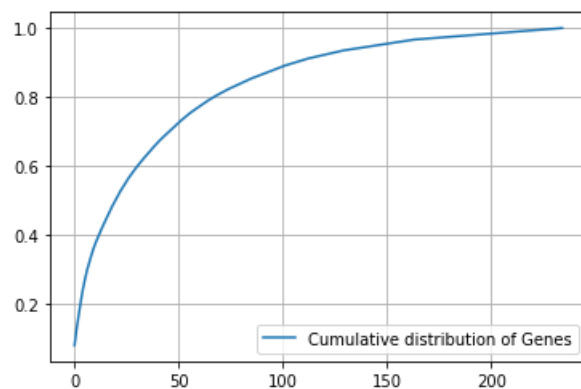
```

```
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 sha
```

```
pe 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 = CountVectorizer()
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

In [23]:

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

Out[23]:

```
2731      BRAF
2235      PTEN
1829  PPP2R1A
2229      PTEN
2101      CDK12
Name: Gene, dtype: object
```

In [24]:

```
gene_vectorizer.get_feature_names()
```

Out[24]:

```
['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ar',
 'araf',
 'arid1a',
 'arid1b',
 'arid2',
 'asx11',
 'asx12',
 'atm',
 'atrx',
 'aurka',
 'aurkb',
 'axin1',
 'axl',
 'b2m',
 'bap1',
 'bard1',
 'bcl10',
 'bcl2',
 'bcl2l11',
 'bcor',
 'braf',
 'brca1',
 'brca2',
 'brd4',
 'brip1',
 'btk',
 'card11',
 'carm1',
 'casp8',
 'cb1',
 'ccnd1',
 'ccnd3',
 'ccne1',
```

'cdh1',
'cdk12',
'cdk4',
'cdk6',
'cdkn1a',
'cdkn1b',
'cdkn2a',
'cdkn2b',
'cdkn2c',
'chek2',
'cic',
'crebbp',
'ctcf',
'ctnnb1',
'ddr2',
'dicer1',
'dnmt3a',
'dnmt3b',
'dusp4',
'egfr',
'eif1ax',
'elf3',
'ep300',
'epas1',
'epcam',
'erbb2',
'erbb3',
'erbb4',
'ercc2',
'ercc3',
'ercc4',
'erg',
'esr1',
'etv1',
'etv6',
'ewsr1',
'ezh2',
'fam58a',
'fanca',
'fancc',
'fat1',
'fbxw7',
'fgf19',
'fgf4',
'fgfr1',
'fgfr2',
'fgfr3',
'fgfr4',
'flt1',
'flt3',
'foxa1',
'foxl2',
'foxo1',
'foxp1',
'fubp1',
'gata3',
'gna11',
'gnaq',
'gnas',
'h3f3a',
'hist1h1c',
'hla',
'hnf1a',
'hras',
'idh1',
'idh2',
'igflr',
'ikzf1',
'il7r',
'inpp4b',
'jak1',
'jak2',
'jun',
'kdm5a',
'kdm5c',
'kdm6a',
'kdr',

'keap1',
'kit',
'klf4',
'kmt2a',
'kmt2c',
'kmt2d',
'knstrn',
'kras',
'lats1',
'lats2',
'map2k1',
'map2k2',
'map2k4',
'map3k1',
'mapk1',
'mdm4',
'med12',
'mef2b',
'men1',
'met',
'mga',
'mlh1',
'mpl',
'msh2',
'msh6',
'mtor',
'myc',
'myd88',
'myod1',
'nfl1',
'nf2',
'nfe2l2',
'nfkbia',
'nkx2',
'notch1',
'notch2',
'npm1',
'nras',
'nsd1',
'ntrk1',
'ntrk2',
'ntrk3',
'nup93',
'pax8',
'pbrm1',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pim1',
'pms2',
'pole',
'ppm1d',
'ppp2r1a',
'prdm1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rab35',
'rac1',
'rad51b',
'rad51c',
'rad51d',
'raf1',
'rara',
'rasa1',
'rb1',
'rbm10',
'ret',
'rheb',
'rhoa',
'riCTOR',
'rit1',


```

'ros1',
'rras2',
'runx1',
'rxra',
'rybp',
'setd2',
'sf3b1',
'shoc2',
'smad2',
'smad3',
'smad4',
'smarca4',
'smarca1',
'smo',
'sos1',
'sox9',
'spop',
'src',
'srsf2',
'stag2',
'stat3',
'stk11',
'tcf3',
'tcf7l2',
'tert',
'tet1',
'tet2',
'tgfbr1',
'tgfbr2',
'tmprss2',
'tp53',
'tsc1',
'tsc2',
'u2af1',
'vhl',
'whsc1',
'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, 234)

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))

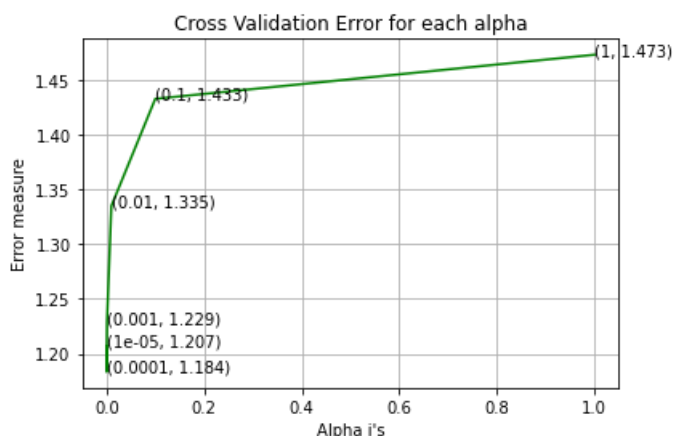
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error 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: 1.2069667368466748
 For values of alpha = 0.0001 The log loss is: 1.1835156544027445
 For values of alpha = 0.001 The log loss is: 1.228509161341078
 For values of alpha = 0.01 The log loss is: 1.3348419420608695
 For values of alpha = 0.1 The log loss is: 1.432556171878459
 For values of alpha = 1 The log loss is: 1.4727832817149002



For values of best alpha = 0.0001 The train log loss is: 1.0083122746528554
 For values of best alpha = 0.0001 The cross validation log loss is: 1.1835156544027445
 For values of best alpha = 0.0001 The test log loss is: 1.1509229268005292

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

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

In [27]:

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

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

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

Q6. How many data points in Test and CV datasets are covered by the 235 genes in train dataset?
Ans
1. In test data 644 out of 665 : 96.84210526315789
2. In cross validation data 515 out of 532 : 96.80451127819549

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 : 1928
Truncating_Mutations    58
Deletion                 53
Amplification            43
Fusions                  25
Q61L                     3
Overexpression           3
G12S                     2
Q61H                     2
A146T                    2
G13C                     2
Name: Variation, dtype: int64
```

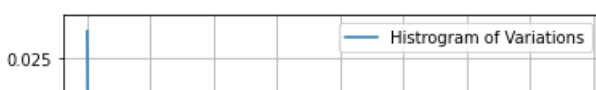
In [29]:

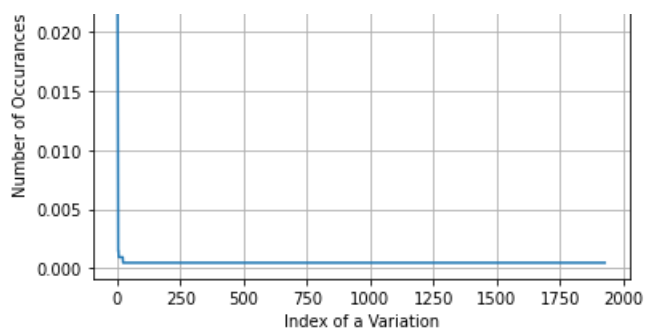
```
print("Ans: There are", unique_variations.shape[0], "different categories of variations in the train data, and they are distributed as follows",)
```

Ans: There are 1928 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()
```

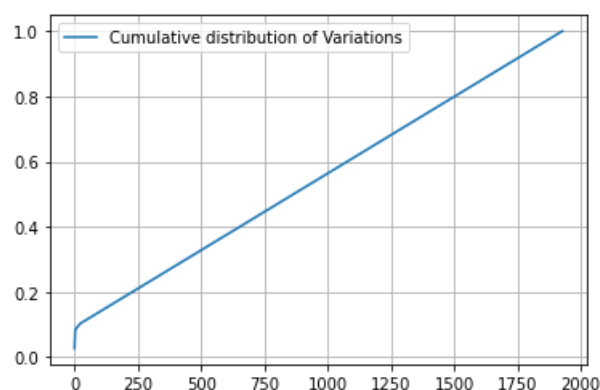




In [31]:

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

```
[0.02730697 0.05225989 0.07250471 ... 0.99905838 0.99952919 1.          ]
```



Q9. How to featurize this Variation feature ?

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

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

1. One hot Encoding
2. Response coding

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

In [32]:

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

In [33]:

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

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

In [34]:

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

In [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, 1961)

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()

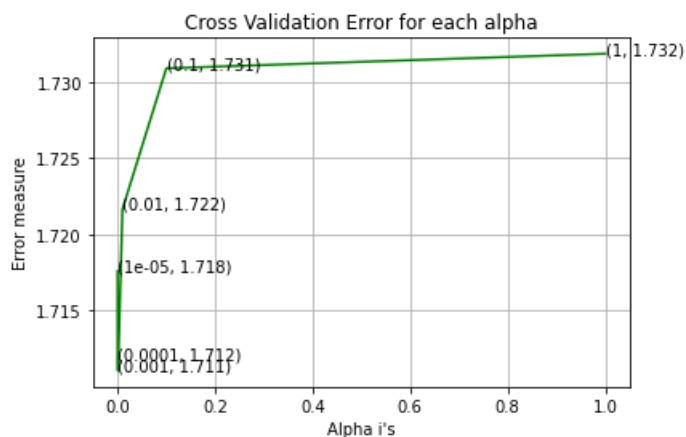
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_variation_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_variation_feature_onehotCoding, y_train)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

```

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

```

For values of alpha = 1e-05 The log loss is: 1.7176017836345205
 For values of alpha = 0.0001 The log loss is: 1.7118223558781864
 For values of alpha = 0.001 The log loss is: 1.7110465236911627
 For values of alpha = 0.01 The log loss is: 1.7216632702333743
 For values of alpha = 0.1 The log loss is: 1.7309149939294546
 For values of alpha = 1 The log loss is: 1.731871010517095



For values of best alpha = 0.001 The train log loss is: 1.0386392759485736
 For values of best alpha = 0.001 The cross validation log loss is: 1.7110465236911627
 For values of best alpha = 0.001 The test log loss is: 1.721847011143125

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

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

In [37]:

```

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

```

Q12. How many data points are covered by total 1928 genes in test and cross validation data sets?

Ans

1. In test data 72 out of 665 : 10.827067669172932
 2. In cross validation data 49 out of 532 : 9.210526315789473

3.2.3 Univariate Analysis on Text Feature

1. How many unique words are present in train data?
2. How are word frequencies distributed?
3. How to featurize text field?
4. Is the text feature useful in 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 [41]:

```

# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = CountVectorizer(min_df=5)
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 : 39400

In [42]:

```

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 [43]:

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

In [44]:

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

In [45]:

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

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

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

In [46]:

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

In [47]:

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

```
Counter({5: 2212, 7: 1721, 8: 1521, 6: 1500, 11: 1465, 9: 1192, 10: 1172, 12: 1026, 14: 956, 16:
743, 13: 721, 15: 677, 20: 618, 18: 604, 21: 530, 19: 516, 17: 501, 22: 500, 24: 456, 33: 434, 23:
392, 30: 375, 25: 335, 26: 304, 32: 303, 28: 301, 27: 299, 29: 266, 42: 261, 35: 248, 52: 245, 31:
239, 36: 232, 34: 231, 40: 214, 37: 205, 39: 185, 41: 180, 44: 177, 38: 177, 43: 168, 45: 165, 46:
159, 48: 158, 49: 157, 50: 153, 54: 151, 55: 150, 47: 150, 57: 132, 56: 132, 60: 129, 53: 129, 51:
124, 61: 123, 59: 121, 58: 117, 77: 116, 66: 111, 67: 110, 84: 109, 65: 108, 72: 103, 62: 103, 64:
99, 71: 96, 63: 92, 68: 91, 73: 90, 70: 85, 76: 81, 81: 80, 101: 77, 79: 77, 78: 77, 75: 72, 104:
71, 88: 70, 74: 70, 69: 70, 86: 69, 80: 69, 92: 68, 82: 66, 96: 63, 100: 60, 97: 60, 110: 59, 89:
59, 85: 59, 99: 58, 87: 57, 93: 56, 107: 55, 95: 55, 83: 55, 120: 53, 105: 53, 102: 53, 91: 53, 9
0: 53, 94: 51, 103: 49, 119: 48, 112: 47, 114: 46, 106: 46, 126: 45, 149: 44, 115: 44, 156: 43, 13
4: 43, 132: 43, 127: 43, 111: 43, 108: 43, 128: 42, 118: 42, 109: 42, 153: 39, 133: 39, 143: 38, 1
16: 38, 123: 37, 117: 37, 113: 37, 157: 36, 98: 36, 150: 35, 146: 35, 144: 35, 137: 35, 135: 35, 1
21: 35, 130: 34, 122: 34, 155: 33, 145: 33, 125: 33, 174: 32, 171: 32, 164: 31, 154: 31, 124: 31,
211: 30, 165: 30, 161: 30, 152: 30, 140: 30, 136: 30, 168: 29, 148: 29, 139: 29, 129: 29, 172: 27,
170: 27, 167: 27, 166: 27, 141: 27, 192: 26, 180: 26, 162: 26, 147: 26, 142: 26, 138: 26, 131: 26,
158: 25, 208: 24, 182: 24, 163: 24, 216: 23, 215: 23, 210: 23, 184: 23, 160: 23, 151: 23, 228: 22,
201: 22, 200: 22, 197: 22, 177: 22, 176: 22, 300: 21, 278: 21, 212: 21, 193: 21, 186: 21, 276: 20,
264: 20, 260: 20, 243: 20, 240: 20, 222: 20, 221: 20, 188: 20, 320: 19, 248: 19, 244: 19, 238: 19,
230: 19, 214: 19, 181: 19, 179: 19, 159: 19, 252: 18, 229: 18, 225: 18, 224: 18, 209: 18, 207: 18,
206: 18, 202: 18, 198: 18, 194: 18, 189: 18, 183: 18, 169: 18, 301: 17, 290: 17, 262: 17, 242: 17,
237: 17, 223: 17, 196: 17, 195: 17, 185: 17, 175: 17, 360: 16, 349: 16, 309: 16, 299: 16, 283: 16,
277: 16, 257: 16, 246: 16, 226: 16, 218: 16, 217: 16, 213: 16, 204: 16, 191: 16, 319: 15, 291: 15,
271: 15, 250: 15, 245: 15, 232: 15, 231: 15, 190: 15, 476: 14, 305: 14, 280: 14, 261: 14, 254: 14,
251: 14, 234: 14, 219: 14, 203: 14, 429: 13, 346: 13, 331: 13, 288: 13, 285: 13, 249: 13, 241: 13,
205: 13, 187: 13, 411: 12, 399: 12, 356: 12, 340: 12, 332: 12, 330: 12, 325: 12, 324: 12, 323: 12,
314: 12, 312: 12, 306: 12, 304: 12, 303: 12, 287: 12, 282: 12, 275: 12, 272: 12, 269: 12, 267: 12,
258: 12, 255: 12, 253: 12, 247: 12, 239: 12, 236: 12, 235: 12, 178: 12, 465: 11, 453: 11, 405: 11,
```


395: 11, 351: 11, 348: 11, 333: 11, 315: 11, 286: 11, 268: 11, 265: 11, 263: 11, 259: 11, 256: 11,
233: 11, 227: 11, 173: 11, 604: 10, 436: 10, 392: 10, 382: 10, 366: 10, 345: 10, 334: 10, 329: 10,
326: 10, 311: 10, 310: 10, 297: 10, 293: 10, 284: 10, 274: 10, 270: 10, 220: 10, 199: 10, 542: 9,
538: 9, 447: 9, 444: 9, 440: 9, 434: 9, 427: 9, 424: 9, 418: 9, 415: 9, 414: 9, 364: 9, 363: 9,
362: 9, 342: 9, 336: 9, 335: 9, 328: 9, 322: 9, 318: 9, 298: 9, 296: 9, 295: 9, 281: 9, 478: 8,
472: 8, 470: 8, 469: 8, 463: 8, 461: 8, 452: 8, 437: 8, 435: 8, 401: 8, 394: 8, 375: 8, 373: 8,
370: 8, 359: 8, 353: 8, 350: 8, 302: 8, 294: 8, 292: 8, 266: 8, 1033: 7, 729: 7, 698: 7, 691: 7,
678: 7, 655: 7, 639: 7, 564: 7, 545: 7, 531: 7, 526: 7, 497: 7, 485: 7, 454: 7, 448: 7, 438: 7,
421: 7, 420: 7, 398: 7, 388: 7, 377: 7, 367: 7, 365: 7, 352: 7, 344: 7, 338: 7, 327: 7, 317: 7,
308: 7, 279: 7, 1414: 6, 946: 6, 922: 6, 838: 6, 826: 6, 816: 6, 813: 6, 732: 6, 722: 6, 711: 6,
683: 6, 660: 6, 654: 6, 635: 6, 627: 6, 626: 6, 603: 6, 587: 6, 574: 6, 525: 6, 513: 6, 511: 6,
506: 6, 493: 6, 483: 6, 466: 6, 459: 6, 445: 6, 428: 6, 425: 6, 422: 6, 416: 6, 402: 6, 391: 6,
380: 6, 379: 6, 372: 6, 361: 6, 358: 6, 355: 6, 354: 6, 347: 6, 313: 6, 307: 6, 273: 6, 1773: 5,
1559: 5, 1154: 5, 1139: 5, 1005: 5, 936: 5, 911: 5, 904: 5, 884: 5, 846: 5, 825: 5, 801: 5, 800: 5,
799: 5, 795: 5, 794: 5, 750: 5, 740: 5, 724: 5, 712: 5, 707: 5, 679: 5, 676: 5, 675: 5, 658: 5,
633: 5, 631: 5, 618: 5, 615: 5, 595: 5, 592: 5, 591: 5, 589: 5, 584: 5, 583: 5, 580: 5, 577: 5,
572: 5, 571: 5, 566: 5, 565: 5, 561: 5, 555: 5, 551: 5, 550: 5, 549: 5, 536: 5, 534: 5, 529: 5,
519: 5, 517: 5, 510: 5, 499: 5, 494: 5, 488: 5, 475: 5, 473: 5, 471: 5, 467: 5, 462: 5, 460: 5,
456: 5, 455: 5, 451: 5, 449: 5, 439: 5, 410: 5, 408: 5, 400: 5, 397: 5, 393: 5, 390: 5, 389: 5,
386: 5, 385: 5, 384: 5, 383: 5, 378: 5, 374: 5, 371: 5, 357: 5, 339: 5, 337: 5, 321: 5, 1989: 4,
1815: 4, 1777: 4, 1657: 4, 1607: 4, 1420: 4, 1381: 4, 1343: 4, 1341: 4, 1326: 4, 1313: 4, 1291: 4,
1248: 4, 1241: 4, 1200: 4, 1195: 4, 1167: 4, 1163: 4, 1160: 4, 1148: 4, 1110: 4, 1076: 4, 1049: 4,
1040: 4, 1024: 4, 1006: 4, 982: 4, 919: 4, 909: 4, 907: 4, 895: 4, 889: 4, 886: 4, 868: 4, 858: 4,
857: 4, 847: 4, 836: 4, 827: 4, 814: 4, 811: 4, 790: 4, 786: 4, 784: 4, 772: 4, 766: 4, 758: 4,
751: 4, 735: 4, 733: 4, 723: 4, 718: 4, 708: 4, 700: 4, 693: 4, 684: 4, 680: 4, 672: 4, 671: 4,
670: 4, 669: 4, 661: 4, 656: 4, 652: 4, 643: 4, 638: 4, 637: 4, 630: 4, 628: 4, 625: 4, 623: 4,
622: 4, 620: 4, 617: 4, 616: 4, 610: 4, 609: 4, 608: 4, 607: 4, 605: 4, 601: 4, 599: 4, 598: 4,
597: 4, 596: 4, 594: 4, 593: 4, 582: 4, 581: 4, 573: 4, 567: 4, 562: 4, 556: 4, 547: 4, 546: 4,
541: 4, 540: 4, 539: 4, 535: 4, 532: 4, 530: 4, 528: 4, 523: 4, 521: 4, 520: 4, 516: 4, 509: 4,
505: 4, 498: 4, 496: 4, 495: 4, 491: 4, 489: 4, 487: 4, 486: 4, 481: 4, 480: 4, 477: 4, 458: 4,
450: 4, 433: 4, 431: 4, 430: 4, 423: 4, 419: 4, 417: 4, 412: 4, 409: 4, 404: 4, 396: 4, 387: 4,
381: 4, 369: 4, 341: 4, 316: 4, 3853: 3, 3573: 3, 3541: 3, 3261: 3, 2657: 3, 2489: 3, 2432: 3, 2418
: 3, 2361: 3, 2326: 3, 2241: 3, 2239: 3, 2198: 3, 2084: 3, 1980: 3, 1979: 3, 1973: 3, 1970: 3, 1936
: 3, 1901: 3, 1867: 3, 1825: 3, 1809: 3, 1735: 3, 1727: 3, 1713: 3, 1700: 3, 1692: 3, 1675: 3, 1539
: 3, 1477: 3, 1454: 3, 1451: 3, 1438: 3, 1412: 3, 1397: 3, 1386: 3, 1385: 3, 1364: 3, 1348: 3, 1334
: 3, 1328: 3, 1305: 3, 1301: 3, 1293: 3, 1289: 3, 1283: 3, 1274: 3, 1272: 3, 1269: 3, 1260: 3, 1236
: 3, 1213: 3, 1203: 3, 1194: 3, 1192: 3, 1184: 3, 1180: 3, 1175: 3, 1157: 3, 1151: 3, 1122: 3, 1088
: 3, 1074: 3, 1073: 3, 1064: 3, 1058: 3, 1050: 3, 1038: 3, 1036: 3, 1031: 3, 1029: 3, 1015: 3, 1014
: 3, 1003: 3, 998: 3, 995: 3, 994: 3, 986: 3, 984: 3, 974: 3, 971: 3, 970: 3, 968: 3, 957: 3, 950:
3, 940: 3, 929: 3, 920: 3, 917: 3, 901: 3, 900: 3, 888: 3, 883: 3, 874: 3, 872: 3, 871: 3, 866: 3,
864: 3, 861: 3, 860: 3, 855: 3, 853: 3, 852: 3, 849: 3, 845: 3, 841: 3, 840: 3, 833: 3, 831: 3,
828: 3, 824: 3, 823: 3, 819: 3, 818: 3, 806: 3, 802: 3, 797: 3, 793: 3, 791: 3, 787: 3, 782: 3,
778: 3, 773: 3, 771: 3, 770: 3, 769: 3, 768: 3, 762: 3, 757: 3, 749: 3, 746: 3, 742: 3, 737: 3,
721: 3, 717: 3, 713: 3, 710: 3, 696: 3, 681: 3, 677: 3, 673: 3, 668: 3, 666: 3, 651: 3, 645: 3,
632: 3, 619: 3, 614: 3, 590: 3, 588: 3, 579: 3, 575: 3, 570: 3, 568: 3, 563: 3, 558: 3, 554: 3,
553: 3, 552: 3, 533: 3, 524: 3, 522: 3, 515: 3, 514: 3, 508: 3, 504: 3, 503: 3, 502: 3, 500: 3,
492: 3, 490: 3, 484: 3, 479: 3, 474: 3, 468: 3, 464: 3, 457: 3, 426: 3, 407: 3, 403: 3, 368: 3,
289: 3, 12561: 2, 10405: 2, 8190: 2, 6300: 2, 6263: 2, 5478: 2, 5471: 2, 5233: 2, 5000: 2, 4940: 2,
4874: 2, 4596: 2, 4356: 2, 4068: 2, 3999: 2, 3968: 2, 3926: 2, 3906: 2, 3767: 2, 3635: 2, 3599: 2,
3590: 2, 3588: 2, 3570: 2, 3514: 2, 3510: 2, 3509: 2, 3507: 2, 3499: 2, 3408: 2, 3384: 2, 3378: 2,
3321: 2, 3289: 2, 3131: 2, 3127: 2, 3120: 2, 3079: 2, 3035: 2, 3032: 2, 3030: 2, 3016: 2, 2963: 2,
2925: 2, 2827: 2, 2738: 2, 2721: 2, 2706: 2, 2697: 2, 2663: 2, 2645: 2, 2642: 2, 2636: 2, 2624: 2,
2623: 2, 2595: 2, 2593: 2, 2570: 2, 2567: 2, 2544: 2, 2537: 2, 2529: 2, 2516: 2, 2510: 2, 2509: 2,
2446: 2, 2426: 2, 2314: 2, 2280: 2, 2265: 2, 2189: 2, 2181: 2, 2176: 2, 2171: 2, 2166: 2, 2160: 2,
2151: 2, 2147: 2, 2142: 2, 2114: 2, 2113: 2, 2088: 2, 2082: 2, 2068: 2, 2058: 2, 2027: 2, 2024: 2,
2023: 2, 2019: 2, 2005: 2, 2004: 2, 1954: 2, 1949: 2, 1937: 2, 1935: 2, 1928: 2, 1926: 2, 1918: 2,
1912: 2, 1892: 2, 1878: 2, 1875: 2, 1872: 2, 1868: 2, 1843: 2, 1837: 2, 1832: 2, 1831: 2, 1811: 2,
1810: 2, 1794: 2, 1790: 2, 1785: 2, 1781: 2, 1767: 2, 1749: 2, 1747: 2, 1738: 2, 1725: 2, 1720: 2,
1715: 2, 1711: 2, 1705: 2, 1685: 2, 1682: 2, 1677: 2, 1638: 2, 1619: 2, 1603: 2, 1600: 2, 1598: 2,
1593: 2, 1590: 2, 1589: 2, 1581: 2, 1573: 2, 1571: 2, 1565: 2, 1562: 2, 1560: 2, 1557: 2, 1547: 2,
1540: 2, 1535: 2, 1531: 2, 1528: 2, 1526: 2, 1524: 2, 1518: 2, 1517: 2, 1510: 2, 1499: 2, 1498: 2,
1497: 2, 1495: 2, 1494: 2, 1490: 2, 1489: 2, 1482: 2, 1467: 2, 1466: 2, 1461: 2, 1459: 2, 1453: 2,
1450: 2, 1448: 2, 1436: 2, 1433: 2, 1431: 2, 1423: 2, 1417: 2, 1416: 2, 1405: 2, 1404: 2, 1401: 2,
1398: 2, 1389: 2, 1388: 2, 1380: 2, 1376: 2, 1374: 2, 1372: 2, 1367: 2, 1357: 2, 1354: 2, 1353: 2,
1352: 2, 1337: 2, 1332: 2, 1324: 2, 1323: 2, 1320: 2, 1316: 2, 1312: 2, 1307: 2, 1303: 2, 1299: 2,
1298: 2, 1297: 2, 1286: 2, 1285: 2, 1281: 2, 1278: 2, 1273: 2, 1270: 2, 1266: 2, 1265: 2, 1252: 2,
1250: 2, 1239: 2, 1237: 2, 1235: 2, 1221: 2, 1209: 2, 1193: 2, 1189: 2, 1187: 2, 1181: 2, 1179: 2,
1168: 2, 1153: 2, 1152: 2, 1150: 2, 1146: 2, 1143: 2, 1140: 2, 1136: 2, 1135: 2, 1131: 2, 1129: 2,
1124: 2, 1123: 2, 1118: 2, 1113: 2, 1105: 2, 1099: 2, 1096: 2, 1093: 2, 1087: 2, 1086: 2, 1085: 2,
1083: 2, 1070: 2, 1067: 2, 1060: 2, 1057: 2, 1056: 2, 1055: 2, 1054: 2, 1052: 2, 1045: 2, 1044: 2,
1042: 2, 1039: 2, 1035: 2, 1030: 2, 1026: 2, 1023: 2, 1022: 2, 1019: 2, 1016: 2, 1012: 2, 1002: 2,
999: 2, 993: 2, 992: 2, 990: 2, 989: 2, 983: 2, 978: 2, 977: 2, 976: 2, 975: 2, 972: 2, 969: 2,
959: 2, 956: 2, 953: 2, 948: 2, 947: 2, 945: 2, 944: 2, 942: 2, 941: 2, 939: 2, 934: 2, 933: 2,
932: 2, 930: 2, 926: 2, 925: 2, 918: 2, 913: 2, 910: 2, 908: 2, 903: 2, 899: 2, 898: 2, 896: 2,
892: 2, 885: 2, 881: 2, 880: 2, 879: 2, 876: 2, 873: 2, 851: 2, 843: 2, 837: 2, 830: 2, 829: 2,
822: 2, 820: 2, 812: 2, 808: 2, 804: 2, 796: 2, 781: 2, 780: 2, 779: 2, 777: 2, 776: 2, 775: 2,
763: 2, 761: 2, 759: 2, 756: 2, 755: 2, 754: 2, 744: 2, 743: 2, 741: 2, 738: 2, 734: 2, 730: 2,

726: 2, 725: 2, 716: 2, 715: 2, 709: 2, 706: 2, 705: 2, 703: 2, 702: 2, 697: 2, 694: 2, 692: 2,
690: 2, 689: 2, 688: 2, 687: 2, 686: 2, 685: 2, 682: 2, 674: 2, 665: 2, 663: 2, 662: 2, 649: 2,
647: 2, 644: 2, 641: 2, 629: 2, 612: 2, 611: 2, 606: 2, 602: 2, 600: 2, 578: 2, 576: 2, 569: 2,
560: 2, 559: 2, 557: 2, 548: 2, 543: 2, 537: 2, 527: 2, 507: 2, 501: 2, 482: 2, 446: 2, 443: 2,
442: 2, 432: 2, 413: 2, 406: 2, 376: 2, 343: 2, 156486: 1, 123226: 1, 82221: 1, 69288: 1, 69052: 1
, 68888: 1, 67153: 1, 66527: 1, 63481: 1, 57803: 1, 54230: 1, 50994: 1, 50696: 1, 48578: 1, 46616:
1, 45275: 1, 43681: 1, 42734: 1, 42399: 1, 42284: 1, 41735: 1, 40874: 1, 40126: 1, 39430: 1, 39056
: 1, 38612: 1, 37838: 1, 37246: 1, 36623: 1, 36210: 1, 35617: 1, 34217: 1, 34161: 1, 33842: 1, 334
03: 1, 32466: 1, 30319: 1, 28767: 1, 28659: 1, 27046: 1, 26437: 1, 26287: 1, 26224: 1, 26092: 1, 2
5510: 1, 25307: 1, 25289: 1, 25227: 1, 25131: 1, 25003: 1, 24813: 1, 24458: 1, 24167: 1, 23472: 1,
22814: 1, 22458: 1, 22187: 1, 22142: 1, 21812: 1, 21670: 1, 21501: 1, 21317: 1, 20890: 1, 20682: 1
, 20212: 1, 20049: 1, 20006: 1, 19960: 1, 19727: 1, 19554: 1, 19293: 1, 19286: 1, 19153: 1, 19133:
1, 19066: 1, 19011: 1, 18895: 1, 18713: 1, 18468: 1, 18384: 1, 18334: 1, 18160: 1, 18146: 1, 18057
: 1, 18012: 1, 17921: 1, 17871: 1, 17840: 1, 17751: 1, 17741: 1, 17667: 1, 17629: 1, 17191: 1, 171
83: 1, 16985: 1, 16869: 1, 16711: 1, 16698: 1, 16670: 1, 16336: 1, 16327: 1, 16166: 1, 16133: 1, 1
5963: 1, 15852: 1, 15780: 1, 15693: 1, 15668: 1, 15559: 1, 15331: 1, 15325: 1, 15228: 1, 15055: 1,
14970: 1, 14877: 1, 14740: 1, 14724: 1, 14654: 1, 14522: 1, 14511: 1, 14502: 1, 14407: 1, 14314: 1
, 14094: 1, 13970: 1, 13852: 1, 13744: 1, 13682: 1, 13649: 1, 13606: 1, 13442: 1, 13437: 1, 13403:
1, 13236: 1, 13194: 1, 13191: 1, 13169: 1, 13100: 1, 13032: 1, 13005: 1, 13001: 1, 12905: 1, 12828
: 1, 12782: 1, 12718: 1, 12697: 1, 12651: 1, 12630: 1, 12628: 1, 12593: 1, 12592: 1, 12584: 1, 125
36: 1, 12517: 1, 12494: 1, 12471: 1, 12445: 1, 12420: 1, 12359: 1, 12350: 1, 12341: 1, 12326: 1, 1
2303: 1, 12200: 1, 12163: 1, 12154: 1, 12135: 1, 12110: 1, 11913: 1, 11883: 1, 11832: 1, 11816: 1,
11777: 1, 11566: 1, 11536: 1, 11417: 1, 11351: 1, 11298: 1, 11284: 1, 11262: 1, 11233: 1, 11227: 1
, 11131: 1, 11114: 1, 10977: 1, 10900: 1, 10733: 1, 10712: 1, 10699: 1, 10673: 1, 10619: 1, 10605:
1, 10577: 1, 10434: 1, 10381: 1, 10311: 1, 10307: 1, 10267: 1, 10260: 1, 10253: 1, 10247: 1, 10213
: 1, 10172: 1, 10146: 1, 10121: 1, 10053: 1, 10040: 1, 9924: 1, 9833: 1, 9808: 1, 9787: 1, 9746: 1,
9730: 1, 9721: 1, 9703: 1, 9690: 1, 9678: 1, 9625: 1, 9532: 1, 9504: 1, 9434: 1, 9432: 1, 9380: 1,
9377: 1, 9374: 1, 9363: 1, 9360: 1, 9357: 1, 9344: 1, 9315: 1, 9228: 1, 9209: 1, 9129: 1, 9105: 1,
9081: 1, 9056: 1, 8983: 1, 8975: 1, 8936: 1, 8922: 1, 8917: 1, 8863: 1, 8803: 1, 8770: 1, 8737: 1,
8701: 1, 8650: 1, 8593: 1, 8586: 1, 8537: 1, 8525: 1, 8519: 1, 8511: 1, 8495: 1, 8469: 1, 8456: 1,
8446: 1, 8438: 1, 8426: 1, 8424: 1, 8423: 1, 8335: 1, 8327: 1, 8278: 1, 8271: 1, 8263: 1, 8228: 1,
8204: 1, 8194: 1, 8181: 1, 8175: 1, 8169: 1, 8167: 1, 8108: 1, 8100: 1, 8072: 1, 8067: 1, 8057: 1,
8003: 1, 7970: 1, 7959: 1, 7887: 1, 7885: 1, 7877: 1, 7834: 1, 7786: 1, 7757: 1, 7744: 1, 7717: 1,
7710: 1, 7693: 1, 7652: 1, 7591: 1, 7570: 1, 7560: 1, 7559: 1, 7526: 1, 7507: 1, 7506: 1, 7447: 1,
7436: 1, 7426: 1, 7424: 1, 7416: 1, 7408: 1, 7404: 1, 7386: 1, 7384: 1, 7369: 1, 7336: 1, 7292: 1,
7285: 1, 7264: 1, 7262: 1, 7237: 1, 7218: 1, 7214: 1, 7201: 1, 7164: 1, 7143: 1, 7142: 1, 7137: 1,
7127: 1, 7115: 1, 7114: 1, 7099: 1, 7098: 1, 7080: 1, 7079: 1, 7072: 1, 7064: 1, 7062: 1, 7056: 1,
6946: 1, 6937: 1, 6927: 1, 6924: 1, 6915: 1, 6868: 1, 6862: 1, 6847: 1, 6834: 1, 6833: 1, 6826: 1,
6795: 1, 6772: 1, 6762: 1, 6761: 1, 6741: 1, 6727: 1, 6721: 1, 6705: 1, 6694: 1, 6690: 1, 6659: 1,
6658: 1, 6637: 1, 6601: 1, 6586: 1, 6576: 1, 6548: 1, 6531: 1, 6526: 1, 6525: 1, 6506: 1, 6483: 1,
6478: 1, 6467: 1, 6458: 1, 6456: 1, 6452: 1, 6447: 1, 6446: 1, 6440: 1, 6432: 1, 6417: 1, 6365: 1,
6347: 1, 6317: 1, 6313: 1, 6305: 1, 6304: 1, 6303: 1, 6282: 1, 6276: 1, 6241: 1, 6228: 1, 6225: 1,
6223: 1, 6220: 1, 6214: 1, 6175: 1, 6174: 1, 6123: 1, 6107: 1, 6088: 1, 6081: 1, 6049: 1, 6039: 1,
6019: 1, 6015: 1, 6014: 1, 6003: 1, 5998: 1, 5997: 1, 5961: 1, 5928: 1, 5923: 1, 5921: 1, 5881: 1,
5871: 1, 5869: 1, 5822: 1, 5820: 1, 5776: 1, 5775: 1, 5770: 1, 5760: 1, 5757: 1, 5751: 1, 5741: 1,
5733: 1, 5712: 1, 5670: 1, 5653: 1, 5616: 1, 5606: 1, 5602: 1, 5601: 1, 5600: 1, 5594: 1, 5591: 1,
5586: 1, 5584: 1, 5574: 1, 5570: 1, 5566: 1, 5565: 1, 5558: 1, 5557: 1, 5547: 1, 5522: 1, 5511: 1,
5505: 1, 5483: 1, 5460: 1, 5456: 1, 5441: 1, 5391: 1, 5370: 1, 5365: 1, 5364: 1, 5357: 1, 5344: 1,
5342: 1, 5320: 1, 5311: 1, 5299: 1, 5298: 1, 5285: 1, 5280: 1, 5273: 1, 5250: 1, 5227: 1, 5219: 1,
5198: 1, 5187: 1, 5149: 1, 5146: 1, 5140: 1, 5132: 1, 5130: 1, 5126: 1, 5122: 1, 5108: 1, 5107: 1,
5106: 1, 5091: 1, 5088: 1, 5086: 1, 5081: 1, 5067: 1, 5059: 1, 5030: 1, 5020: 1, 5017: 1, 5011: 1,
4990: 1, 4980: 1, 4976: 1, 4975: 1, 4964: 1, 4924: 1, 4914: 1, 4910: 1, 4908: 1, 4903: 1, 4896: 1,
4890: 1, 4856: 1, 4813: 1, 4803: 1, 4782: 1, 4770: 1, 4767: 1, 4761: 1, 4755: 1, 4754: 1, 4746: 1,
4738: 1, 4730: 1, 4727: 1, 4722: 1, 4721: 1, 4704: 1, 4700: 1, 4697: 1, 4688: 1, 4660: 1, 4656: 1,
4655: 1, 4646: 1, 4640: 1, 4633: 1, 4606: 1, 4595: 1, 4589: 1, 4575: 1, 4572: 1, 4559: 1, 4556: 1,
4549: 1, 4548: 1, 4540: 1, 4539: 1, 4529: 1, 4527: 1, 4526: 1, 4520: 1, 4511: 1, 4507: 1, 4487: 1,
4462: 1, 4456: 1, 4454: 1, 4450: 1, 4448: 1, 4444: 1, 4438: 1, 4434: 1, 4426: 1, 4425: 1, 4418: 1,
4407: 1, 4406: 1, 4402: 1, 4393: 1, 4386: 1, 4372: 1, 4368: 1, 4347: 1, 4345: 1, 4342: 1, 4339: 1,
4332: 1, 4327: 1, 4320: 1, 4317: 1, 4312: 1, 4304: 1, 4300: 1, 4298: 1, 4296: 1, 4283: 1, 4280: 1,
4273: 1, 4270: 1, 4266: 1, 4254: 1, 4249: 1, 4248: 1, 4247: 1, 4244: 1, 4237: 1, 4216: 1, 4212: 1,
4209: 1, 4203: 1, 4197: 1, 4193: 1, 4190: 1, 4186: 1, 4169: 1, 4157: 1, 4154: 1, 4153: 1, 4152: 1,
4151: 1, 4139: 1, 4133: 1, 4132: 1, 4125: 1, 4115: 1, 4113: 1, 4109: 1, 4107: 1, 4098: 1, 4090: 1,
4088: 1, 4079: 1, 4073: 1, 4064: 1, 4062: 1, 4056: 1, 4048: 1, 4036: 1, 4021: 1, 4017: 1, 4003: 1,
3994: 1, 3991: 1, 3981: 1, 3980: 1, 3976: 1, 3972: 1, 3966: 1, 3954: 1, 3936: 1, 3934: 1, 3929: 1,
3920: 1, 3914: 1, 3912: 1, 3903: 1, 3896: 1, 3895: 1, 3890: 1, 3888: 1, 3886: 1, 3885: 1, 3878: 1,
3877: 1, 3861: 1, 3858: 1, 3850: 1, 3841: 1, 3824: 1, 3819: 1, 3817: 1, 3814: 1, 3807: 1, 3796: 1,
3793: 1, 3783: 1, 3781: 1, 3779: 1, 3776: 1, 3773: 1, 3766: 1, 3765: 1, 3759: 1, 3750: 1, 3745: 1,
3733: 1, 3730: 1, 3722: 1, 3720: 1, 3704: 1, 3701: 1, 3696: 1, 3690: 1, 3688: 1, 3684: 1, 3682: 1,
3678: 1, 3677: 1, 3676: 1, 3675: 1, 3674: 1, 3672: 1, 3664: 1, 3662: 1, 3661: 1, 3656: 1, 3653: 1,
3647: 1, 3646: 1, 3645: 1, 3644: 1, 3641: 1, 3639: 1, 3630: 1, 3628: 1, 3623: 1, 3604: 1, 3583: 1,
3569: 1, 3565: 1, 3561: 1, 3559: 1, 3552: 1, 3545: 1, 3542: 1, 3525: 1, 3524: 1, 3521: 1, 3518: 1,
3513: 1, 3508: 1, 3504: 1, 3498: 1, 3496: 1, 3489: 1, 3485: 1, 3484: 1, 3449: 1, 3440: 1, 3437: 1,
3435: 1, 3431: 1, 3427: 1, 3425: 1, 3419: 1, 3418: 1, 3413: 1, 3402: 1, 3400: 1, 3397: 1, 3393: 1,
3392: 1, 3388: 1, 3371: 1, 3369: 1, 3353: 1, 3352: 1, 3350: 1, 3348: 1, 3347: 1, 3345: 1, 3342: 1,
3337: 1, 3335: 1, 3328: 1, 3326: 1, 3324: 1, 3313: 1, 3309: 1, 3308: 1, 3304: 1, 3302: 1, 3298: 1,
3293: 1, 3286: 1, 3285: 1, 3281: 1, 3277: 1, 3272: 1, 3262: 1, 3256: 1, 3254: 1, 3248: 1, 3239: 1,
3233: 1, 3227: 1, 3224: 1, 3215: 1, 3213: 1, 3212: 1, 3211: 1, 3202: 1, 3198: 1, 3195: 1, 3189: 1.

```

3188: 1, 3172: 1, 3171: 1, 3167: 1, 3164: 1, 3155: 1, 3142: 1, 3141: 1, 3136: 1, 3128: 1, 3122: 1,
3109: 1, 3104: 1, 3100: 1, 3097: 1, 3096: 1, 3095: 1, 3090: 1, 3089: 1, 3084: 1, 3083: 1, 3070: 1,
3065: 1, 3060: 1, 3056: 1, 3054: 1, 3053: 1, 3045: 1, 3043: 1, 3040: 1, 3033: 1, 3031: 1, 3009: 1,
3008: 1, 3007: 1, 3005: 1, 3002: 1, 2998: 1, 2991: 1, 2988: 1, 2986: 1, 2981: 1, 2974: 1, 2970: 1,
2958: 1, 2956: 1, 2946: 1, 2932: 1, 2921: 1, 2920: 1, 2916: 1, 2890: 1, 2885: 1, 2883: 1, 2870: 1,
2866: 1, 2865: 1, 2848: 1, 2847: 1, 2825: 1, 2823: 1, 2821: 1, 2819: 1, 2810: 1, 2796: 1, 2789: 1,
2786: 1, 2778: 1, 2777: 1, 2776: 1, 2769: 1, 2768: 1, 2767: 1, 2766: 1, 2761: 1, 2760: 1, 2757: 1,
2753: 1, 2752: 1, 2751: 1, 2737: 1, 2736: 1, 2735: 1, 2730: 1, 2727: 1, 2715: 1, 2712: 1, 2710: 1,
2705: 1, 2704: 1, 2698: 1, 2692: 1, 2688: 1, 2687: 1, 2686: 1, 2685: 1, 2681: 1, 2676: 1, 2675: 1,
2669: 1, 2667: 1, 2666: 1, 2664: 1, 2660: 1, 2659: 1, 2658: 1, 2655: 1, 2654: 1, 2652: 1, 2650: 1,
2649: 1, 2647: 1, 2629: 1, 2625: 1, 2622: 1, 2621: 1, 2618: 1, 2617: 1, 2605: 1, 2601: 1, 2597: 1,
2596: 1, 2594: 1, 2585: 1, 2584: 1, 2583: 1, 2582: 1, 2575: 1, 2569: 1, 2564: 1, 2563: 1, 2556: 1,
2553: 1, 2550: 1, 2548: 1, 2545: 1, 2536: 1, 2531: 1, 2530: 1, 2524: 1, 2522: 1, 2515: 1, 2513: 1,
2511: 1, 2507: 1, 2500: 1, 2497: 1, 2495: 1, 2493: 1, 2488: 1, 2487: 1, 2486: 1, 2480: 1, 2476: 1,
2473: 1, 2472: 1, 2470: 1, 2467: 1, 2466: 1, 2464: 1, 2456: 1, 2449: 1, 2448: 1, 2447: 1, 2444: 1,
2443: 1, 2442: 1, 2436: 1, 2434: 1, 2429: 1, 2419: 1, 2414: 1, 2413: 1, 2409: 1, 2405: 1, 2402: 1,
2401: 1, 2397: 1, 2394: 1, 2393: 1, 2391: 1, 2390: 1, 2387: 1, 2384: 1, 2383: 1, 2380: 1, 2378: 1,
2376: 1, 2375: 1, 2374: 1, 2373: 1, 2371: 1, 2370: 1, 2366: 1, 2363: 1, 2362: 1, 2357: 1, 2354: 1,
2352: 1, 2350: 1, 2347: 1, 2342: 1, 2336: 1, 2335: 1, 2334: 1, 2333: 1, 2331: 1, 2329: 1, 2328: 1,
2327: 1, 2322: 1, 2306: 1, 2305: 1, 2304: 1, 2300: 1, 2297: 1, 2296: 1, 2288: 1, 2276: 1, 2273: 1,
2267: 1, 2263: 1, 2256: 1, 2252: 1, 2246: 1, 2245: 1, 2242: 1, 2235: 1, 2233: 1, 2228: 1, 2226: 1,
2224: 1, 2220: 1, 2217: 1, 2216: 1, 2215: 1, 2213: 1, 2212: 1, 2210: 1, 2209: 1, 2206: 1, 2199: 1,
2192: 1, 2187: 1, 2186: 1, 2182: 1, 2180: 1, 2179: 1, 2174: 1, 2168: 1, 2167: 1, 2163: 1, 2162: 1,
2161: 1, 2156: 1, 2155: 1, 2152: 1, 2146: 1, 2144: 1, 2143: 1, 2141: 1, 2139: 1, 2132: 1, 2131: 1,
2128: 1, 2126: 1, 2125: 1, 2124: 1, 2110: 1, 2108: 1, 2106: 1, 2103: 1, 2102: 1, 2101: 1, 2099: 1,
2098: 1, 2097: 1, 2093: 1, 2091: 1, 2087: 1, 2080: 1, 2078: 1, 2070: 1, 2067: 1, 2065: 1, 2064: 1,
2062: 1, 2055: 1, 2048: 1, 2047: 1, 2046: 1, 2042: 1, 2040: 1, 2038: 1, 2034: 1, 2033: 1, 2032: 1,
2030: 1, 2025: 1, 2014: 1, 2013: 1, 2012: 1, 2010: 1, 2007: 1, 2006: 1, 2001: 1, 1996: 1, 1992: 1,
1986: 1, 1985: 1, 1983: 1, 1976: 1, 1971: 1, 1968: 1, 1966: 1, 1965: 1, 1964: 1, 1963: 1, 1962: 1,
1961: 1, 1957: 1, 1953: 1, 1952: 1, 1948: 1, 1946: 1, 1944: 1, 1942: 1, 1940: 1, 1938: 1, 1932: 1,
1930: 1, 1925: 1, 1920: 1, 1919: 1, 1914: 1, 1910: 1, 1906: 1, 1905: 1, 1904: 1, 1896: 1, 1895: 1,
1889: 1, 1888: 1, 1887: 1, 1885: 1, 1883: 1, 1881: 1, 1879: 1, 1877: 1, 1873: 1, 1870: 1, 1869: 1,
1865: 1, 1864: 1, 1862: 1, 1860: 1, 1859: 1, 1857: 1, 1851: 1, 1850: 1, 1847: 1, 1842: 1, 1841: 1,
1839: 1, 1838: 1, 1834: 1, 1827: 1, 1823: 1, 1819: 1, 1818: 1, 1816: 1, 1813: 1, 1808: 1, 1807: 1,
1803: 1, 1802: 1, 1800: 1, 1797: 1, 1796: 1, 1792: 1, 1791: 1, 1784: 1, 1776: 1, 1774: 1, 1770: 1,
1769: 1, 1768: 1, 1764: 1, 1762: 1, 1758: 1, 1755: 1, 1752: 1, 1751: 1, 1746: 1, 1744: 1, 1742: 1,
1739: 1, 1723: 1, 1722: 1, 1719: 1, 1718: 1, 1717: 1, 1716: 1, 1712: 1, 1709: 1, 1708: 1, 1704: 1,
1703: 1, 1702: 1, 1701: 1, 1699: 1, 1698: 1, 1696: 1, 1691: 1, 1689: 1, 1688: 1, 1686: 1, 1681: 1,
1680: 1, 1674: 1, 1672: 1, 1671: 1, 1670: 1, 1666: 1, 1665: 1, 1663: 1, 1662: 1, 1658: 1, 1655: 1,
1654: 1, 1652: 1, 1649: 1, 1648: 1, 1647: 1, 1645: 1, 1644: 1, 1642: 1, 1640: 1, 1636: 1, 1631: 1,
1629: 1, 1625: 1, 1624: 1, 1622: 1, 1621: 1, 1620: 1, 1618: 1, 1617: 1, 1613: 1, 1612: 1, 1610: 1,
1604: 1, 1602: 1, 1601: 1, 1597: 1, 1594: 1, 1587: 1, 1584: 1, 1583: 1, 1576: 1, 1574: 1, 1572: 1,
1570: 1, 1569: 1, 1563: 1, 1558: 1, 1556: 1, 1555: 1, 1554: 1, 1553: 1, 1552: 1, 1550: 1, 1546: 1,
1543: 1, 1542: 1, 1537: 1, 1534: 1, 1530: 1, 1525: 1, 1523: 1, 1516: 1, 1515: 1, 1513: 1, 1509: 1,
1506: 1, 1505: 1, 1502: 1, 1501: 1, 1500: 1, 1488: 1, 1487: 1, 1480: 1, 1476: 1, 1474: 1, 1473: 1,
1471: 1, 1470: 1, 1469: 1, 1465: 1, 1463: 1, 1460: 1, 1457: 1, 1452: 1, 1447: 1, 1446: 1, 1445: 1,
1444: 1, 1443: 1, 1442: 1, 1440: 1, 1437: 1, 1434: 1, 1429: 1, 1426: 1, 1425: 1, 1422: 1, 1421: 1,
1418: 1, 1415: 1, 1413: 1, 1411: 1, 1410: 1, 1407: 1, 1403: 1, 1402: 1, 1400: 1, 1399: 1, 1395: 1,
1394: 1, 1392: 1, 1391: 1, 1387: 1, 1379: 1, 1368: 1, 1365: 1, 1363: 1, 1362: 1, 1360: 1, 1359: 1,
1358: 1, 1356: 1, 1355: 1, 1350: 1, 1347: 1, 1345: 1, 1344: 1, 1338: 1, 1336: 1, 1335: 1, 1333: 1,
1331: 1, 1330: 1, 1329: 1, 1319: 1, 1318: 1, 1317: 1, 1315: 1, 1310: 1, 1309: 1, 1300: 1, 1296: 1,
1295: 1, 1292: 1, 1288: 1, 1284: 1, 1282: 1, 1280: 1, 1279: 1, 1271: 1, 1268: 1, 1263: 1, 1261: 1,
1259: 1, 1258: 1, 1257: 1, 1256: 1, 1254: 1, 1249: 1, 1247: 1, 1244: 1, 1242: 1, 1233: 1, 1230: 1,
1229: 1, 1228: 1, 1225: 1, 1222: 1, 1220: 1, 1217: 1, 1208: 1, 1207: 1, 1206: 1, 1204: 1, 1202: 1,
1199: 1, 1198: 1, 1197: 1, 1186: 1, 1185: 1, 1177: 1, 1174: 1, 1173: 1, 1171: 1, 1170: 1, 1169: 1,
1166: 1, 1162: 1, 1161: 1, 1159: 1, 1158: 1, 1149: 1, 1145: 1, 1144: 1, 1142: 1, 1138: 1, 1137: 1,
1134: 1, 1133: 1, 1130: 1, 1127: 1, 1125: 1, 1120: 1, 1119: 1, 1117: 1, 1116: 1, 1114: 1, 1111: 1,
1108: 1, 1107: 1, 1104: 1, 1101: 1, 1098: 1, 1094: 1, 1092: 1, 1090: 1, 1084: 1, 1082: 1, 1081: 1,
1078: 1, 1075: 1, 1068: 1, 1066: 1, 1063: 1, 1062: 1, 1061: 1, 1059: 1, 1051: 1, 1047: 1, 1043: 1,
1041: 1, 1037: 1, 1032: 1, 1028: 1, 1027: 1, 1025: 1, 1018: 1, 1013: 1, 1010: 1, 1009: 1, 1008: 1,
1007: 1, 1004: 1, 1000: 1, 997: 1, 991: 1, 988: 1, 987: 1, 980: 1, 979: 1, 967: 1, 964: 1, 963: 1,
962: 1, 961: 1, 958: 1, 955: 1, 952: 1, 943: 1, 931: 1, 928: 1, 924: 1, 921: 1, 916: 1, 912: 1,
906: 1, 905: 1, 897: 1, 893: 1, 891: 1, 890: 1, 882: 1, 878: 1, 877: 1, 875: 1, 869: 1, 865: 1,
863: 1, 862: 1, 856: 1, 850: 1, 848: 1, 844: 1, 839: 1, 835: 1, 834: 1, 832: 1, 817: 1, 810: 1,
809: 1, 803: 1, 798: 1, 789: 1, 785: 1, 774: 1, 767: 1, 760: 1, 753: 1, 752: 1, 748: 1, 747: 1,
745: 1, 736: 1, 731: 1, 728: 1, 720: 1, 719: 1, 714: 1, 704: 1, 701: 1, 699: 1, 667: 1, 664: 1,
653: 1, 650: 1, 648: 1, 642: 1, 640: 1, 636: 1, 634: 1, 624: 1, 621: 1, 613: 1, 586: 1, 585: 1,
544: 1, 518: 1, 512: 1, 441: 1})

```

In [48]:

```

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

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

```

```

#-----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_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)

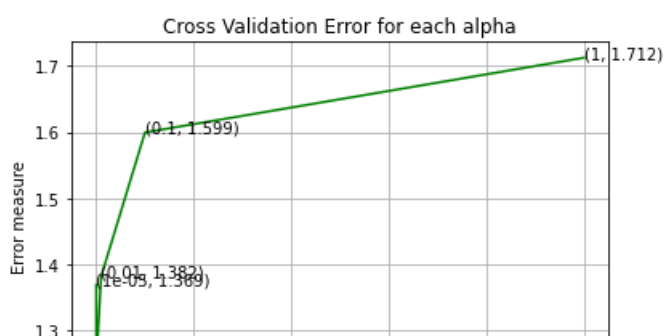
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best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

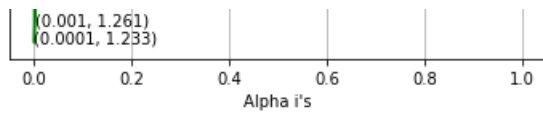
```

```

For values of alpha = 1e-05 The log loss is: 1.3691897299976423
For values of alpha = 0.0001 The log loss is: 1.2334552336725613
For values of alpha = 0.001 The log loss is: 1.2606914367659425
For values of alpha = 0.01 The log loss is: 1.3816721084507282
For values of alpha = 0.1 The log loss is: 1.5994369166054632
For values of alpha = 1 The log loss is: 1.7122477520941206

```





For values of best alpha = 0.0001 The train log loss is: 0.6997543386146584
 For values of best alpha = 0.0001 The cross validation log loss is: 1.2334552336725613
 For values of best alpha = 0.0001 The test log loss is: 1.0728654207979529

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

Ans. Yes, it seems like!

In [49]:

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

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

In [50]:

```
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")
```

93.258 % of word of test data appeared in train data
 95.176 % of word of Cross Validation appeared in train data

4. Machine Learning Models

In [51]:

```
#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 willl 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 [52]:

```
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 [58]:

```

# this function will be used just for naive bayes
# for the given indices, we will print the name of the features
# and we will check whether the feature present in the test point text or not
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = CountVectorizer()
    var_count_vec = CountVectorizer()
    text_count_vec = CountVectorizer(min_df=5)

    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")

```

In [87]:

```

gene_variation = []

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

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

```

Stacking the three types of features

In [90]:

```

tfidfVectorizer = TfidfVectorizer(max_features=3000)
text2 = tfidfVectorizer.fit_transform(gene_variation)
gene_variation_features = tfidfVectorizer.get_feature_names()

train_text = tfidfVectorizer.transform(train_df['TEXT'])
test_text = tfidfVectorizer.transform(test_df['TEXT'])
cv_text = tfidfVectorizer.transform(cv_df['TEXT'])

```

In [92]:

```

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

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

cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text))
cv_x_onehotCoding = hstack((cv_x_onehotCoding, cv_text_feature_onehotCoding)).tocsr()
cv_y = np.array(list(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 [93]:

```
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, 44595)
(number of data points * number of features) in test data = (665, 44595)
(number of data points * number of features) in cross validation data = (532, 44595)
```

In [94]:

```
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 [95]:

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

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

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

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

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

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

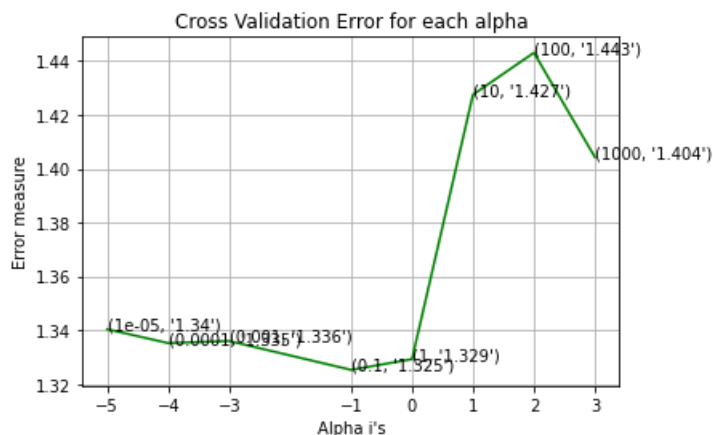
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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



```

for alpha = 1e-05
Log Loss : 1.3403425579459463
for alpha = 0.0001
Log Loss : 1.3352033611671004
for alpha = 0.001
Log Loss : 1.335969289198716
for alpha = 0.1
Log Loss : 1.325259351230726
for alpha = 1
Log Loss : 1.3292191823083355
for alpha = 10
Log Loss : 1.427483583692554
for alpha = 100
Log Loss : 1.4430691062504672
for alpha = 1000
Log Loss : 1.404294132563168

```



For values of best alpha = 0.1 The train log loss is: 0.8762052850971482
 For values of best alpha = 0.1 The cross validation log loss is: 1.325259351230726
 For values of best alpha = 0.1 The test log loss is: 1.207206816348683

4.1.1.2. Testing the model with best hyper paramters

In [96]:

```

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

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

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

clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)

```

```

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

```

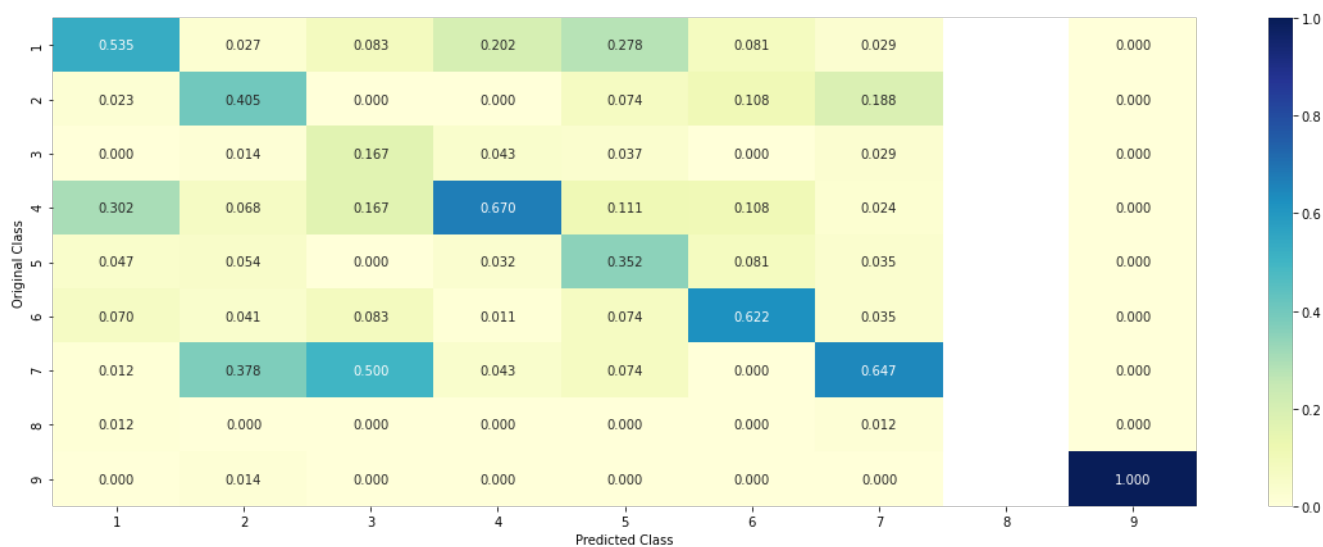
Log Loss : 1.325259351230726

Number of missclassified point : 0.4398496240601504

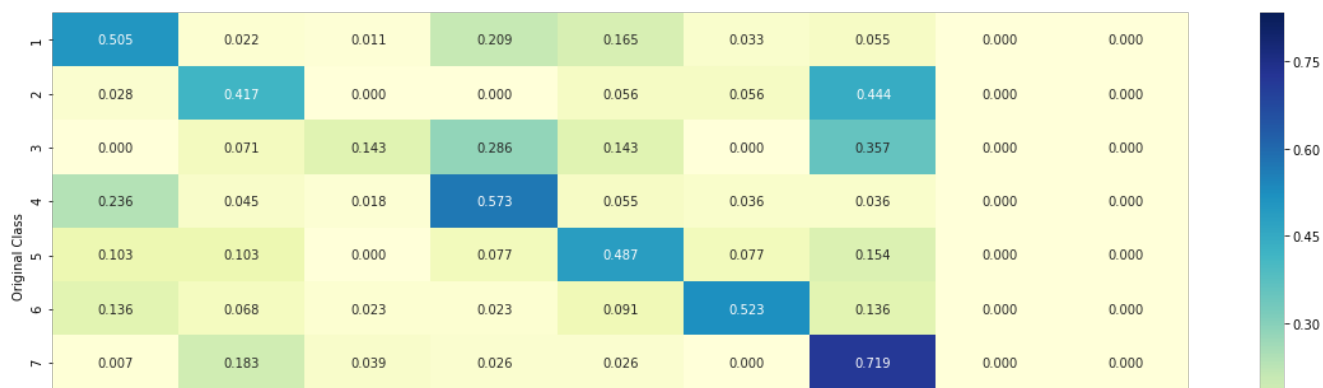
----- Confusion matrix -----

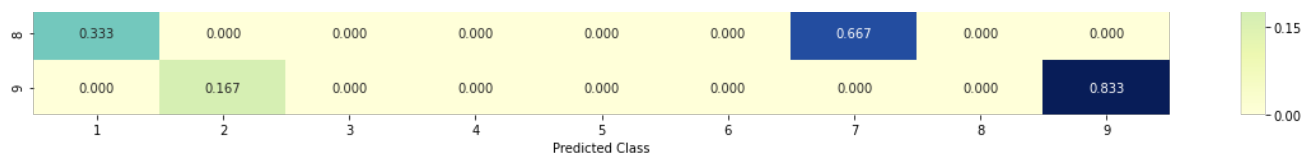


----- Precision matrix (Column Sum=1) -----



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





4.1.1.3. Feature Importance, Correctly classified point

In [97]:

```
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(-1*abs(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 : 5

Predicted Class Probabilities: [[0.2689 0.0841 0.0162 0.1307 0.3144 0.0455 0.1314 0.006 0.0029]]

Actual Class : 5

32 Text feature [methods] present in test data point [True]
38 Text feature [method] present in test data point [True]
Out of the top 100 features 2 are present in query point

4.1.1.4. Feature Importance, Incorrectly classified point

In [98]:

```
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(-1*abs(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: [[0.5678 0.0785 0.0149 0.1225 0.0439 0.0421 0.122 0.0056 0.0027]]

Actual Class : 1

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

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

In [67]:

```
# 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.appliedaiaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-geometric-intuition-with-a-toy-example-1/
#-----

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

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

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

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

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

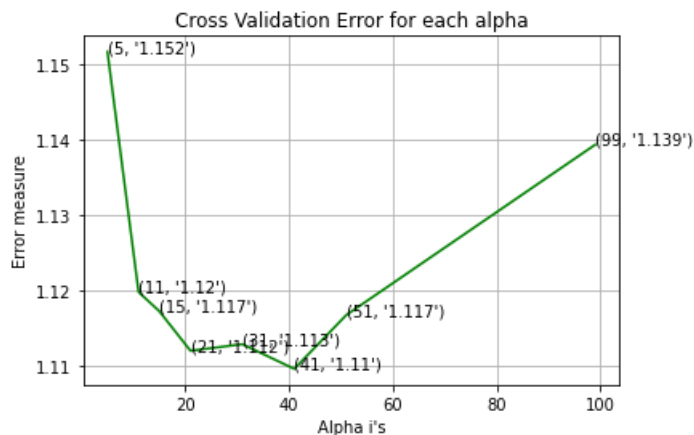
```

```

for alpha = 5
Log Loss : 1.1517546469107443
for alpha = 11
Log Loss : 1.1197933032523453
for alpha = 15
Log Loss : 1.1172319206531056
for alpha = 21
Log Loss : 1.1119742388648874
for alpha = 31
Log Loss : 1.1128493318544335

```

```
for alpha = 41
Log Loss : 1.109569252534094
for alpha = 51
Log Loss : 1.1167132437238476
for alpha = 99
Log Loss : 1.1393948588458642
```



```
For values of best alpha = 41 The train log loss is: 0.8515553363314864
For values of best alpha = 41 The cross validation log loss is: 1.109569252534094
For values of best alpha = 41 The test log loss is: 1.0865516932706785
```

4.2.2. Testing the model with best hyper paramters

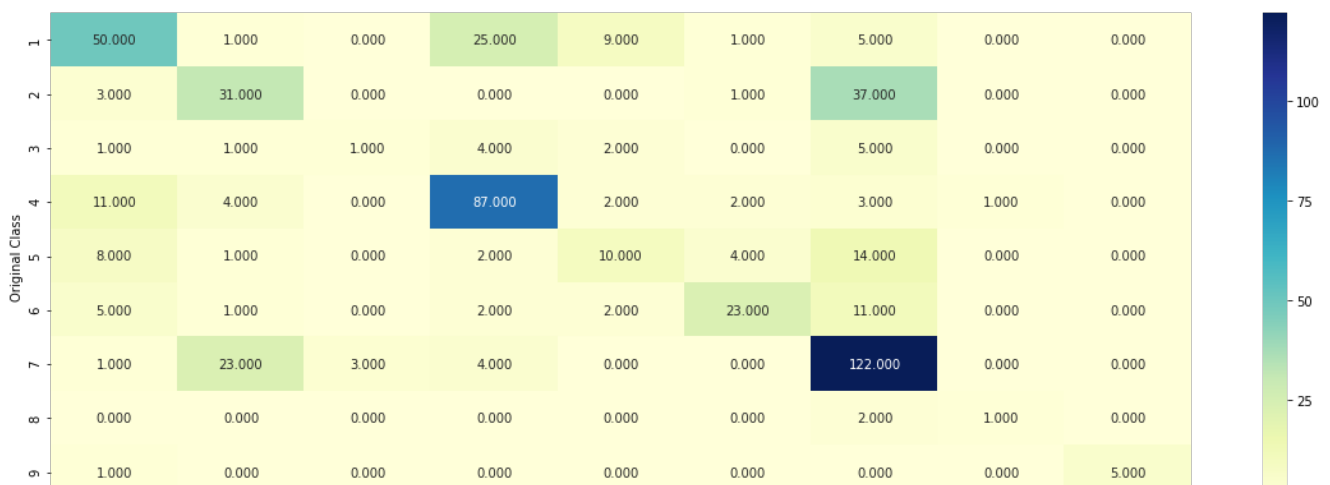
In [68]:

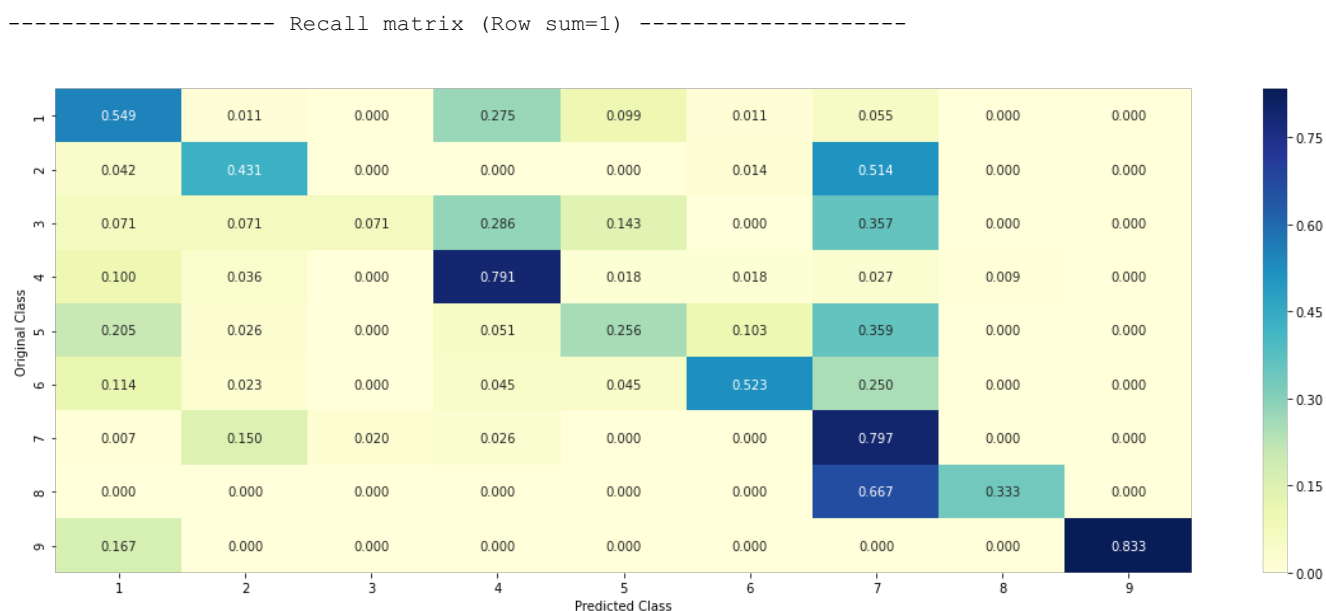
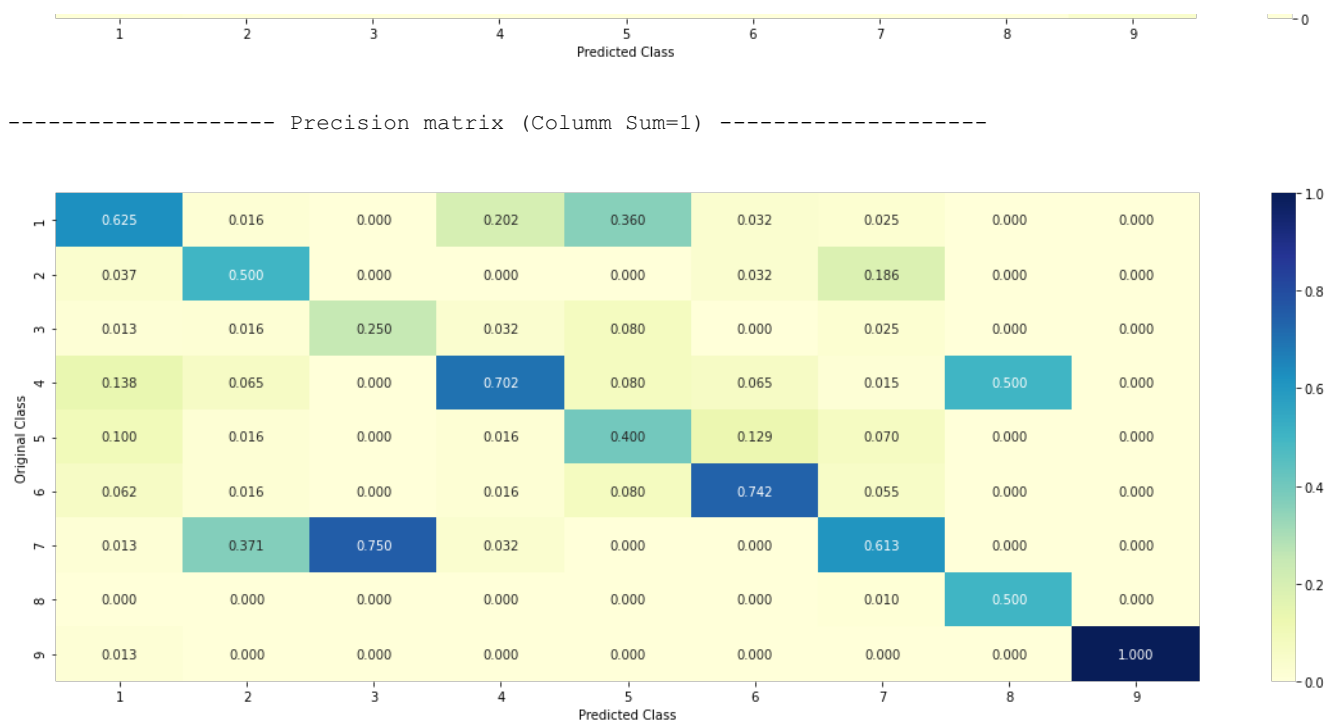
```
# find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None, 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 : 1.109569252534094
Number of mis-classified points : 0.37969924812030076
```

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





4.2.3. Sample Query point -1

In [69]:

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

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

Predicted Class : 7

Actual Class : 5

The 41 nearest neighbours of the test points belongs to classes [1 1 5 5 1 4 5 4 5 5 5 1 1 1 1 1 1 5 5 5 5 5 5 1 1 4 4 5 4 1 5 1 5 1 1 5

$$\begin{bmatrix} 1 & 1 & 1 & 6 \end{bmatrix}$$

4.2.4. Sample Query Point-2

In [70]:

```
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 belong to classes",train_y[neighbors[1][0]])
print("Frequency of nearest points :",Counter(train_y[neighbors[1][0]]))
```

[illegible]

4.3. Logistic Regression

4.3.1. With Class balancing

4.3.1.1. Hyper paramter tuning

In [99]:

```
# 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
```

```

# predict_proba() -> posterior probabilities of classification
#-----
# video link:
#-----

alpha = [10 ** x for x in range(-6, 3)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)

    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilities we use log-probability estimates
    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

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

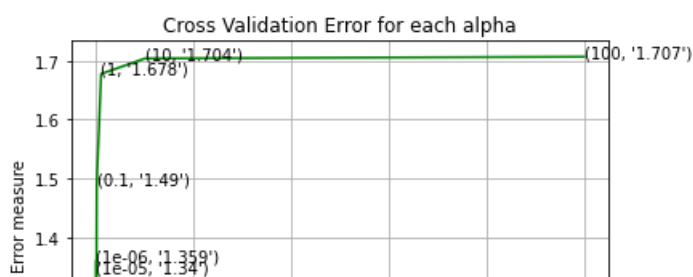
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

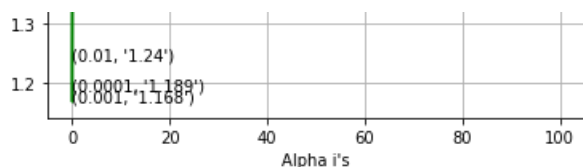
```

```

for alpha = 1e-06
Log Loss : 1.3594447834565797
for alpha = 1e-05
Log Loss : 1.3400679247846996
for alpha = 0.0001
Log Loss : 1.1886254498018929
for alpha = 0.001
Log Loss : 1.1679744423667313
for alpha = 0.01
Log Loss : 1.2402376917451046
for alpha = 0.1
Log Loss : 1.4899194255908614
for alpha = 1
Log Loss : 1.6775930620964115
for alpha = 10
Log Loss : 1.7038258700952287
for alpha = 100
Log Loss : 1.7067726556350924

```





For values of best alpha = 0.001 The train log loss is: 0.5375028601782635
 For values of best alpha = 0.001 The cross validation log loss is: 1.1679744423667313
 For values of best alpha = 0.001 The test log loss is: 0.9903096911558612

4.3.1.2. Testing the model with best hyper paramters

In [72]:

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

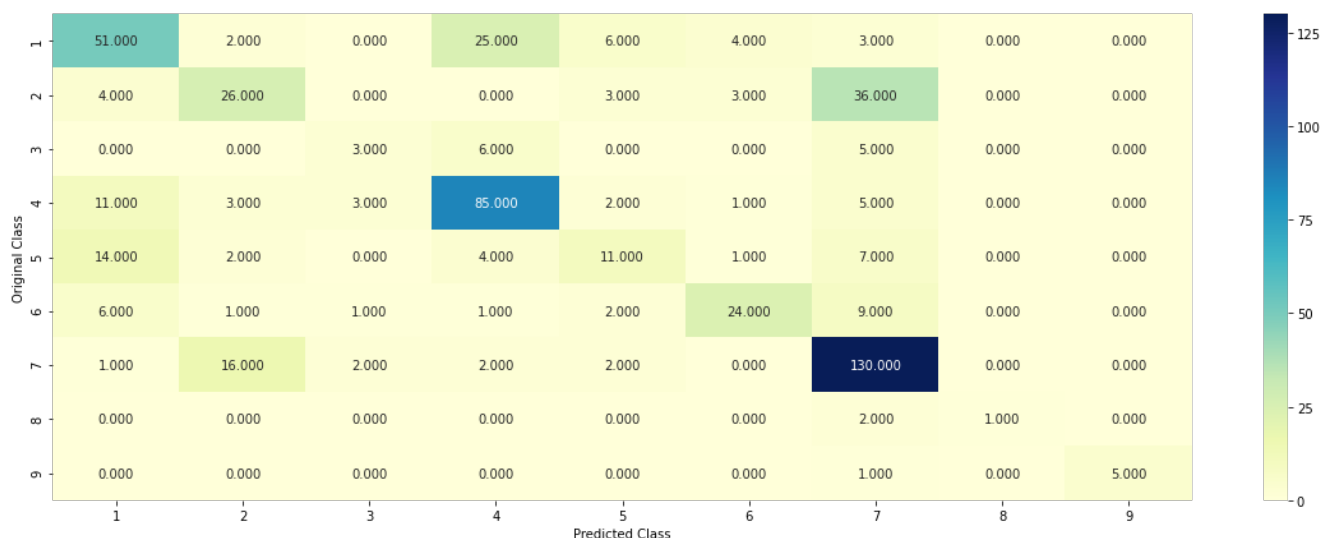
# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

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

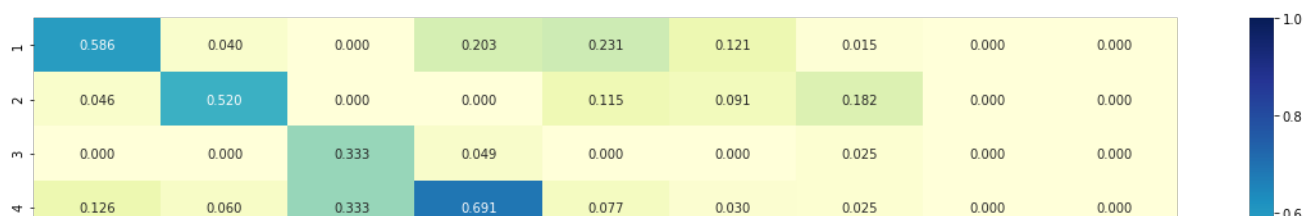
Log loss : 1.1845970019673246

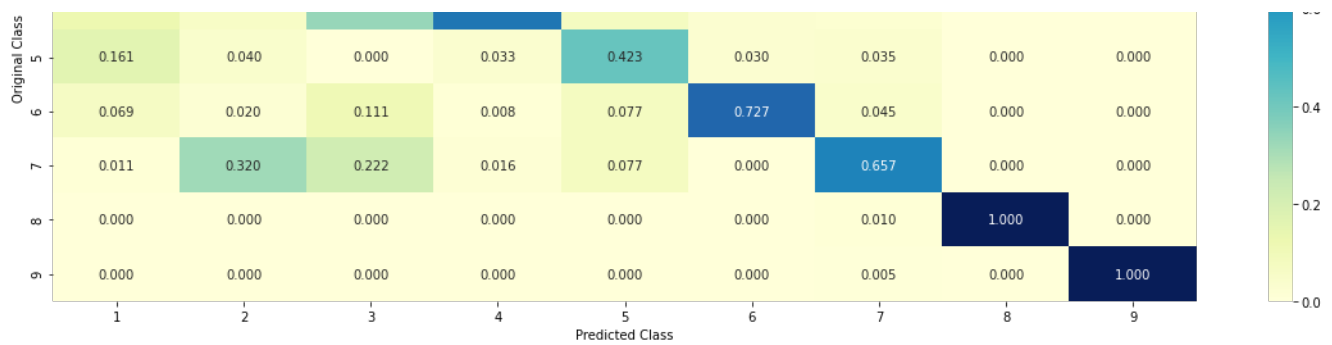
Number of mis-classified points : 0.3684210526315789

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

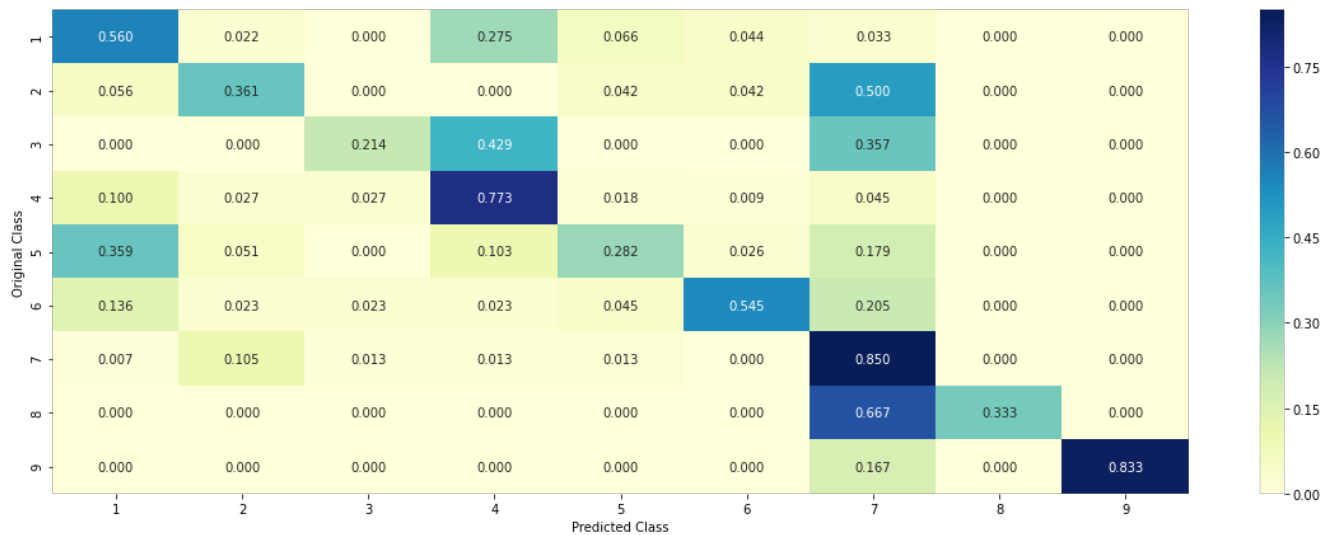


----- Precision matrix (Column Sum=1) -----





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



4.3.1.3. Feature Importance

In [73]:

```
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 [74]:

```
# from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
test_point_index = 1
no_feature = 500
```

```

predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(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: [[0.5014 0.0207 0.0038 0.2023 0.2419 0.0113 0.0108 0.0056 0.002 ]]
Actual Class : 5
-----
234 Text feature [kinase] present in test data point [True]
318 Text feature [nssnvs] present in test data point [True]
450 Text feature [mutants] present in test data point [True]
Out of the top 500 features 3 are present in query point

```

4.3.1.3.2. Incorrectly Classified point

In [75]:

```

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(-1*abs(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: [[0.5126 0.1061 0.0345 0.094 0.0707 0.0726 0.0911 0.0049 0.0134]]
Actual Class : 1
-----
450 Text feature [mutants] present in test data point [True]
Out of the top 500 features 1 are present in query point

```

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

In [76]:

```

# 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.appliedaicaourse.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)

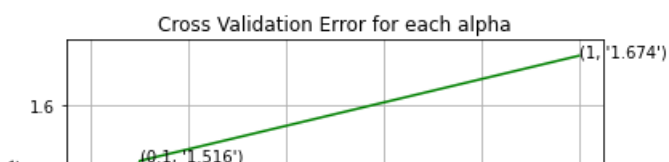
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

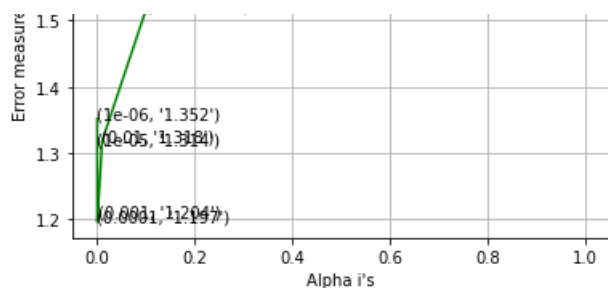
```

```

for alpha = 1e-06
Log Loss : 1.3516860591083677
for alpha = 1e-05
Log Loss : 1.3136250353377172
for alpha = 0.0001
Log Loss : 1.1965201668847973
for alpha = 0.001
Log Loss : 1.203971633795541
for alpha = 0.01
Log Loss : 1.3181847085943195
for alpha = 0.1
Log Loss : 1.5158252384544342
for alpha = 1
Log Loss : 1.6740111224561578

```





For values of best alpha = 0.0001 The train log loss is: 0.5248209573595253
 For values of best alpha = 0.0001 The cross validation log loss is: 1.1965201668847973
 For values of best alpha = 0.0001 The test log loss is: 1.0155795161035606

4.3.2.2. Testing model with best hyper parameters

In [77]:

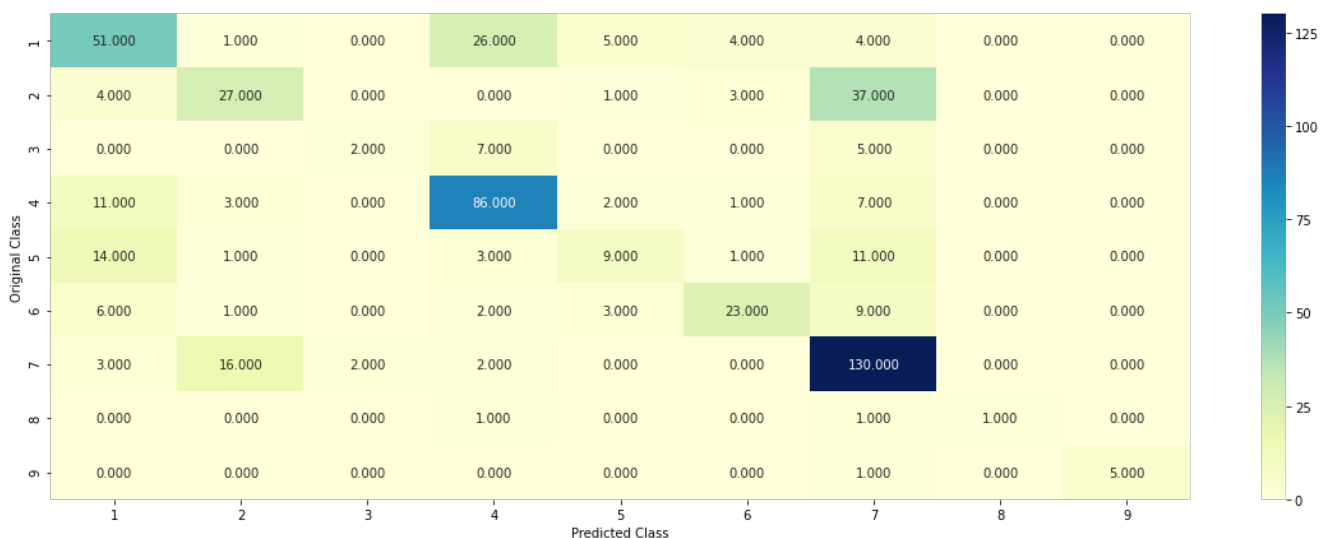
```
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, 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.

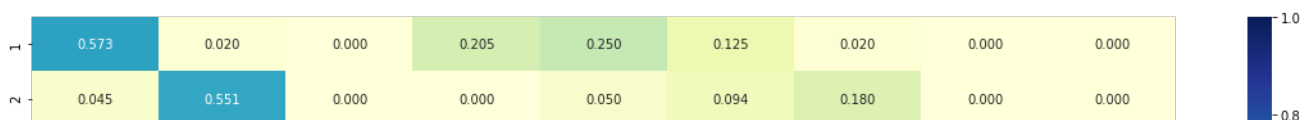
#-----
# video link:
#-----

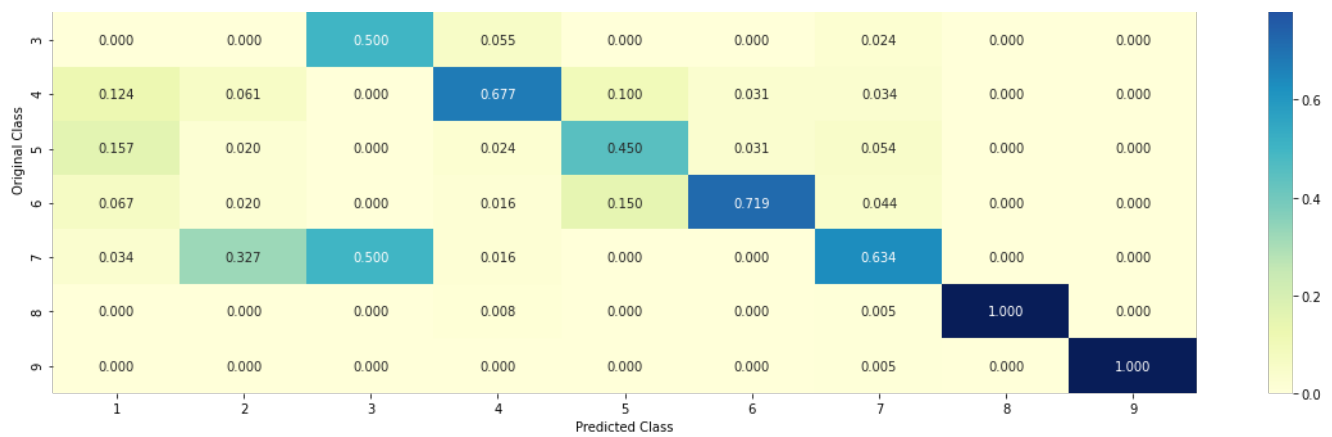
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```

Log loss : 1.1965201668847973
 Number of mis-classified points : 0.37218045112781956
 ----- Confusion matrix -----

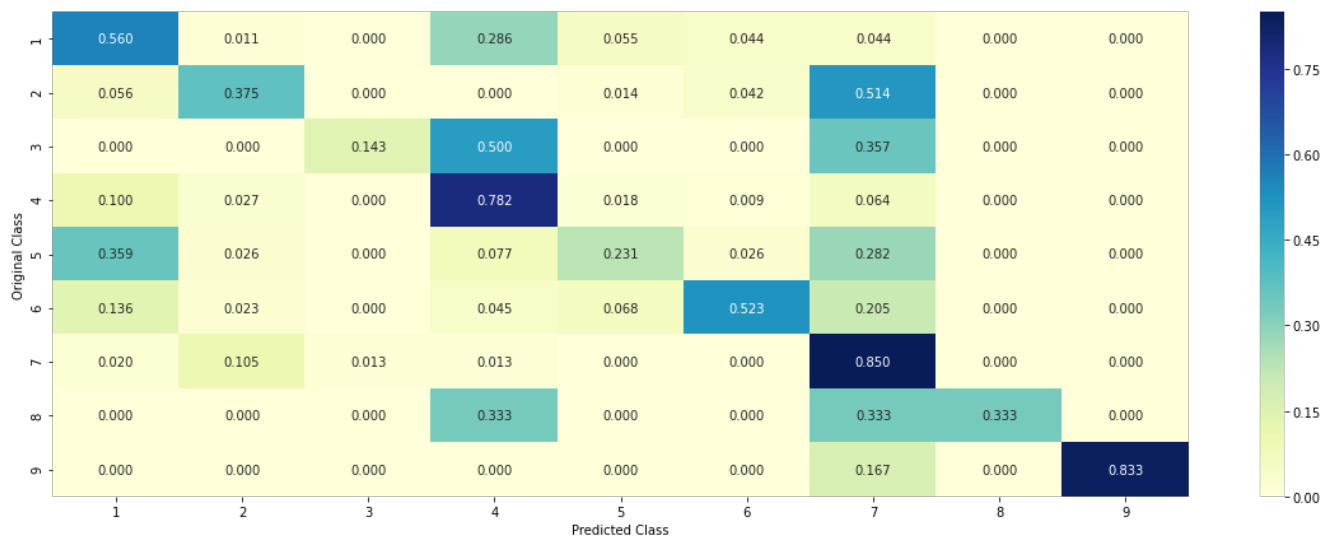


----- Precision matrix (Column Sum=1) -----





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



4.3.2.3. Feature Importance, Correctly Classified point

In [78]:

```
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(-1*abs(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: [[0.4336 0.0365 0.0031 0.3528 0.1236 0.0053 0.0394 0.0041 0.0016]]
Actual Class : 5

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

4.3.2.4. Feature Importance, Incorrectly Classified point

In [79]:

```
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(-1*abs(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: [[0.5059 0.1066 0.0547 0.0664 0.0744 0.0812 0.0831 0.0066 0.021 ]]
Actual Class : 1
-----
Out of the top 500 features 0 are present in query point

```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

In [80]:

```

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

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

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

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

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

fig, ax = plt.subplots()

```

```

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

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

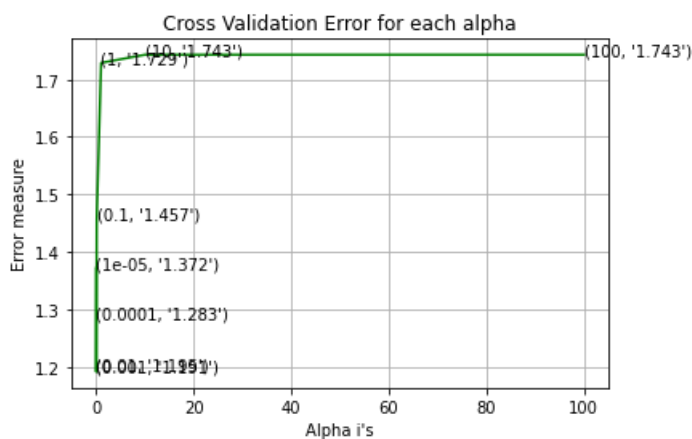
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "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 C = 1e-05
Log Loss : 1.3716000268317061
for C = 0.0001
Log Loss : 1.2833215682997223
for C = 0.001
Log Loss : 1.1910205698377403
for C = 0.01
Log Loss : 1.1949566733654677
for C = 0.1
Log Loss : 1.4572755288019748
for C = 1
Log Loss : 1.7289869906130009
for C = 10
Log Loss : 1.7431263041838327
for C = 100
Log Loss : 1.7431262667026213

```



```

For values of best alpha = 0.001 The train log loss is: 0.5406640799704894
For values of best alpha = 0.001 The cross validation log loss is: 1.1910205698377403
For values of best alpha = 0.001 The test log loss is: 1.048853487671113

```

4.4.2. Testing model with best hyper parameters

In [81]:

```

# 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)

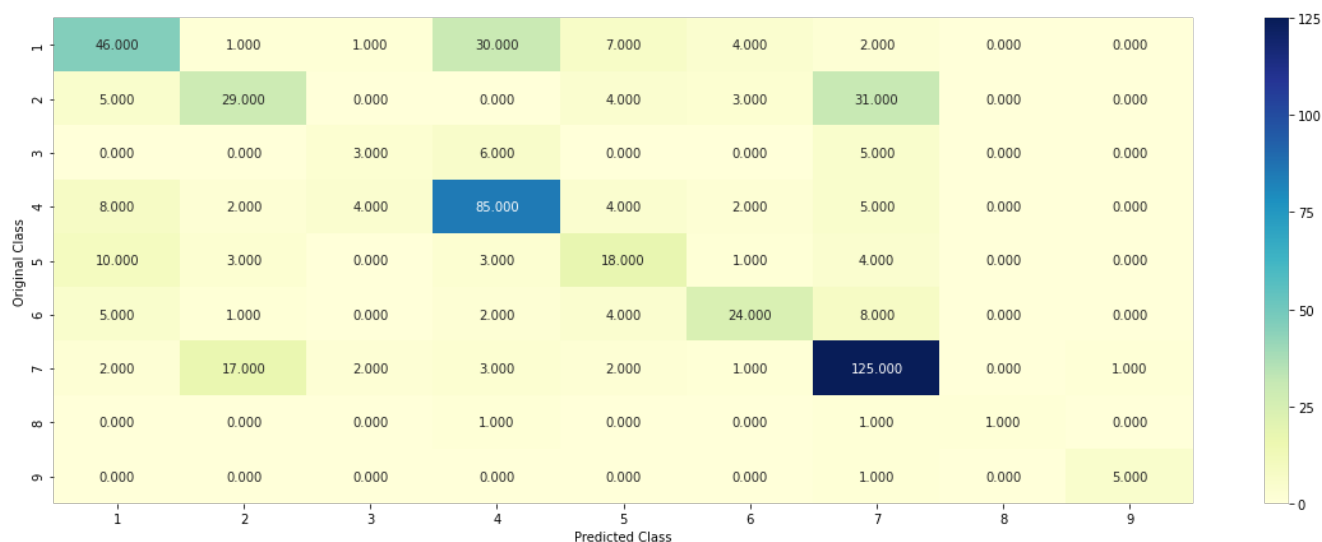
# Some of methods of SVM()
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
# predict(X) Perform classification on samples in X.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

# 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)
```

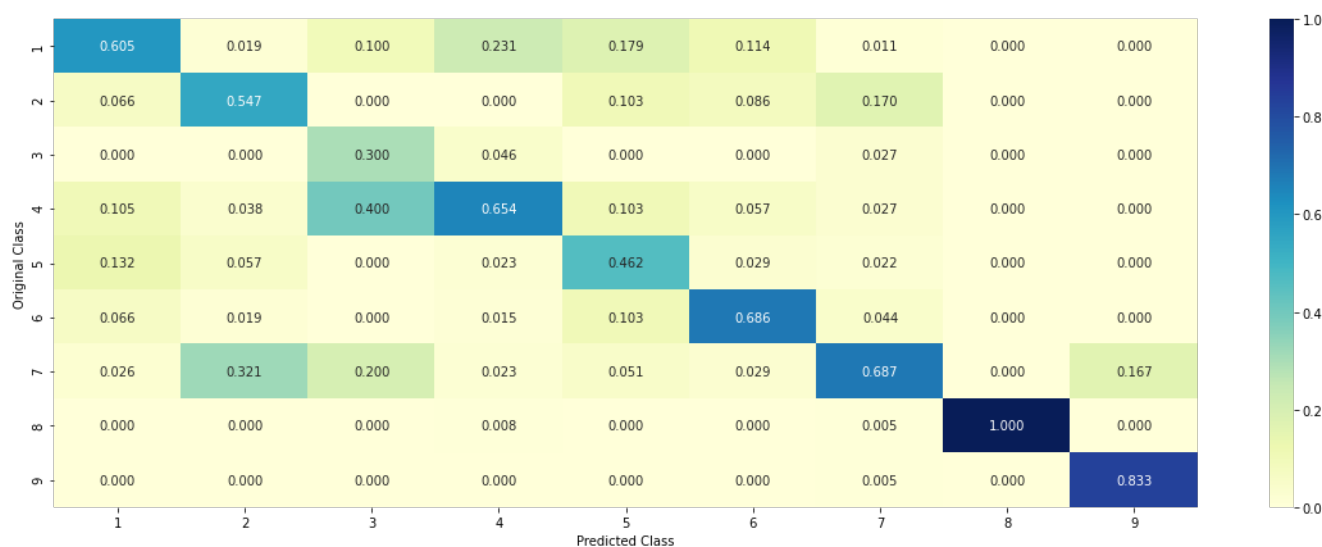
Log loss : 1.1910205698377403

Number of mis-classified points : 0.3684210526315789

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

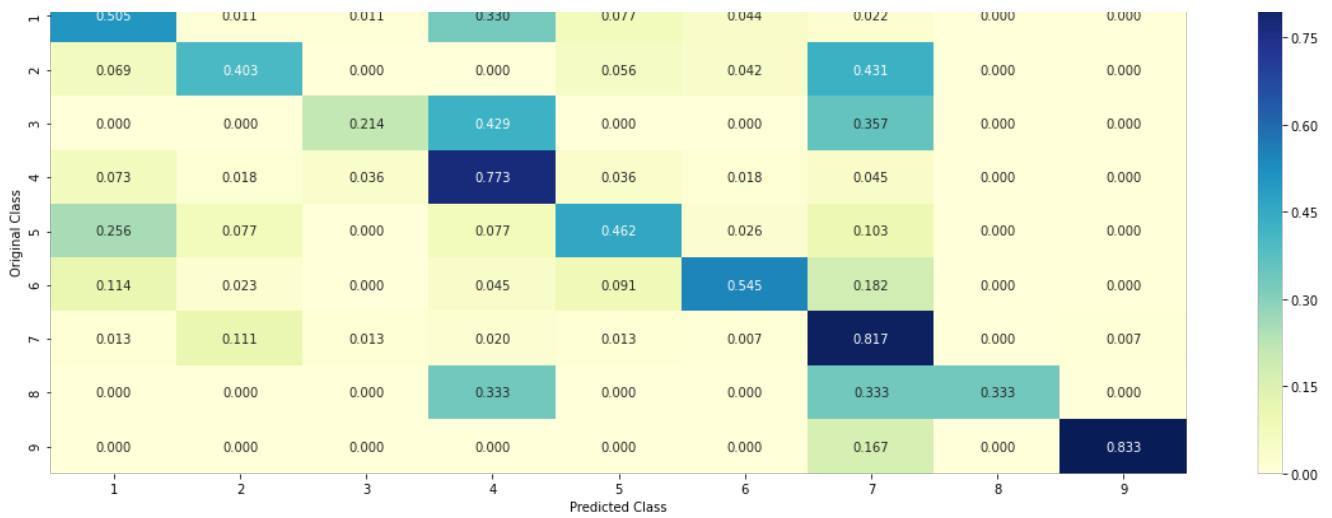


----- Precision matrix (Column Sum=1) -----



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





4.3.3. Feature Importance

4.3.3.1. For Correctly classified point

In [82]:

```
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(-1*abs(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: [[0.3501 0.0619 0.0044 0.2828 0.2021 0.0042 0.0837 0.0059 0.005]]
Actual Class : 5

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

4.3.3.2. For Incorrectly classified point

In [83]:

```
test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(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: [[0.595 0.0932 0.0289 0.0623 0.067 0.0602 0.0808 0.0071 0.0055]]
Actual Class : 1

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

4.5 Random Forest Classifier

4.5.1. Hyper paramter tuning (With One hot Encoding)

In [84]:

```
# -----
# 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-fores
t-and-their-construction-2/
# -----

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

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

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

```

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

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The train log loss
is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The cross validation log loss
is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The test log loss
is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for n_estimators = 100 and max depth = 5
Log Loss : 1.2621646443714885
for n_estimators = 100 and max depth = 10
Log Loss : 1.1882401285490771
for n_estimators = 200 and max depth = 5
Log Loss : 1.2443776365057888
for n_estimators = 200 and max depth = 10
Log Loss : 1.186690381521011
for n_estimators = 500 and max depth = 5
Log Loss : 1.231400332204527
for n_estimators = 500 and max depth = 10
Log Loss : 1.1778291021497949
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2304067512090928
for n_estimators = 1000 and max depth = 10
Log Loss : 1.1758123003369987
for n_estimators = 2000 and max depth = 5
Log Loss : 1.2284175876925612
for n_estimators = 2000 and max depth = 10
Log Loss : 1.1764943526313882
For values of best estimator = 1000 The train log loss is: 0.7051681416517144
For values of best estimator = 1000 The cross validation log loss is: 1.1758123003369987
For values of best estimator = 1000 The test log loss is: 1.1119921233852879

```

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

In [84]:

```

# -----
# 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/
# -----

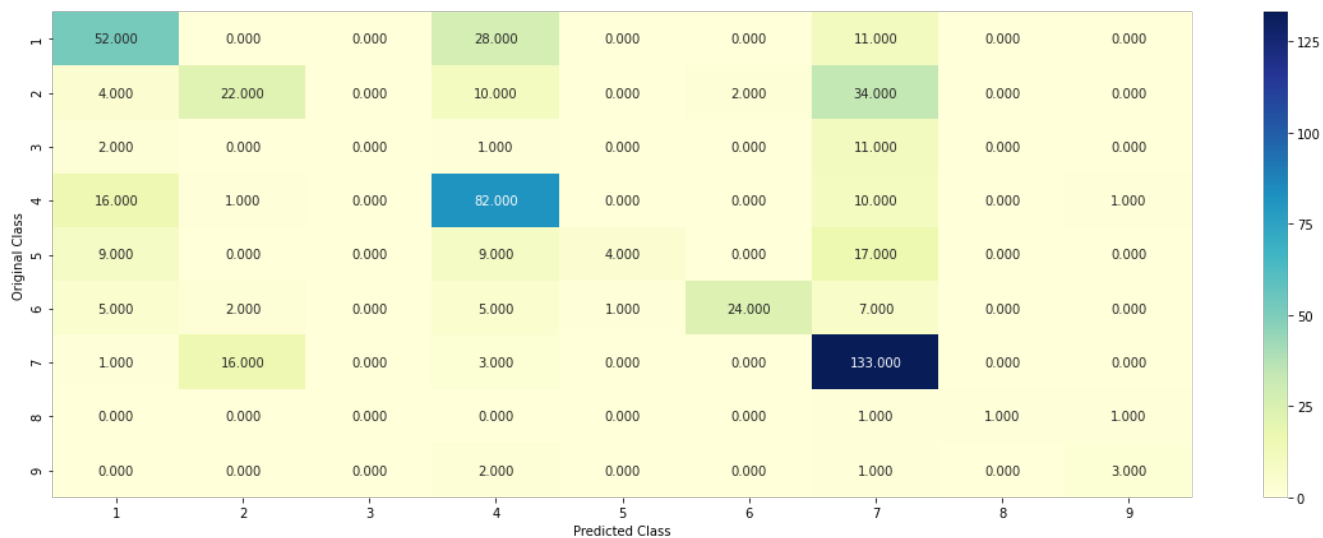
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max
_depth[int(best_alpha%2)], random_state=42, n_jobs=-1)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

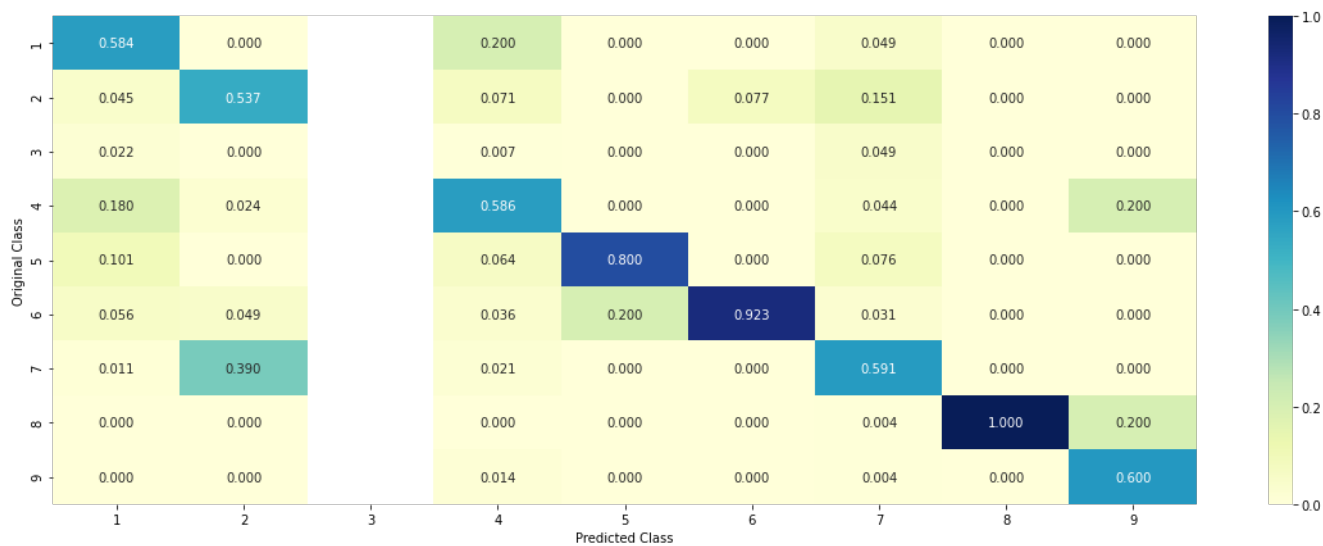
Log loss : 1.1899619916532993

Number of mis-classified points : 0.3966165413533835

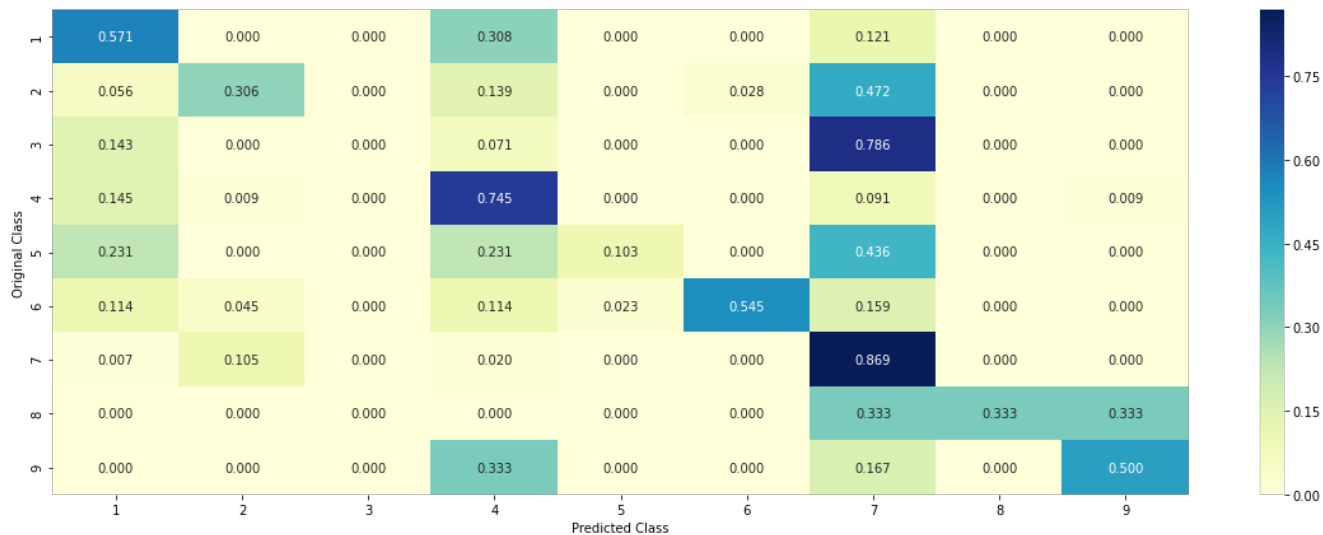
----- Confusion matrix -----



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



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



4.5.3. Feature Importance

4.5.3.1. Correctly Classified point

In [85]:

```
# 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: [[0.2066 0.0347 0.0181 0.5993 0.0462 0.0395 0.0416 0.0056 0.0084]]

Actual Class : 4

```
-----
6 Text feature [inhibitor] present in test data point [True]
7 Text feature [activation] present in test data point [True]
8 Text feature [missense] present in test data point [True]
12 Text feature [nonsense] present in test data point [True]
14 Text feature [function] present in test data point [True]
19 Text feature [functional] present in test data point [True]
28 Text feature [deleterious] present in test data point [True]
29 Text feature [cells] present in test data point [True]
30 Text feature [treated] present in test data point [True]
31 Text feature [expressing] present in test data point [True]
37 Text feature [frameshift] present in test data point [True]
38 Text feature [loss] present in test data point [True]
39 Text feature [months] present in test data point [True]
44 Text feature [pathogenic] present in test data point [True]
46 Text feature [protein] present in test data point [True]
47 Text feature [cell] present in test data point [True]
48 Text feature [proteins] present in test data point [True]
49 Text feature [unstable] present in test data point [True]
51 Text feature [patients] present in test data point [True]
52 Text feature [extracellular] present in test data point [True]
56 Text feature [drug] present in test data point [True]
61 Text feature [variants] present in test data point [True]
63 Text feature [stability] present in test data point [True]
68 Text feature [variant] present in test data point [True]
77 Text feature [serum] present in test data point [True]
79 Text feature [lines] present in test data point [True]
84 Text feature [nuclear] present in test data point [True]
86 Text feature [truncating] present in test data point [True]
91 Text feature [expression] present in test data point [True]
92 Text feature [ligand] present in test data point [True]
94 Text feature [clinical] present in test data point [True]
Out of the top 100 features 31 are present in query point
```

4.5.3.2. Inorrectly Classified point

In [86]:

```
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])
```

```
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_impfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_index], test_df['Variation'].iloc[test_point_index], no_feature)
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0822 0.1516 0.0209 0.0712 0.0465 0.0381 0.5734 0.006 0.01]]

Actual Class : 7

```
-----
0 Text feature [kinase] present in test data point [True]
1 Text feature [activating] present in test data point [True]
3 Text feature [inhibitors] present in test data point [True]
4 Text feature [phosphorylation] present in test data point [True]
5 Text feature [activated] present in test data point [True]
6 Text feature [inhibitor] present in test data point [True]
7 Text feature [activation] present in test data point [True]
8 Text feature [missense] present in test data point [True]
9 Text feature [erk] present in test data point [True]
10 Text feature [suppressor] present in test data point [True]
11 Text feature [constitutive] present in test data point [True]
13 Text feature [signaling] present in test data point [True]
14 Text feature [function] present in test data point [True]
15 Text feature [akt] present in test data point [True]
16 Text feature [kinases] present in test data point [True]
17 Text feature [oncogenic] present in test data point [True]
18 Text feature [treatment] present in test data point [True]
19 Text feature [functional] present in test data point [True]
20 Text feature [therapeutic] present in test data point [True]
25 Text feature [downstream] present in test data point [True]
26 Text feature [growth] present in test data point [True]
29 Text feature [cells] present in test data point [True]
30 Text feature [treated] present in test data point [True]
31 Text feature [expressing] present in test data point [True]
32 Text feature [trials] present in test data point [True]
34 Text feature [therapy] present in test data point [True]
35 Text feature [transforming] present in test data point [True]
36 Text feature [activate] present in test data point [True]
38 Text feature [loss] present in test data point [True]
42 Text feature [inhibition] present in test data point [True]
46 Text feature [protein] present in test data point [True]
47 Text feature [cell] present in test data point [True]
48 Text feature [proteins] present in test data point [True]
51 Text feature [patients] present in test data point [True]
56 Text feature [drug] present in test data point [True]
57 Text feature [resistance] present in test data point [True]
59 Text feature [ras] present in test data point [True]
61 Text feature [variants] present in test data point [True]
66 Text feature [mek] present in test data point [True]
70 Text feature [autophosphorylation] present in test data point [True]
72 Text feature [mapk] present in test data point [True]
77 Text feature [serum] present in test data point [True]
78 Text feature [response] present in test data point [True]
79 Text feature [lines] present in test data point [True]
80 Text feature [3t3] present in test data point [True]
85 Text feature [inhibited] present in test data point [True]
87 Text feature [favor] present in test data point [True]
91 Text feature [expression] present in test data point [True]
94 Text feature [clinical] present in test data point [True]
96 Text feature [proliferation] present in test data point [True]
98 Text feature [survival] present in test data point [True]
Out of the top 100 features 51 are present in query point
```

4.5.3. Hyper paramter tuning (With Response Coding)

In [87]:

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

```

verbose=0, warm_start=False,
# class_weight=None)

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

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

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

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

alpha = [10,50,100,200,500,1000]
max_depth = [2,3,5,10]
cv_log_error_array = []
for i in alpha:
    for j in max_depth:
        print("for n_estimators = ", i,"and max depth = ", j)
        clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
        clf.fit(train_x_responseCoding, train_y)
        sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
        sig_clf.fit(train_x_responseCoding, train_y)
        sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
        cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
        print("Log Loss :",log_loss(cv_y, sig_clf_probs))
'''
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 : 2.1057339861176305
for n_estimators = 10 and max depth = 3
Log Loss : 1.5871212464988946
for n_estimators = 10 and max depth = 5
Log Loss : 1.5039719610007254
for n_estimators = 10 and max depth = 10
Log Loss : 2.0390188712265336
for n_estimators = 50 and max depth = 2
Log Loss : 1.8265828024876034
for n_estimators = 50 and max depth = 3
Log Loss : 1.5099360897726288
for n_estimators = 50 and max depth = 5
Log Loss : 1.4421202029845321
for n_estimators = 50 and max depth = 10
Log Loss : 1.8329407511442466
for n_estimators = 100 and max depth = 2
Log Loss : 1.6532009919902344
for n_estimators = 100 and max depth = 3
Log Loss : 1.5427867315412678
for n_estimators = 100 and max depth = 5
Log Loss : 1.3985713366119918
for n_estimators = 100 and max depth = 10
Log Loss : 1.7898425221792906
for n_estimators = 200 and max depth = 2
Log Loss : 1.717120764036802
for n_estimators = 200 and max depth = 3
Log Loss : 1.4919021188415809
for n_estimators = 200 and max depth = 5
Log Loss : 1.4409783661920599
for n_estimators = 200 and max depth = 10
Log Loss : 1.8005150313744969
for n_estimators = 500 and max depth = 2
Log Loss : 1.7827061236641906
for n_estimators = 500 and max depth = 3
Log Loss : 1.574352304933845
for n_estimators = 500 and max depth = 5
Log Loss : 1.438075885130304
for n_estimators = 500 and max depth = 10
Log Loss : 1.8231343269493834
for n_estimators = 1000 and max depth = 2
Log Loss : 1.7555803972526438
for n_estimators = 1000 and max depth = 3
Log Loss : 1.5803078817896814
for n_estimators = 1000 and max depth = 5
Log Loss : 1.4323125305218847
for n_estimators = 1000 and max depth = 10
Log Loss : 1.8224943498080939
For values of best alpha = 100 The train log loss is: 0.0498442408823619
For values of best alpha = 100 The cross validation log loss is: 1.3985713366119916
For values of best alpha = 100 The test log loss is: 1.394602663940248

```

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

In [88]:

```

# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_s
amples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min
impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None,
verbose=0, warm_start=False,
# class_weight=None)

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

# some of attributes of RandomForestClassifier()
# feature_importances_ : array of shape = [n_features]

```

```
# The feature importances (the higher, the more important the feature).

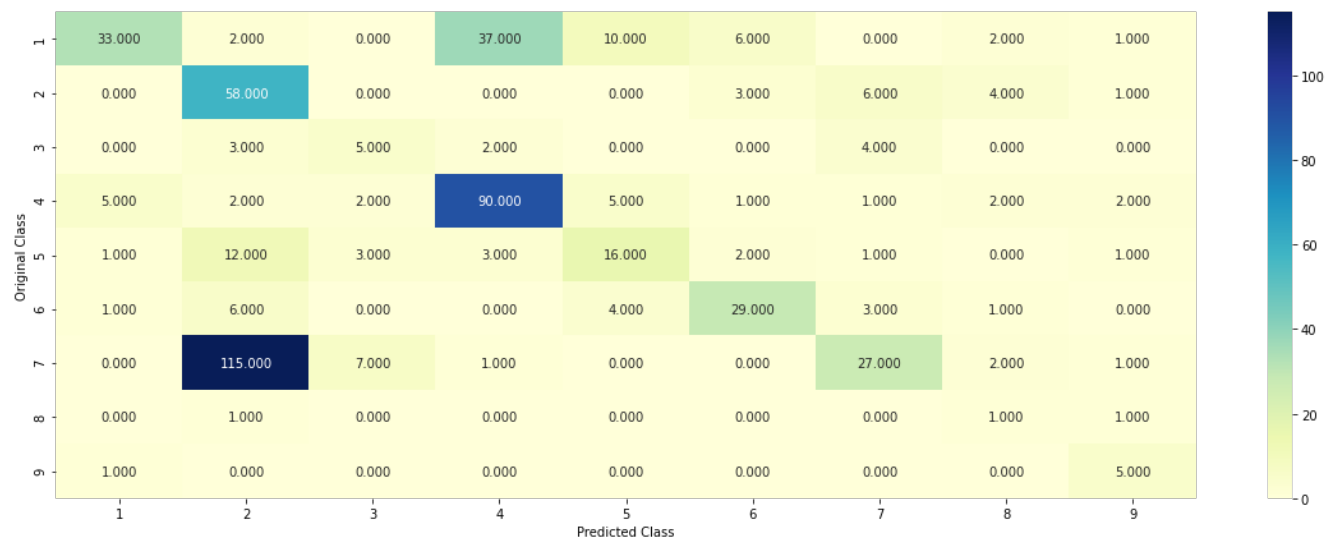
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

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)
```

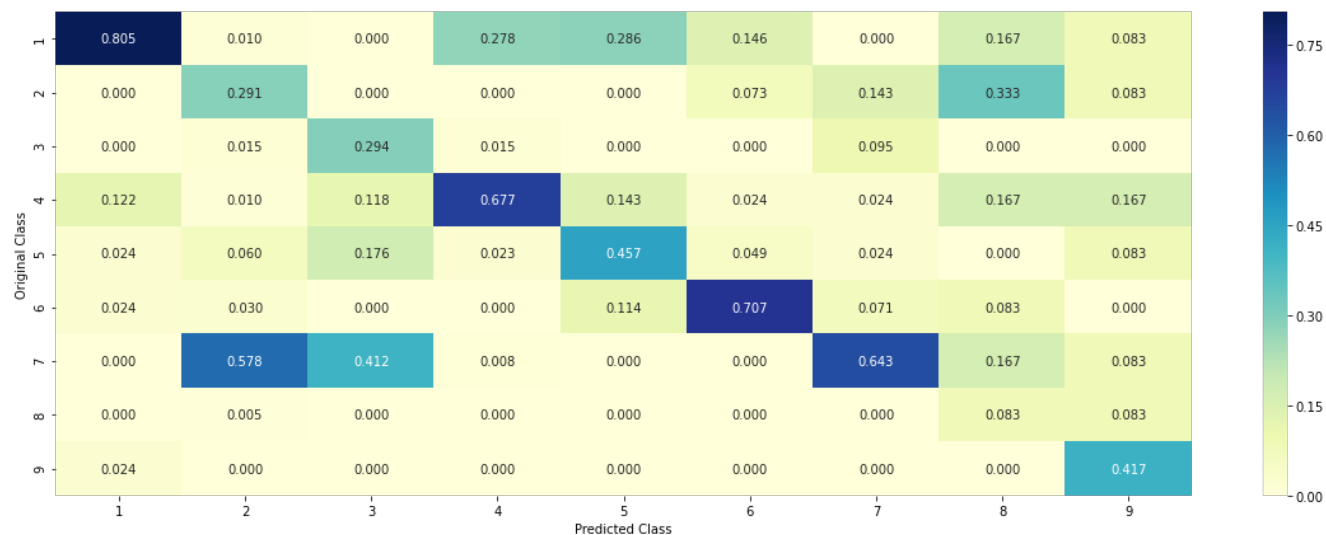
Log loss : 1.3985713366119918

Number of mis-classified points : 0.5037593984962406

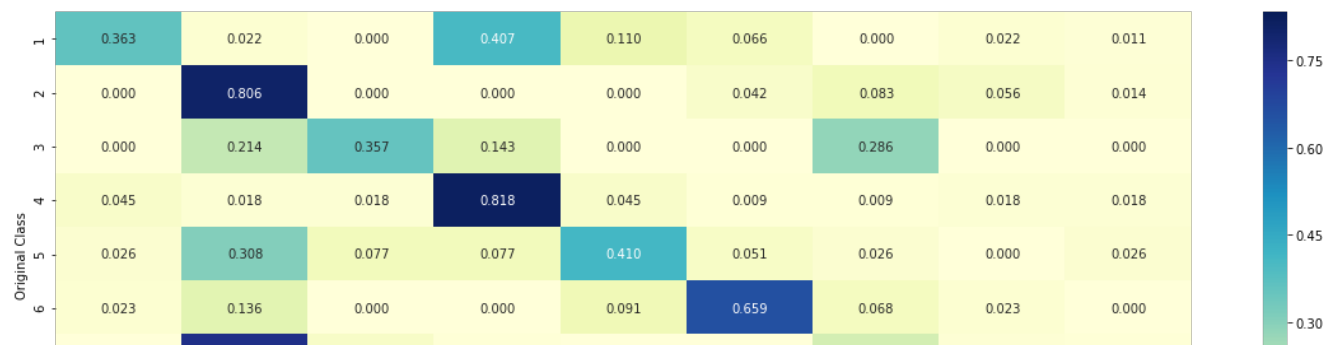
----- Confusion matrix -----

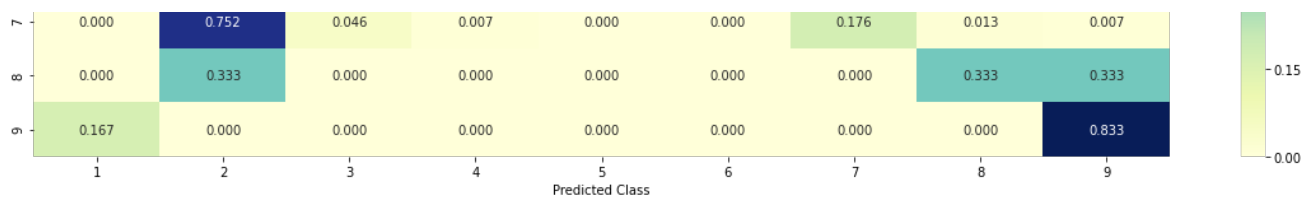


----- Precision matrix (Column Sum=1) -----



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





4.5.5. Feature Importance

4.5.5.1. Correctly Classified point

In [89]:

```
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: [[0.0821 0.0352 0.152  0.5299 0.0288 0.0398 0.0075 0.0517 0.0731]]
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
Text is important feature
Gene is important feature
Gene is important feature
Text is important feature
Variation is important feature
Text is important feature
Variation is important feature
Text is important feature
Gene is important feature
Gene is important feature
Gene is important feature
```

4.5.5.2. Incorrectly Classified point

4.6.2. Incorrectly classified point

In [90]:

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

Predicted Class : 7
Predicted Class Probabilities: [[0.025 0.2817 0.175 0.0305 0.038 0.0662 0.2955 0.0291 0.059]]
Actual Class : 7

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

4.7 Stack the models

4.7.1 testing with hyper parameter tuning

In [91]:

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

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

#-----
```

```

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

# read more about support vector machines with linear kernels here http://scikit-
learn.org/stable/modules/generated/sklearn.svm.SVC.html
# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, 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/
# -----

# 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_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-fores
t-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_onehot
Coding))))
sig_clf2.fit(train_x_onehotCoding, train_y)
print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y,
sig_clf2.predict_proba(cv_x_onehotCoding))))
sig_clf3.fit(train_x_onehotCoding, train_y)
print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predict_proba(cv_x_onehotCoding)))
)
print("-"*50)
alpha = [0.0001,0.001,0.01,0.1,1,10]
best_alpha = 999
for i in alpha:

```

```

lr = LogisticRegression(C=i)
sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_p
robas=True)
sclf.fit(train_x_onehotCoding, train_y)
print("Stacking Classifier : for the value of alpha: %f Log Loss: %0.3f" % (i, log_loss(cv_y, sc
lf.predict_proba(cv_x_onehotCoding))))
log_error = log_loss(cv_y, sclf.predict_proba(cv_x_onehotCoding))
if best_alpha > log_error:
    best_alpha = log_error

```

Logistic Regression : Log Loss: 1.12
Support vector machines : Log Loss: 1.73
Naive Bayes : Log Loss: 1.27

Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.178
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.036
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.509
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.104
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.175
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.419

4.7.2 testing the model with the best hyper parameters

In [92]:

```

lr = LogisticRegression(C=0.1)
sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, use_proba
s=True)
sclf.fit(train_x_onehotCoding, train_y)

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

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

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

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

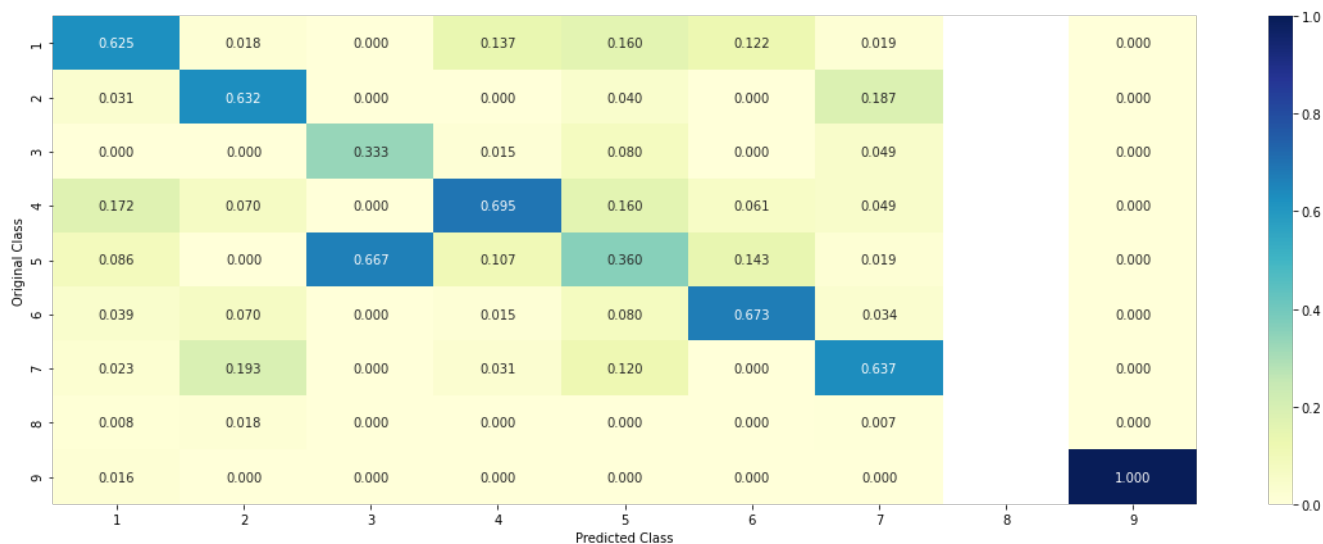
```

Log loss (train) on the stacking classifier : 0.6340791778015482
Log loss (CV) on the stacking classifier : 1.1040946954947586
Log loss (test) on the stacking classifier : 1.144718441629678
Number of missclassified point : 0.3609022556390977

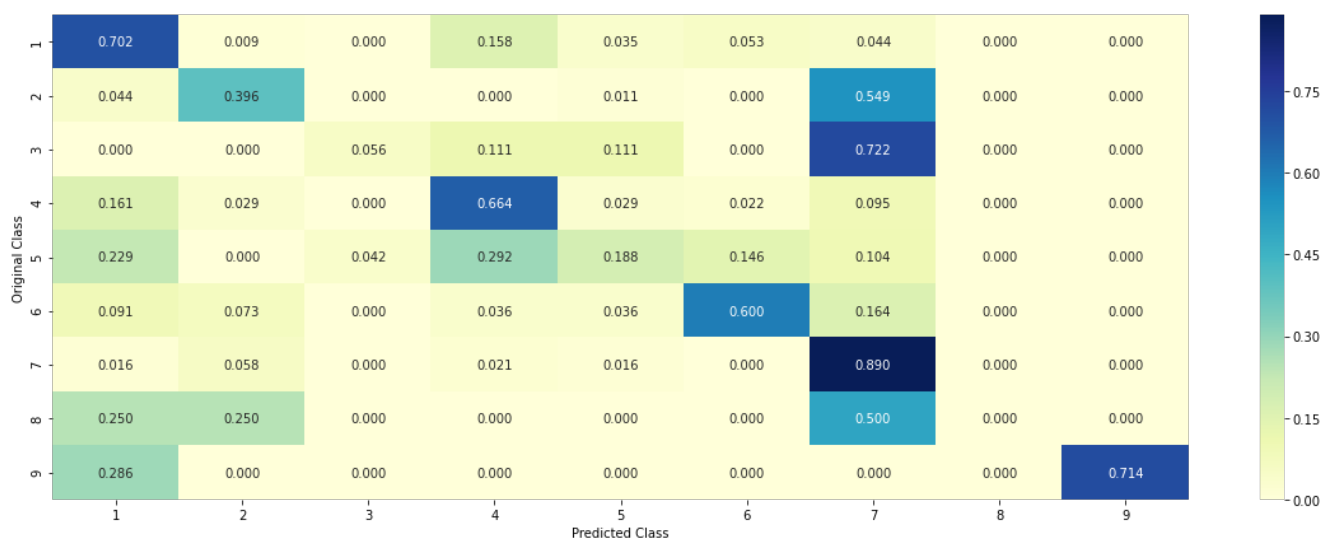
----- Confusion matrix -----



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



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



4.7.3 Maximum Voting classifier

In [93]:

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

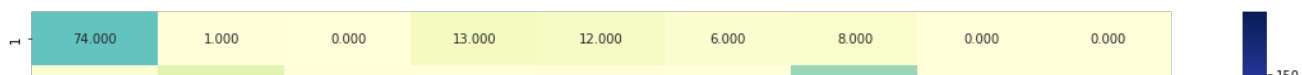
Log loss (train) on the VotingClassifier : 0.8847096843900962

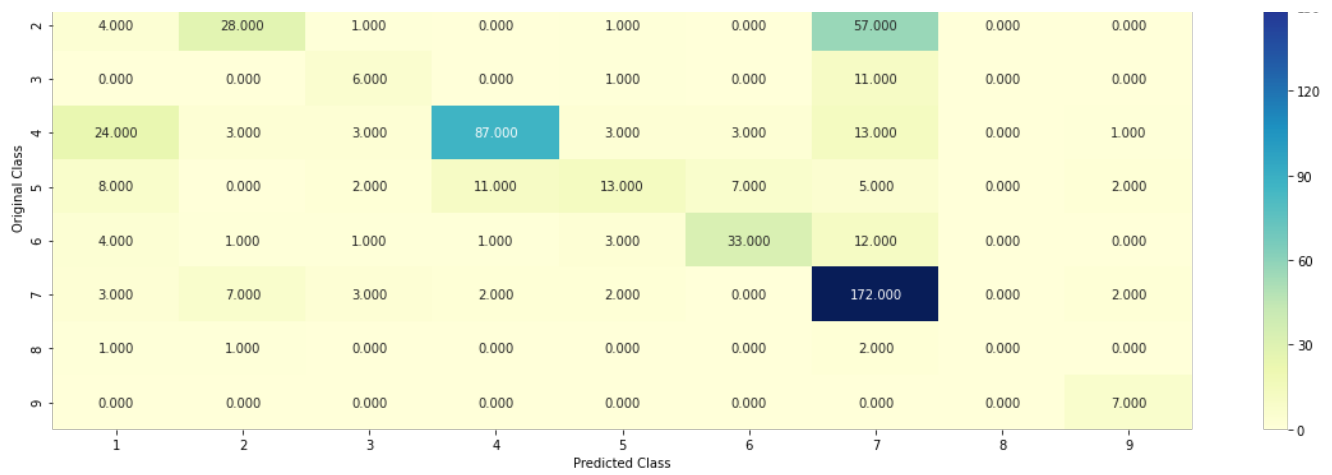
Log loss (CV) on the VotingClassifier : 1.2089743022996116

Log loss (test) on the VotingClassifier : 1.224280685055379

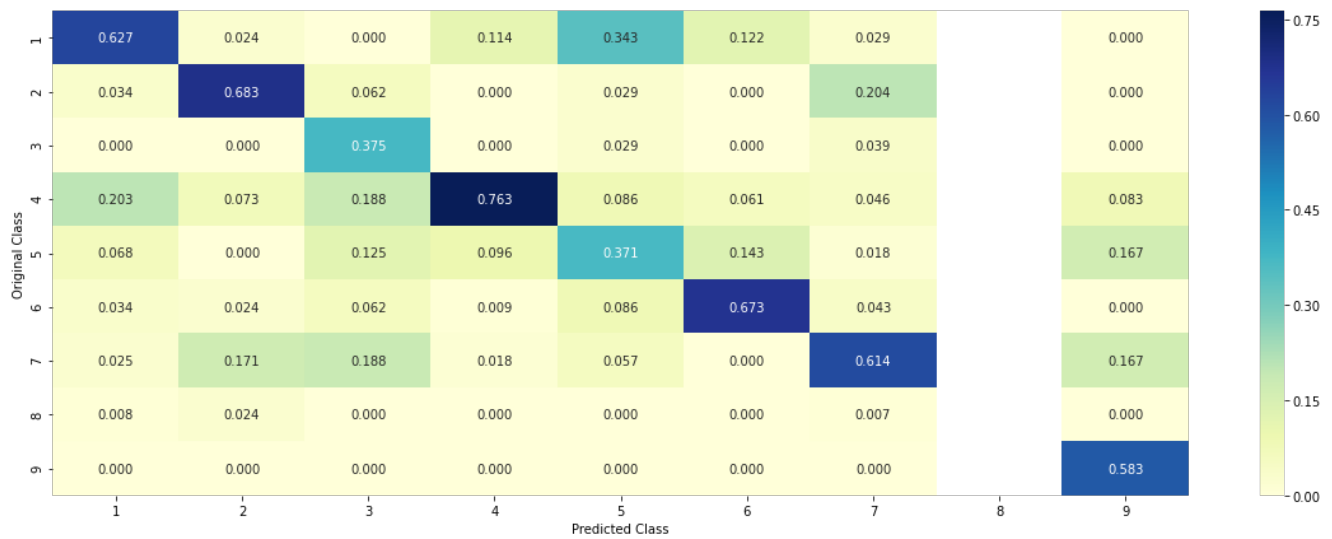
Number of missclassified point : 0.3684210526315789

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

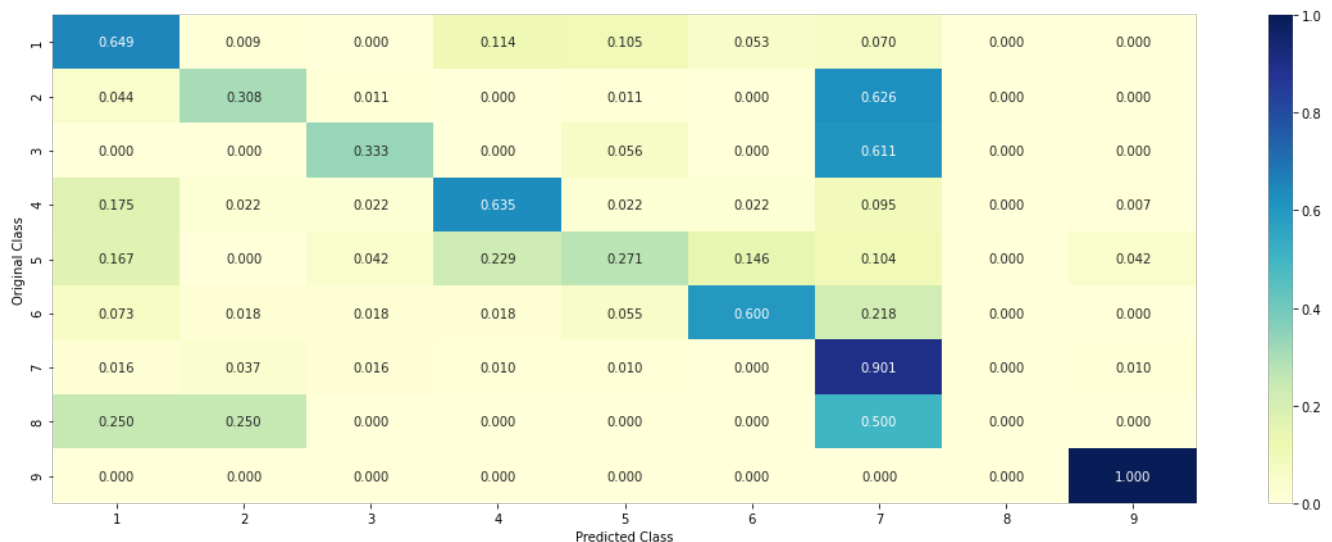




----- Precision matrix (Column Sum=1) -----



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



5. Assignments

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

In [94]:

```
# Apply Logistic regression + BOW + FEATURE ENGINEERING
```

In [95]:

```
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = CountVectorizer(min_df=5, ngram_range=(1, 4))
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 : 1651479

In [97]:

```
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 [98]:

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

In [99]:

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

In [100]:

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

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

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

In [101]:

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

In [102]:

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

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

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

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

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

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

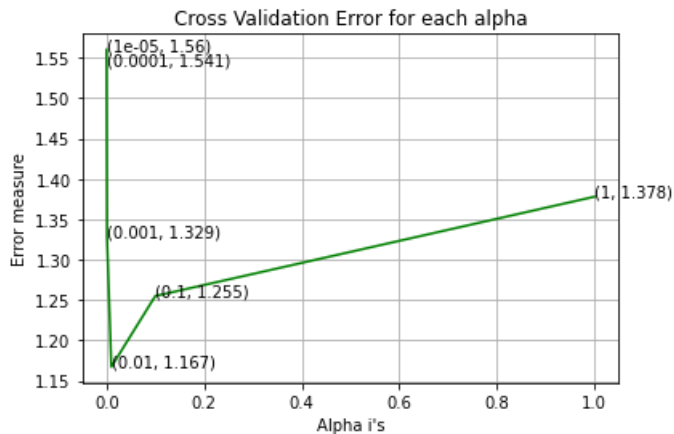
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_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: 1.5599603280867194
 For values of alpha = 0.0001 The log loss is: 1.5407830894889714
 For values of alpha = 0.001 The log loss is: 1.328815193831239
 For values of alpha = 0.01 The log loss is: 1.1672746425556175
 For values of alpha = 0.1 The log loss is: 1.25497854581271
 For values of alpha = 1 The log loss is: 1.3779512778314134



For values of best alpha = 0.01 The train log loss is: 0.7938722731971273
 For values of best alpha = 0.01 The cross validation log loss is: 1.1672746425556175
 For values of best alpha = 0.01 The test log loss is: 1.2379315042072556

In [103]:

```
def get_intersec_text(df):
    df_text_vec = CountVectorizer(min_df=5, ngram_range=(1, 4))
    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 [104]:

```
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")
```

96.525 % of word of test data appeared in train data
 99.583 % of word of Cross Validation appeared in train data

In [105]:

```
# this function will be used just for naive bayes
# for the given indices, we will print the name of the features
# and we will check whether the feature present in the test point text or not
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = CountVectorizer()
    var_count_vec = CountVectorizer()
    text_count_vec = CountVectorizer(min_df=3, ngram_range=(1, 2))

    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")

```

In [106]:

```

# 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)
)

# Adding the train_text feature

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

# Adding the test_text feature

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

# Adding the cv_text feature

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 [107]:

```
# 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 probabilities we use log-probability estimates
    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

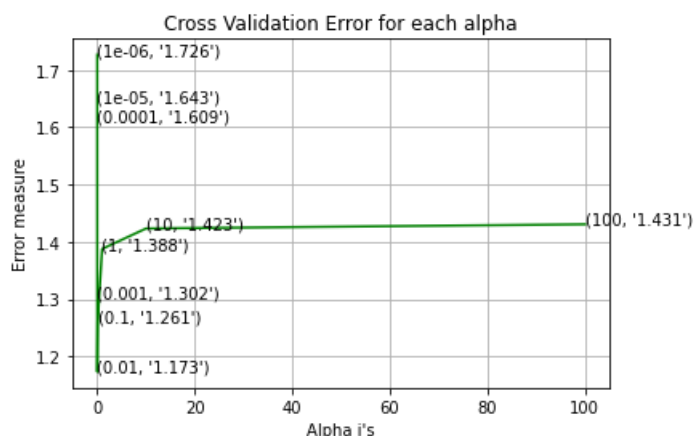
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], str(txt)), (alpha[i], cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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

```
redict_y, labels=clf.classes_, eps=1e-15))
```

```
for alpha = 1e-06
Log Loss : 1.7263704714277415
for alpha = 1e-05
Log Loss : 1.6429507688877656
for alpha = 0.0001
Log Loss : 1.6094218707755106
for alpha = 0.001
Log Loss : 1.3016380450411291
for alpha = 0.01
Log Loss : 1.173108663268534
for alpha = 0.1
Log Loss : 1.260636539708206
for alpha = 1
Log Loss : 1.3875789395392888
for alpha = 10
Log Loss : 1.4232179348799643
for alpha = 100
Log Loss : 1.4305116155396154
```



```
For values of best alpha = 0.01 The train log loss is: 0.7865478277365986
For values of best alpha = 0.01 The cross validation log loss is: 1.173108663268534
For values of best alpha = 0.01 The test log loss is: 1.2291798670450922
```

In [108]:

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

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

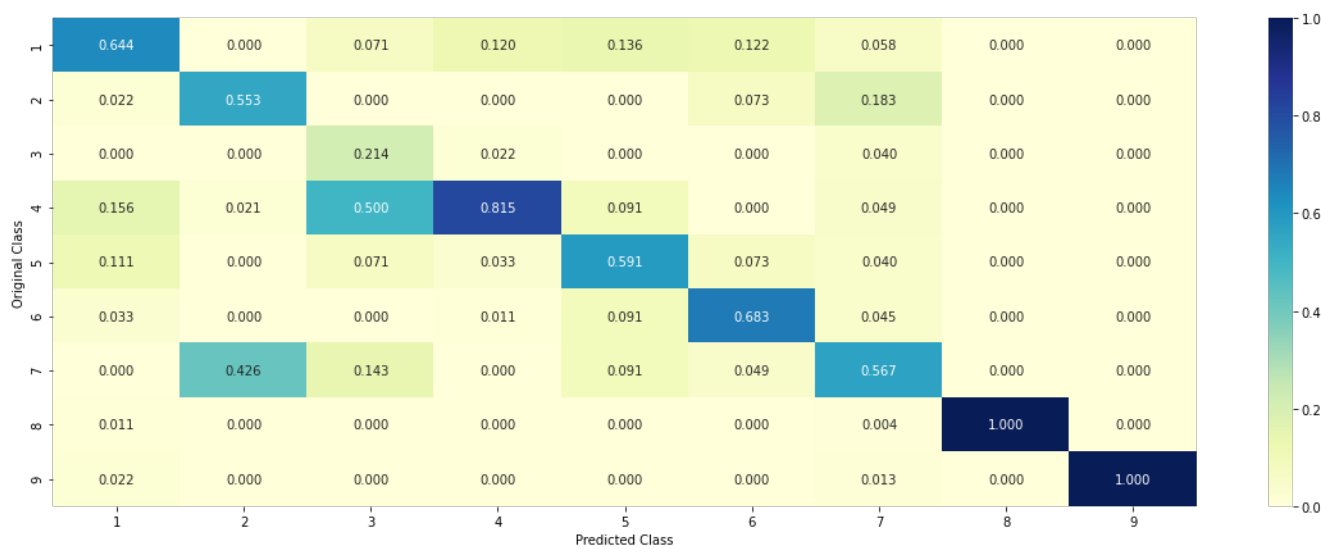
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,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

```
Log loss : 1.173108663268534
Number of mis-classified points : 0.37593984962406013
----- Confusion matrix -----
```

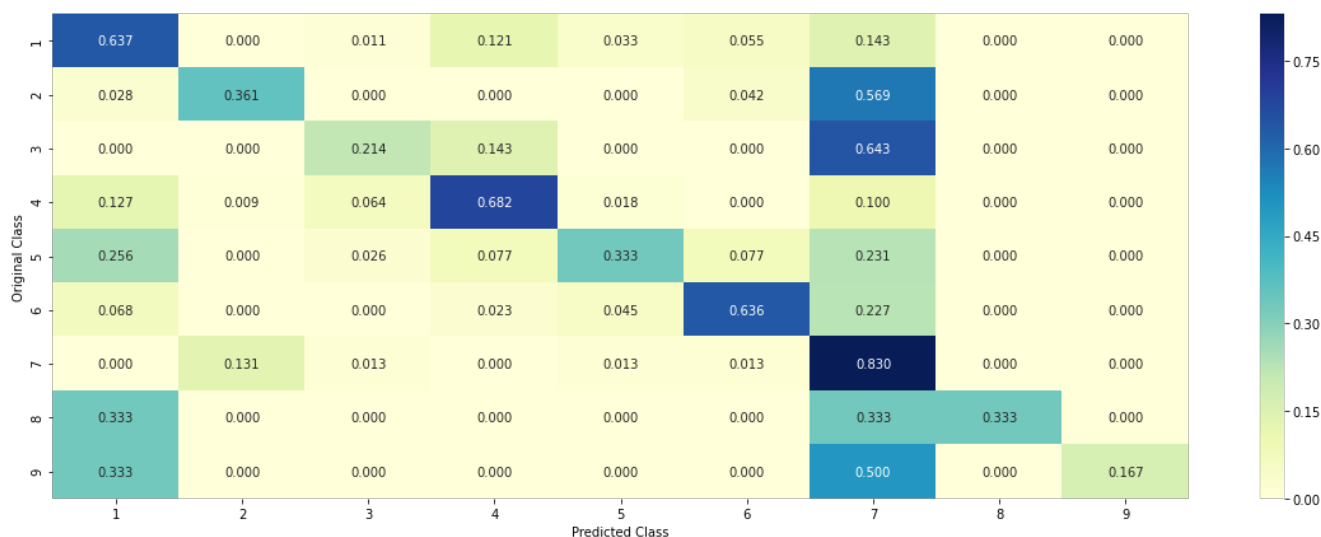
58.000	0.000	1.000	11.000	3.000	5.000	13.000	0.000	0.000
--------	-------	-------	--------	-------	-------	--------	-------	-------



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



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



In [109]:

```
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"])
            incresingorder_ind += 1
    return tabulte_list
```

```

    elif i < 10:
        tabulte_list.append([increasingorder_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([increasingorder_ind, train_text_features[i], yes_no])
    increasingorder_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"]))

```

In [110]:

```
# Logistic regression with unbalanced class weight
```

In [111]:

```

# 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(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=41
)
    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")

```



```

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',)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

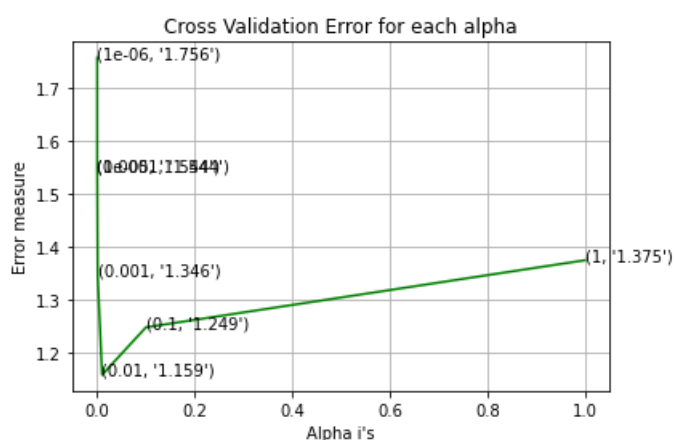
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "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 : 1.7563560808889553
for alpha = 1e-05
Log Loss : 1.5441654765395623
for alpha = 0.0001
Log Loss : 1.5435071569016814
for alpha = 0.001
Log Loss : 1.3464524884531262
for alpha = 0.01
Log Loss : 1.1589057114391617
for alpha = 0.1
Log Loss : 1.2485006389438005
for alpha = 1
Log Loss : 1.3747930827218593

```



```

For values of best alpha = 0.01 The train log loss is: 0.7686832073187695
For values of best alpha = 0.01 The cross validation log loss is: 1.1589057114391617
For values of best alpha = 0.01 The test log loss is: 1.2276314923463043

```

In [112]:

```

# 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.

#-----

```

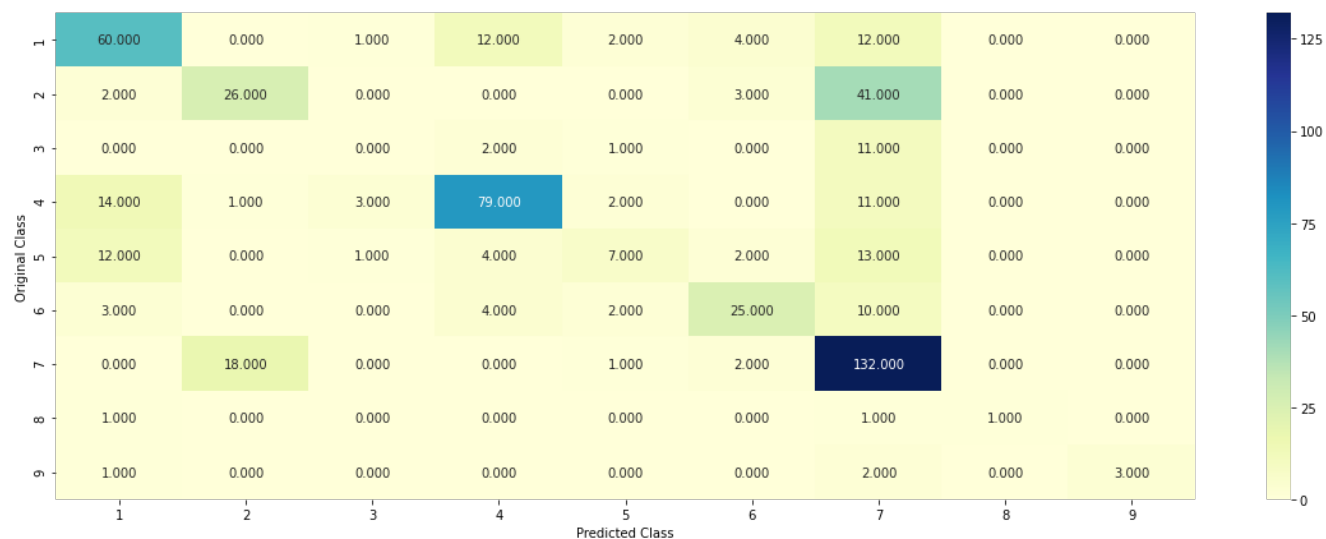
```
# video link:
#-----

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

Log loss : 1.1589057114391617

Number of mis-classified points : 0.37406015037593987

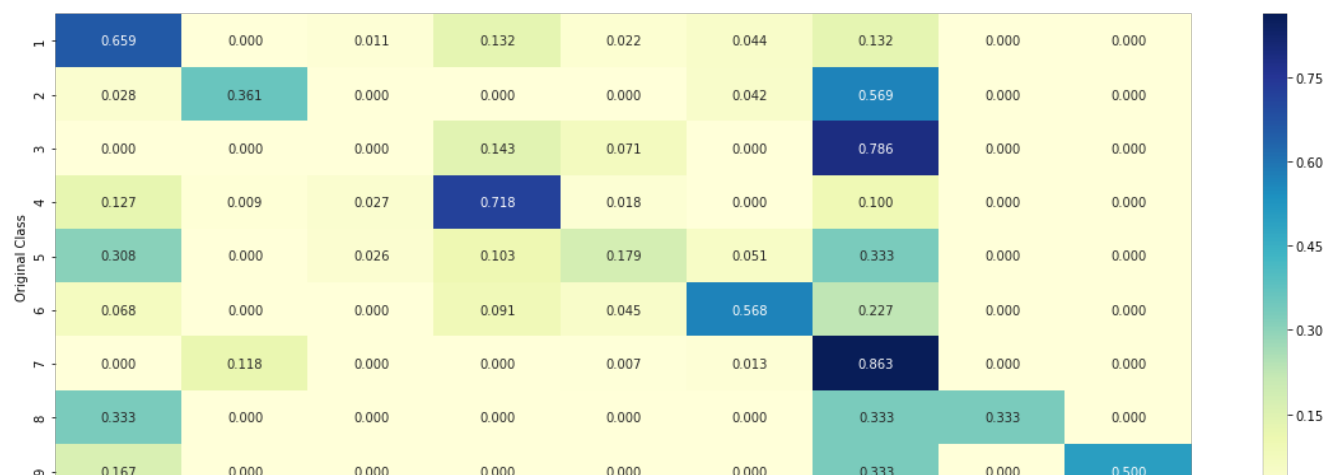
----- Confusion matrix -----

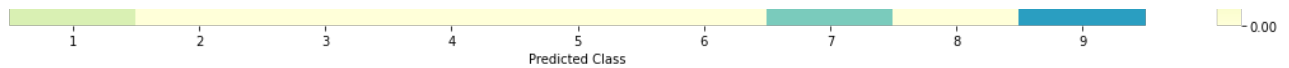


----- Precision matrix (Column Sum=1) -----



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





In [114]:

```
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["MODEL", "FEATURIZATION", "LOG-LOSS", "% MISCLASSIFICATION"]
x.add_row(["NAIVE BAYES", "ONE HOT ENCODING + TF-IDF", 1.18, 35.5])
x.add_row(["KNN", "RESPONSE CODING + TF-IDF", 1.07, 36.09])
x.add_row(["LOGISTIC REGRESSION WITH CLASS BALANCED", "ONE HOT ENCODING + TF-IDF", 0.99, 32.14])
x.add_row(["LOGISTIC REGRESSION IMBALANCED", "ONE HOT ENCODING + TF-IDF", 1.003, 32.3])
x.add_row(["SVM BALANCED", "ONE HOT ENCODING + TF-IDF", 1.05, 30.26])
x.add_row(["RANDOM FOREST", "ONE HOT ENCODING + TF-IDF", 1.108, 34.39])
x.add_row(["RANDOM FOREST", "RESPONSE CODING + TF-IDF", 1.28, 41.72])
x.add_row(["STACKING LOGISTIC REG + SVM + MULTINOMIALNB", "ONE HOT ENCODING + TF-IDF", 1.104, 35.6])
x.add_row(["STACKING LOGISTIC + SVM + MULTINOMIALNB + MAX VOTING", "ONE HOT ENCODING + TF-IDF", 1.17, 35.3])
x.add_row(["LOGISTIC REGRESSION", "ONE HOT ENCODING + BOW + BALANCED ", 1.14, 35.7])
x.add_row(["LOISTIC REGRESSION", "ONE HOT ENCODING + BOW + UNBALANCED ", 1.15, 32.7])

print(x)
```

MODEL			FEATURIZATION		LOG
SS	% MISCLASSIFICATION				

NAIVE BAYES			ONE HOT ENCODING + TF-IDF		1
8	35.5				
KNN			RESPONSE CODING + TF-IDF		1
7	36.09				
LOGISTIC REGRESSION WITH CLASS BALANCED			ONE HOT ENCODING + TF-IDF		0
.99	32.14				
LOGISTIC REGRESSION IMBALANCED			ONE HOT ENCODING + TF-IDF		1.
03	32.3				
SVM BALANCED			ONE HOT ENCODING + TF-IDF		
1.05	30.26				
RANDOM FOREST			ONE HOT ENCODING + TF-IDF		
1.108	34.39				
RANDOM FOREST			RESPONSE CODING + TF-IDF		
1.28	41.72				
STACKING LOGISTIC REG + SVM + MULTINOMIALNB			ONE HOT ENCODING + TF-IDF		1.
104	35.6				
STACKING LOGISTIC + SVM + MULTINOMAILNB + MAX VOTING			ONE HOT ENCODING + TF-IDF		
1.17	35.3				
LOGISTIC REGRESSION			ONE HOT ENCODING + BOW + BALANCED		1
14	35.7				
LOISTIC REGRESSION			ONE HOT ENCODING + BOW + UNBALANCED		1
15	32.7				

In [114]:

```
gene_vectorizer = TfidfVectorizer()
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(x_train['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(x_test['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(x_cv['Gene'])
```

In [115]:

```
variation_vectorizer = TfidfVectorizer()
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(x_train['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(x_test['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(x_cv['Variation'])
```

In [116]:

```
text_vectorizer = TfidfVectorizer()
train_text_feature_onehotCoding = text_vectorizer.fit_transform(x_train['TEXT'])
```

```
# getting all the feature names (words)
train_text_features= text_vectorizer.get_feature_names()

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

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

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

Total number of unique words in train data : 124799

In [117]:

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

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

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

In [118]:

```
train_text_feature_responseCoding = get_text_responsecoding(x_train)
test_text_feature_responseCoding = get_text_responsecoding(x_test)
cv_text_feature_responseCoding = get_text_responsecoding(x_cv)

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

In [119]:

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

In [120]:

```
gene_variation = []

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

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

In [121]:

```

tfidfvectorizer = TfidfVectorizer(max_features=3000)
text2 = tfidfVectorizer.fit_transform(gene_variation)
gene_variation_features = tfidfVectorizer.get_feature_names()

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

```

In [122]:

```

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

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

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

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

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

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

```

In [124]:

```

lpha = [10 ** x for x in range(-4, 1)]
cv_log_error_array = []
for i in lpha:
    print("for alpha =", i)
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=41
)
    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(lpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((lpha[i],str(txt)), (lpha[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=lpha[best_alpha], penalty='l2', loss='log',)

```

```

clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

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

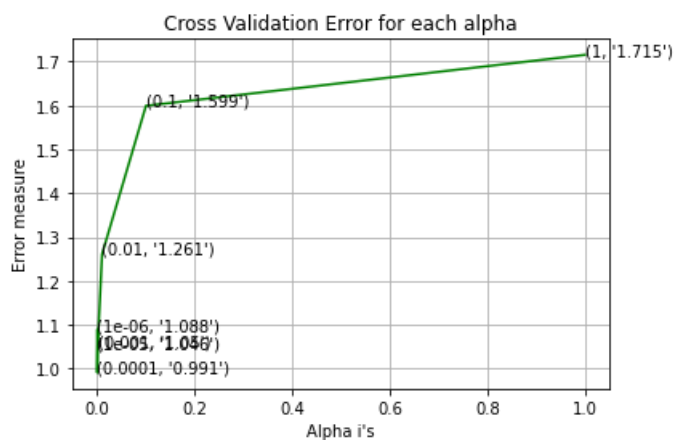
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ',
      alpha[best_alpha], "The test log loss is:",
      log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for alpha = 1e-06
Log Loss : 1.0876648136915068
for alpha = 1e-05
Log Loss : 1.0457670450272398
for alpha = 0.0001
Log Loss : 0.9906657845466823
for alpha = 0.001
Log Loss : 1.0504868991879939
for alpha = 0.01
Log Loss : 1.2611886997241484
for alpha = 0.1
Log Loss : 1.5991381575536485
for alpha = 1
Log Loss : 1.71524128837181

```



```

For values of best alpha = 0.0001 The train log loss is: 0.42578344411377717
For values of best alpha = 0.0001 The cross validation log loss is: 0.9953712908857676
For values of best alpha = 0.0001 The test log loss is: 1.0251880126512387

```

In [125]:

```

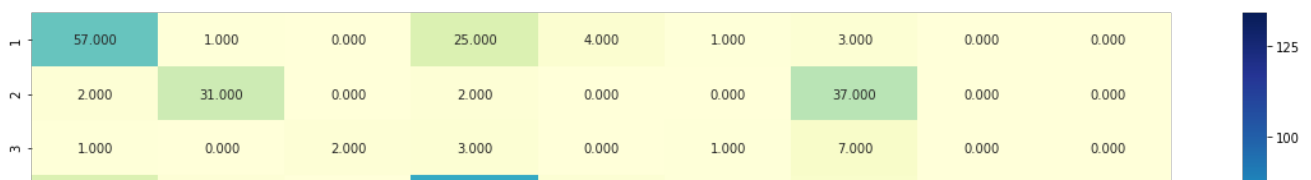
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log',)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

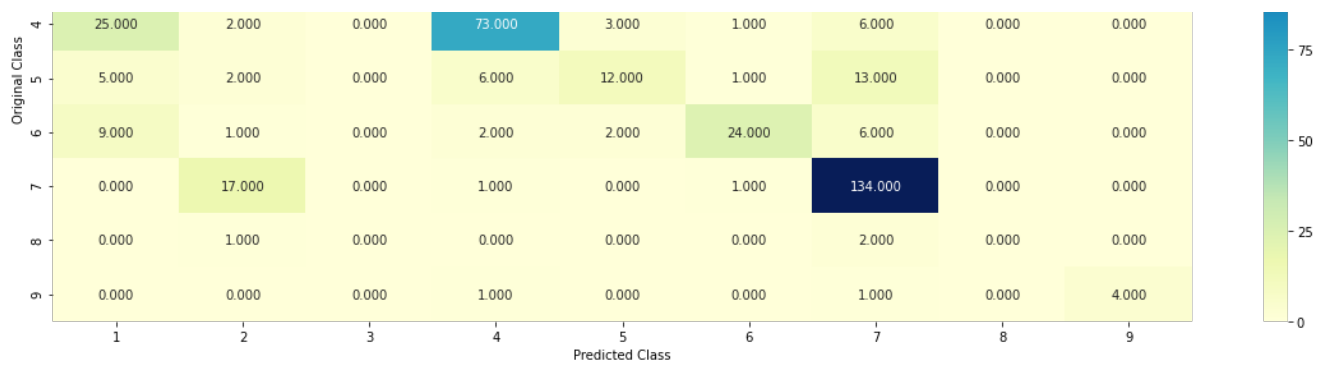
```

```

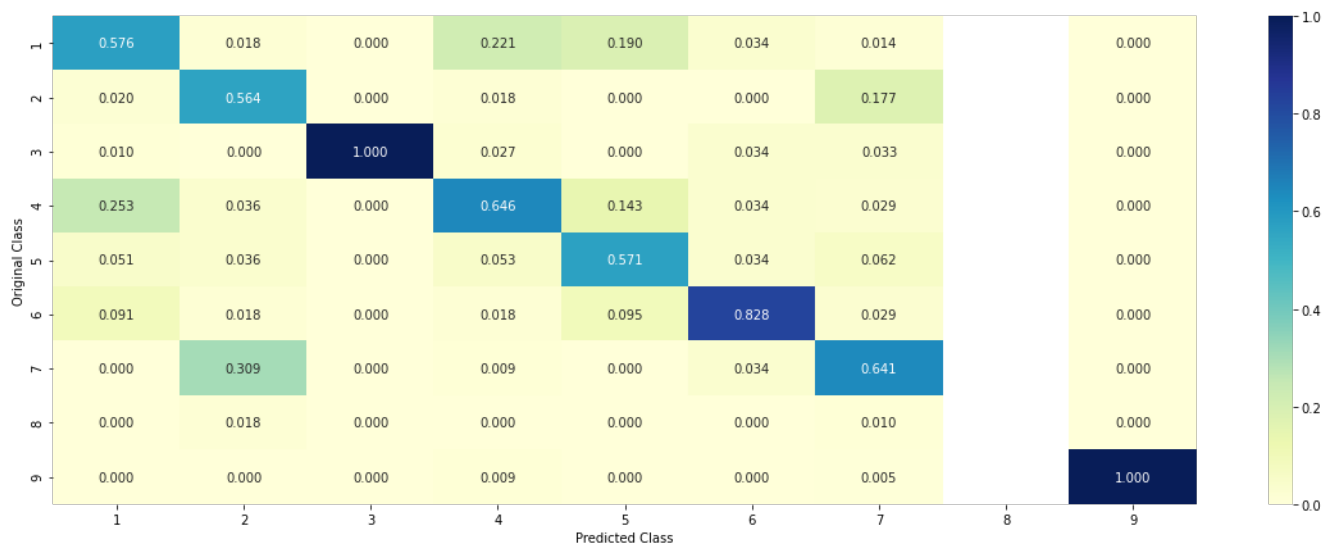
Log loss : 0.9902237185310909
Number of mis-classified points : 0.36654135338345867
----- Confusion matrix -----

```

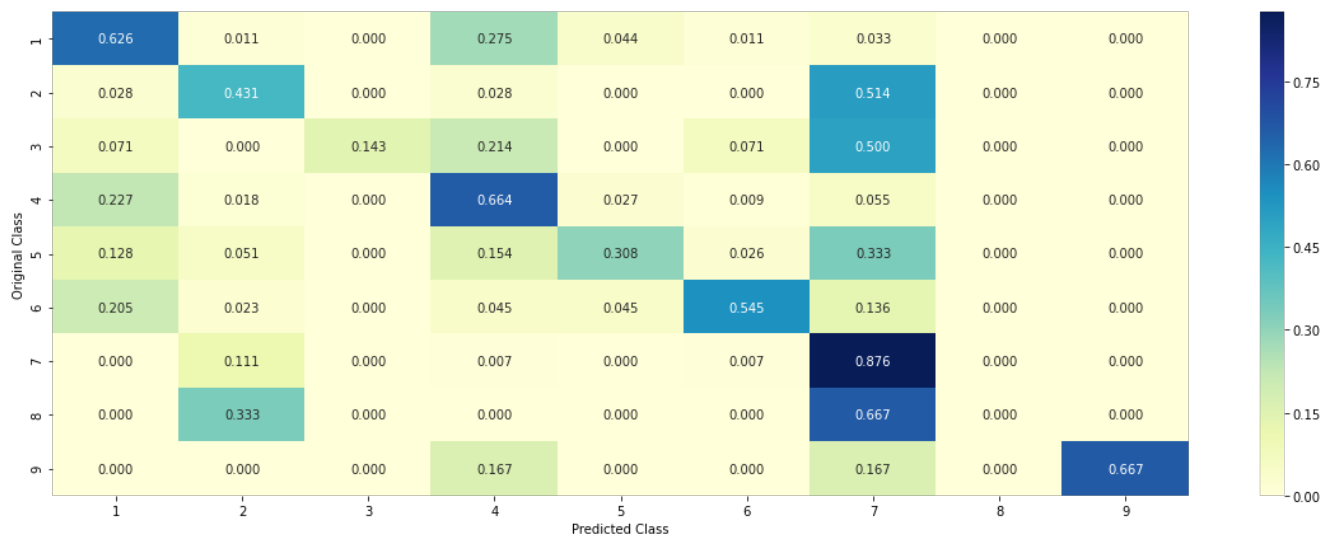




----- Precision matrix (Column Sum=1) -----



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



Logistic Regression With Class balancing

In [100]:

```
alpha = [10 ** x for x in range(-6, 3)]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_x_onehotCoding, train_y)
    sig clf = CalibratedClassifierCV(clf, method="sigmoid")
```

```

sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
# to avoid rounding error while multiplying probabilities we use log-probability estimates
print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

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

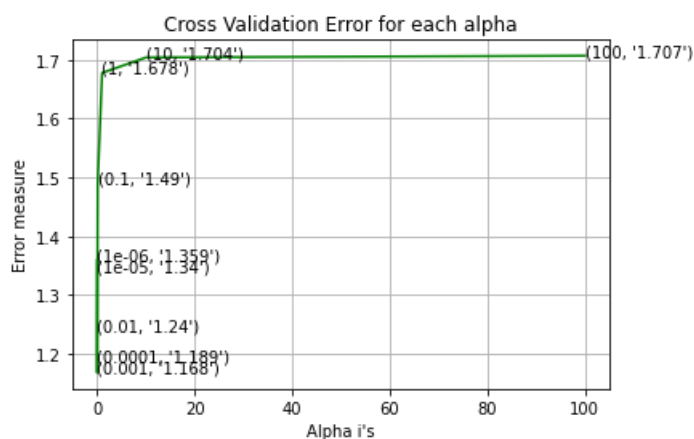
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 1e-06
Log Loss : 1.3594447834565797
for alpha = 1e-05
Log Loss : 1.3400679247846996
for alpha = 0.0001
Log Loss : 1.1886254498018929
for alpha = 0.001
Log Loss : 1.1679744423667313
for alpha = 0.01
Log Loss : 1.2402376917451046
for alpha = 0.1
Log Loss : 1.4899194255908614
for alpha = 1
Log Loss : 1.6775930620964115
for alpha = 10
Log Loss : 1.7038258700952287
for alpha = 100
Log Loss : 1.7067726556350924

```



```

For values of best alpha = 0.001 The train log loss is: 0.5375028601782635
For values of best alpha = 0.001 The cross validation log loss is: 1.1679744423667313
For values of best alpha = 0.001 The test log loss is: 0.9903096911558612

```

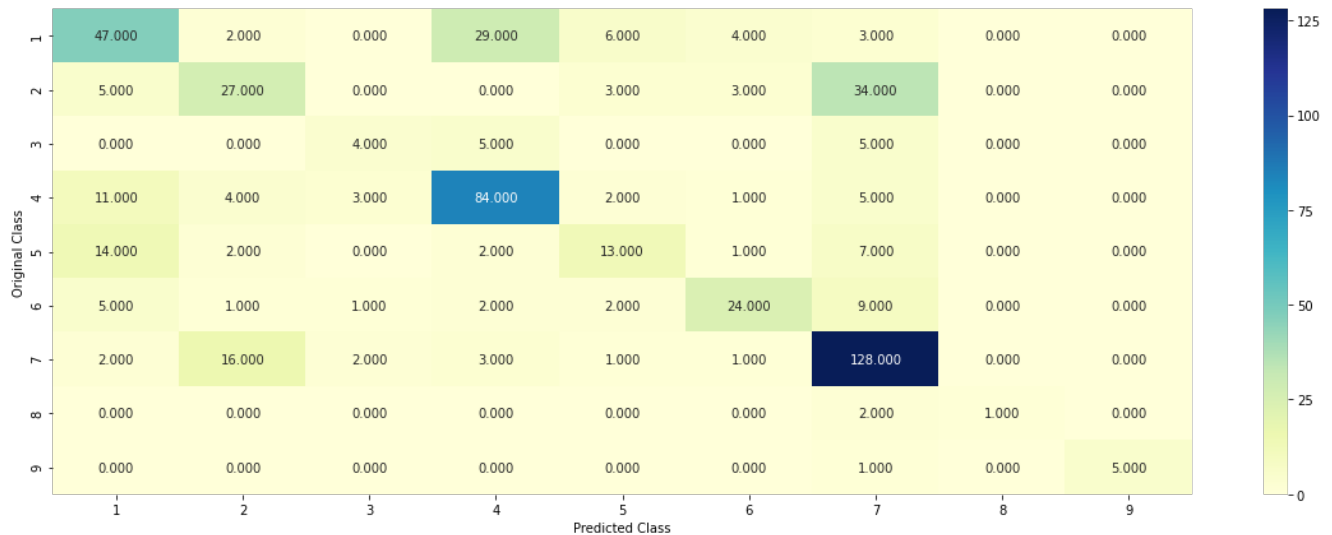

In [101]:

```
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)
```

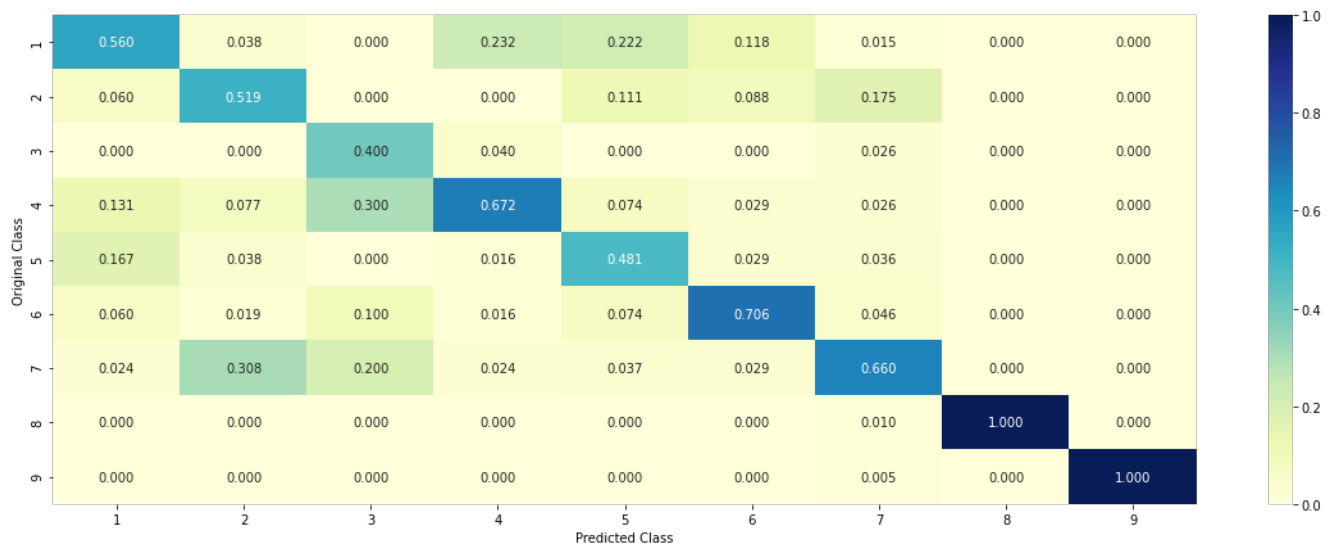
Log loss : 1.1679744423667313

Number of mis-classified points : 0.37406015037593987

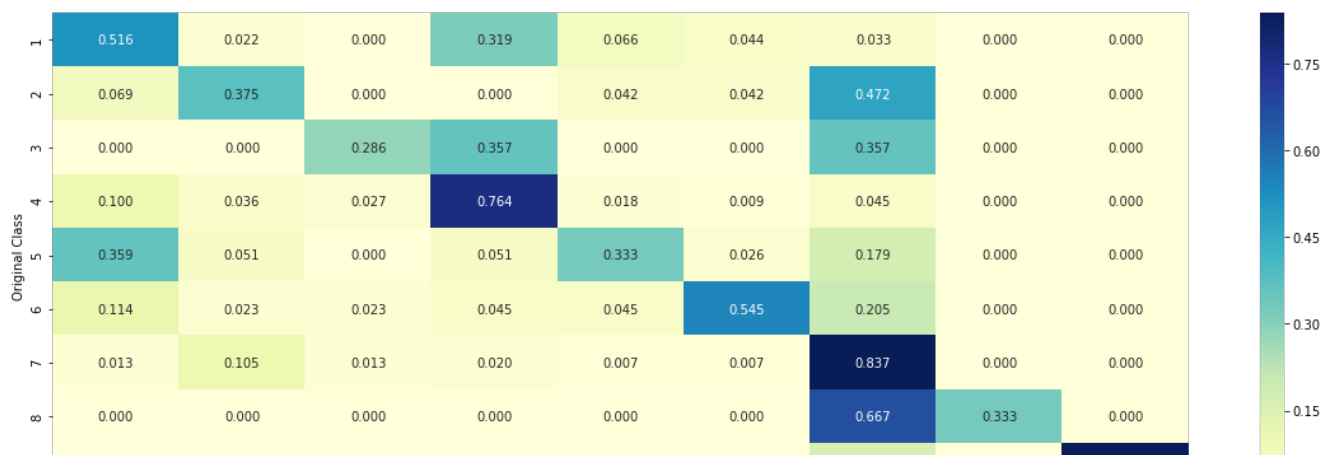
```
----- Confusion matrix -----
```

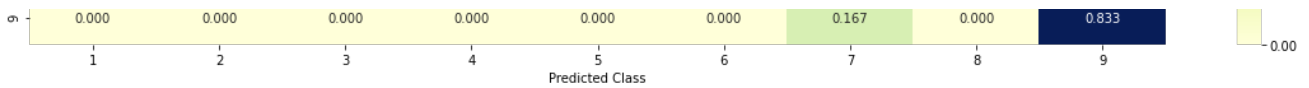


```
----- Precision matrix (Columm Sum=1) -----
```



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





Logistic Regression Without Class balancing

In [103]:

```
# 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, 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)
```

```

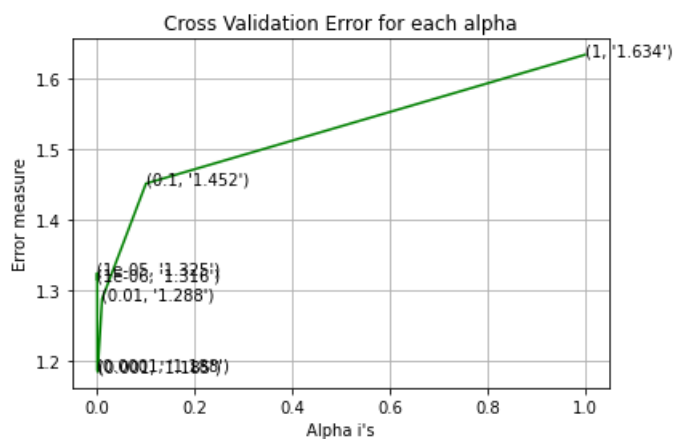
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_lo
ss(y_cv, 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 : 1.3159249754622169
for alpha = 1e-05
Log Loss : 1.324973193056196
for alpha = 0.0001
Log Loss : 1.187899597991038
for alpha = 0.001
Log Loss : 1.1849514236726115
for alpha = 0.01
Log Loss : 1.2882429539924007
for alpha = 0.1
Log Loss : 1.4516565338172593
for alpha = 1
Log Loss : 1.6340400761827116

```



```

For values of best alpha = 0.001 The train log loss is: 0.5423437489191432
For values of best alpha = 0.001 The cross validation log loss is: 1.1849514236726115
For values of best alpha = 0.001 The test log loss is: 0.9957491301603113

```

In [104]:

```

# 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:
#-----

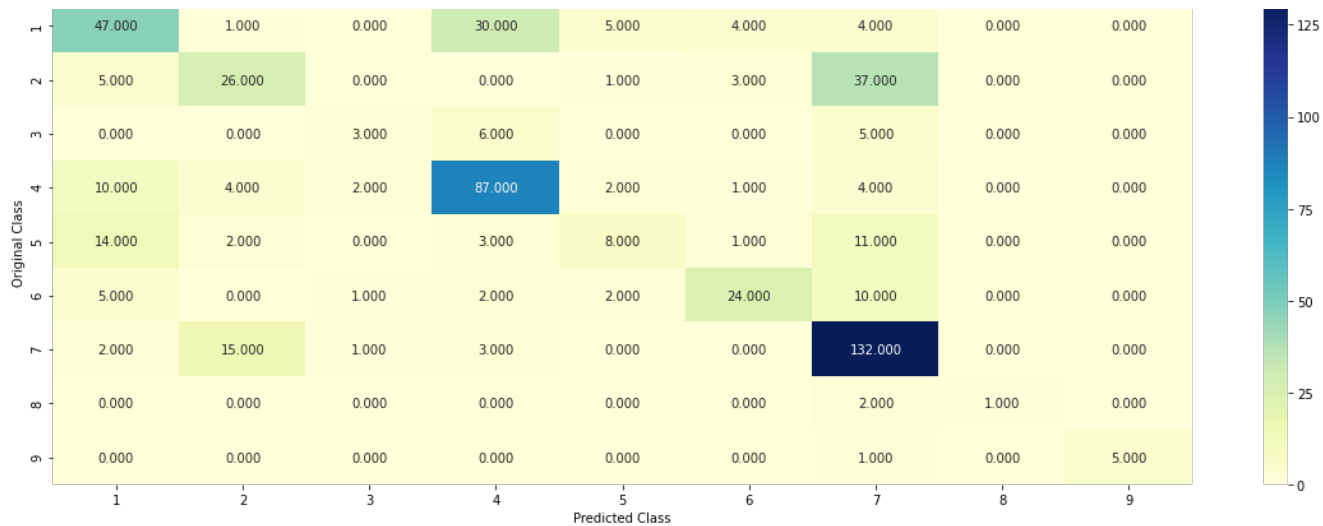
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

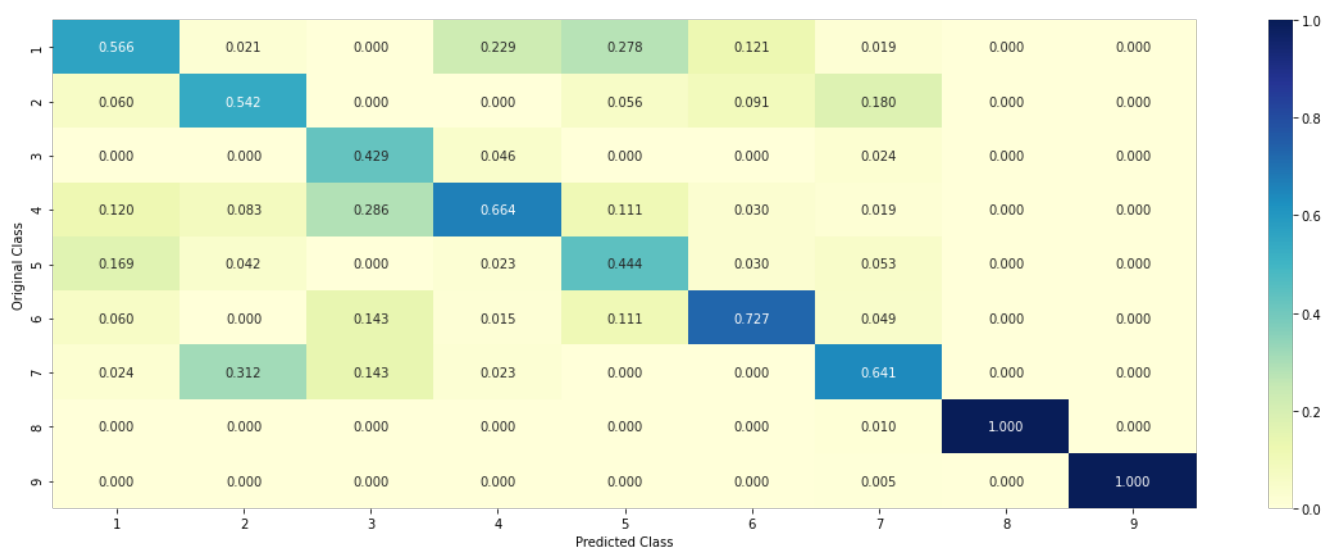
```

Log loss : 1.1849514236726115
Number of mis-classified points : 0.37406015037593987
----- Confusion matrix -----

```



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)

