

# Personalized cancer diagnosis

## 1. Business Problem

### 1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training\_variants.zip and training\_text.zip from Kaggle.

#### **Context:**

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

#### **Problem statement :**

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

### 1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

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

Metric(s):

- Multi class log-loss
- Confusion matrix

### 2.2.3. Machine Learing Objectives and Constraints

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

Constraints:

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

## 2.3. Train, CV and Test Datasets

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

## 3. Exploratory Data Analysis

In [175]:

```
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, CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import SGDClassifier
from imblearn.over_sampling import SMOTE
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from sklearn.model_selection import StratifiedKFold
from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier
warnings.filterwarnings("ignore")

from mlxtend.classifier import StackingClassifier

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

## 3.1. Reading Data

### 3.1.1. Reading Gene and Variation Data

In [2]:

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

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

Out[2]:

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

training/training\_variants is a comma separated file containing the description of the genetic mutations used for training.

Fields are

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

### 3.1.2. Reading Text Data

In [3]:

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

Number of data points : 3321

Number of features : 2

Features : ['ID' 'TEXT']

Out[3]:

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

### 3.1.3. Preprocessing of text

In [4]:

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

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

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

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

In [5]:

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

there is no text description for id: 1109  
 there is no text description for id: 1277  
 there is no text description for id: 1407  
 there is no text description for id: 1639  
 there is no text description for id: 2755  
 Time took for preprocessing the text : 306.3311523 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*                | abstract background non small cell lung cancer...   |
| 2  | 2    | CBL       | Q249E                | abstract background non small cell lung cancer...   |
| 3  | 3    | CBL       | N454D                | recent evidence demonstrated acquired uniparen...   |
| 4  | 4    | CBL       | L399V                | oncogenic mutations monomeric casitas b lineage...  |

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'
X_train, test_df, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_size=0.15)
# split the train data into train and cross validation by maintaining same distribution of 'y'
train_df, cv_df, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.15)
```

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], '(%')

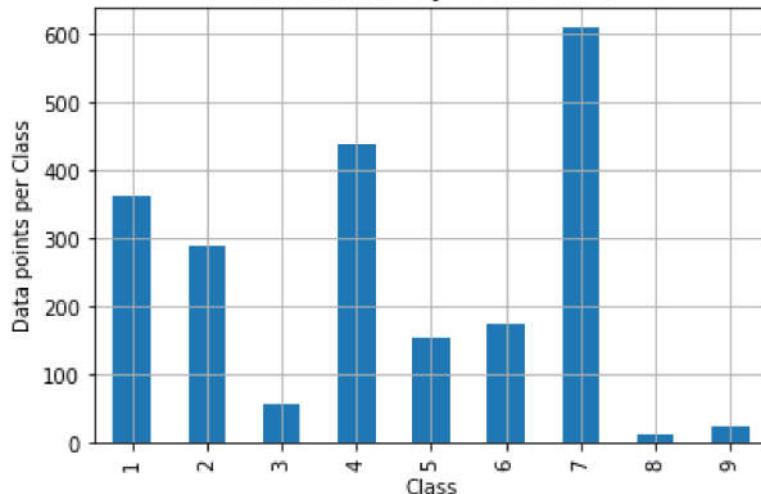
print('*'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

# 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], '(%'

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], '(%'
```

Distribution of yi in train data



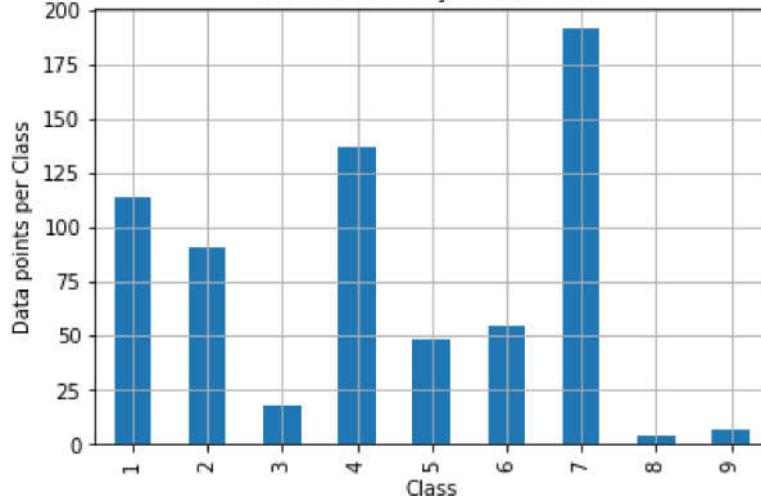
Number of data points in class 7 : 609 ( 28.672 %)  
 Number of data points in class 4 : 439 ( 20.669 %)  
 Number of data points in class 1 : 363 ( 17.09 %)  
 Number of data points in class 2 : 289 ( 13.606 %)  
 Number of data points in class 6 : 176 ( 8.286 %)  
 Number of data points in class 5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

---



---

Distribution of yi in test data

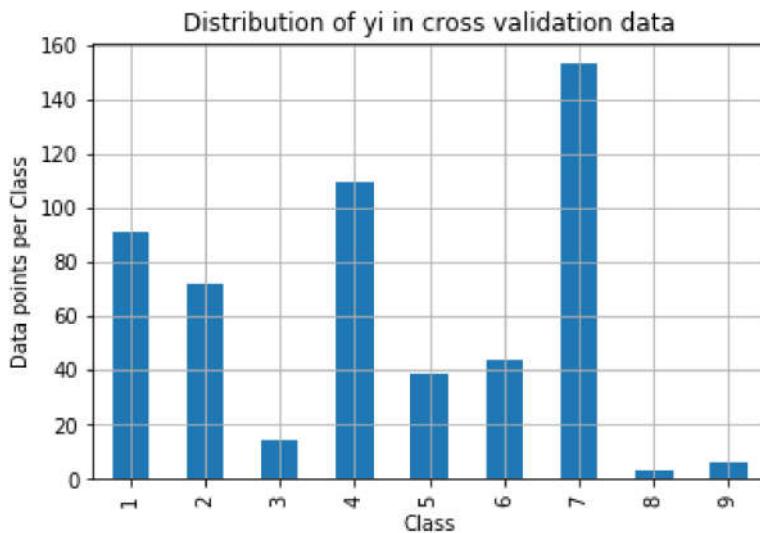


Number of data points in class 7 : 191 ( 28.722 %)  
 Number of data points in class 4 : 137 ( 20.602 %)  
 Number of data points in class 1 : 114 ( 17.143 %)  
 Number of data points in class 2 : 91 ( 13.684 %)  
 Number of data points in class 6 : 55 ( 8.271 %)  
 Number of data points in class 5 :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

---



---



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

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

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

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

    B = (C/C.sum(axis=0))
    #divid each element of the confusion matrix with the sum of elements in that row

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

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

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

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

In [14]:

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

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

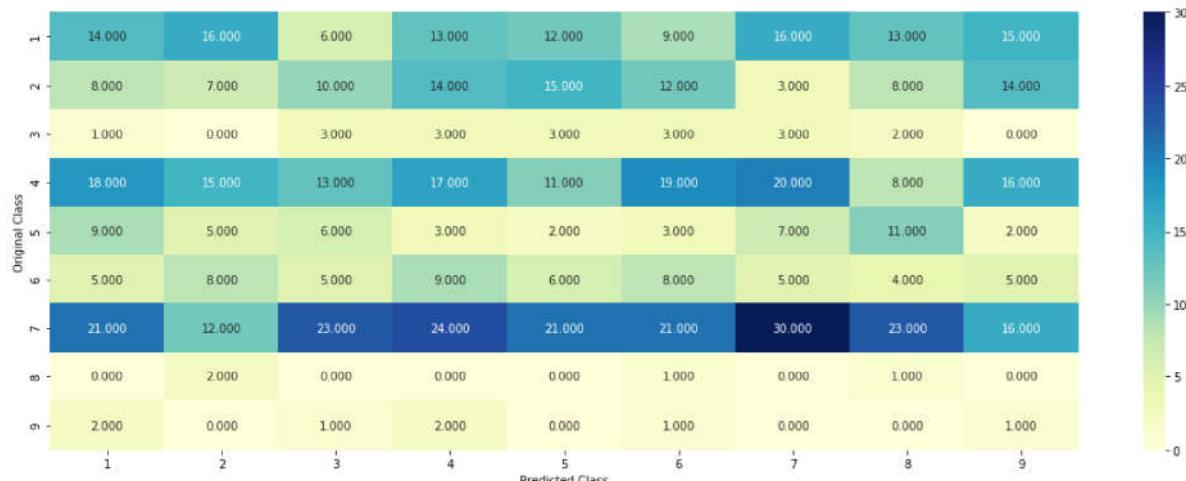
# Test-Set error.
#we create a output array that has exactly same as the test data
test_predicted_y = np.zeros((test_data_len,9))
for i in range(test_data_len):
    rand_probs = np.random.rand(1,9)
    test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random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

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

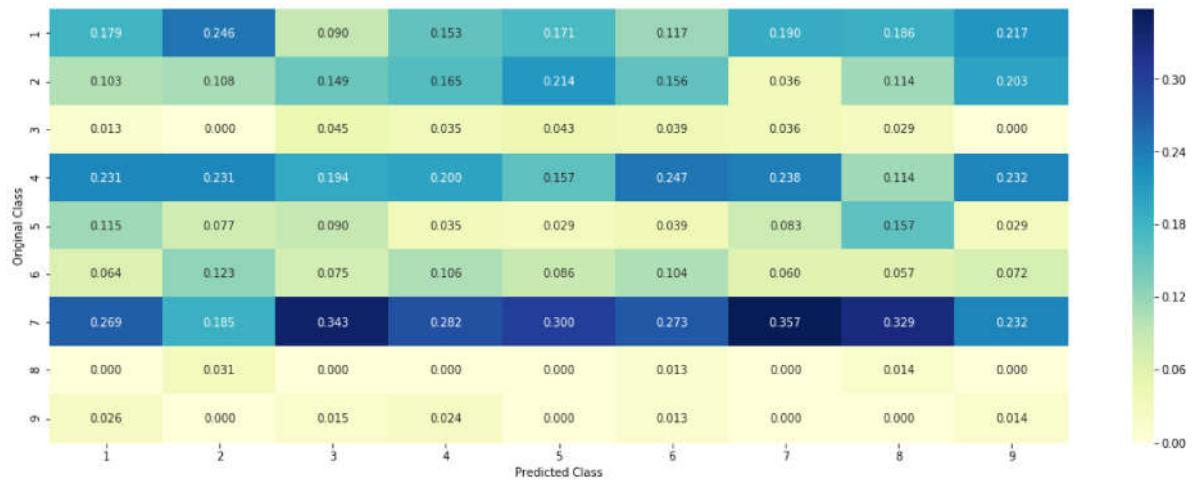
Log loss on Cross Validation Data using Random Model 2.5164445887481466

Log loss on Test Data using Random Model 2.518711330495878

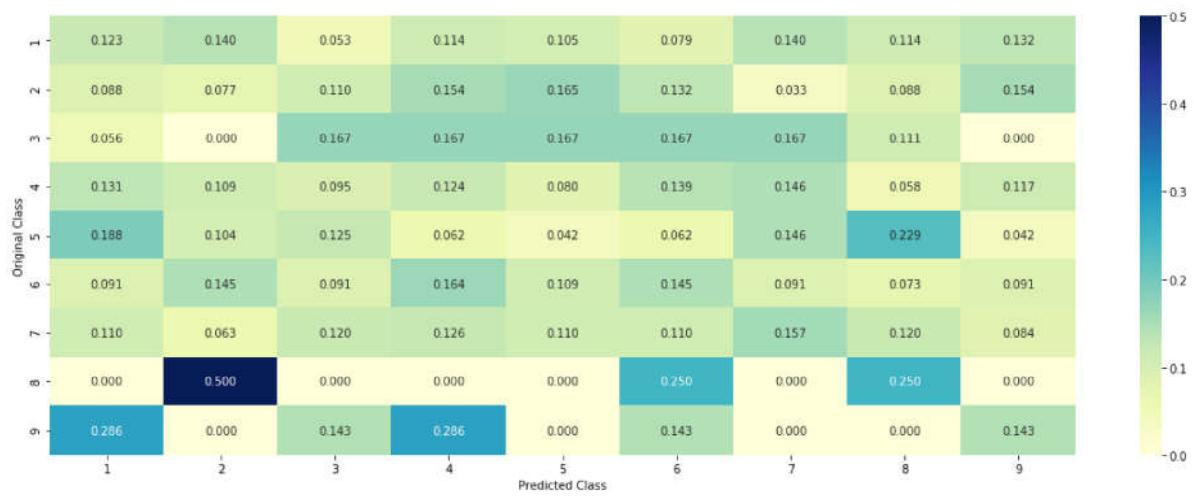
----- Confusion matrix -----



----- 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 da
# build a vector (1*9) , the first element = (number of times it occurred in class1 + 10*alp
# gv_dict is like a look up table, for every gene it stores a (1*9) representation of it
# for a value of feature in df:
# if it is in train data:
# we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
# if it is not there is train:
# we add [1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
# return 'gv_fea'
# -----
# get_gv_fea_dict: Get Gene variation Feature Dict
def get_gv_fea_dict(alpha, feature, df):
    # value_count: it contains a dict like
    # print(train_df['Gene'].value_counts())
    # output:
    #     {BRCA1      174
    #      TP53      106
    #      EGFR      86
    #      BRCA2      75
    #      PTEN      69
    #      KIT       61
    #      BRAF      60
    #      ERBB2      47
    #      PDGFRA     46
    #      ...}
    # print(train_df['Variation'].value_counts())
    # output:
    # {
    # Truncating_Mutations      63
    # Deletion                  43
    # Amplification              43
    # Fusions                   22
    # Overexpression              3
    # E17K                      3
    # Q61L                      3
    # S222D                     2
    # P130S                     2
    # ...
    # }
    value_count = train_df[feature].value_counts()

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

    # denominator will contain the number of time that particular feature occurred in whole
    for i, denominator in value_count.items():
        # vec will contain ( $p(y_i=1/G_i)$ ) probability of gene/variation belongs to particular
        # 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 f
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': [0.200757575757575, 0.037878787878788, 0.0681818181818177, 0.1363
    # 'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.2704
    # 'EGFR': [0.0568181818181816, 0.21590909090909091, 0.0625, 0.0681818181818177
    # 'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.078
    # 'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.465
    # 'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.07284
    # 'BRAF': [0.066666666666666666, 0.17999999999999999, 0.07333333333333334, 0.0733
    # ...
    #     }
    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()

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

```

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

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

### 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 : 243

```
BRCA1    167  
TP53     109  
EGFR     95  
PTEN     79  
BRCA2     77  
KIT      65  
BRAF     57  
PIK3CA    44  
ERBB2     41  
ALK      41  
Name: Gene, dtype: int64
```

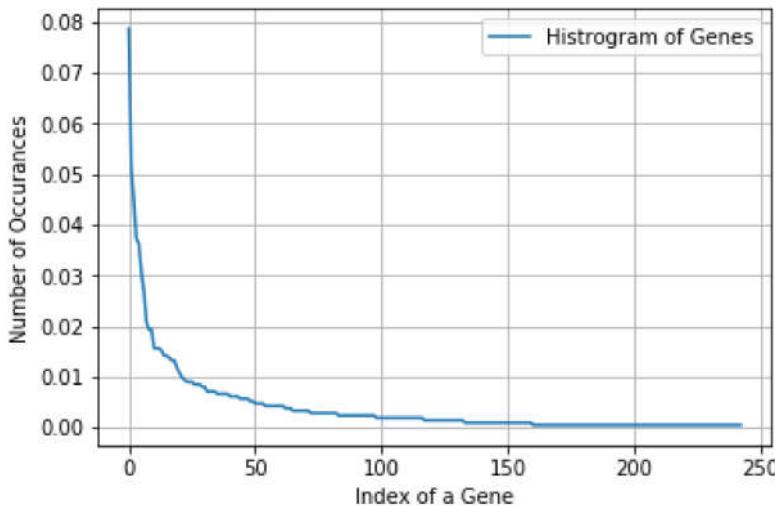
In [17]:

```
print("Ans: There are", unique_genes.shape[0] , "different categories of genes in the train
```

Ans: There are 243 different categories of genes in the train data, and they are distributed as follows

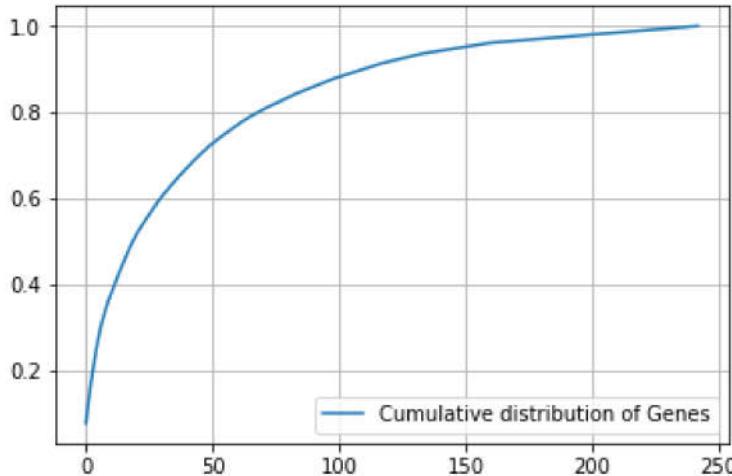
In [18]:

```
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 Occurrences')  
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/> (<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.")
```

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

In [22]:

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

In [23]:

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

Out[23]:

```
1707    PPM1D
252     EGFR
1638    RRAS2
2371    PTPN11
1233    PIM1
Name: Gene, dtype: object
```

In [24]:

```
train_gene_feature_onehotCoding
```

Out[24]:

```
<2124x242 sparse matrix of type '<class 'numpy.float64'>'  
with 2124 stored elements in Compressed Sparse Row format>
```

In [25]:

```
gene_vectorizer.get_feature_names()
```

Out[25]:

```
['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ar',
 'araf',
 'arid1b',
 'arid2',
 'arid5b',
 'asxl1',
 'asxl2',
 'atm',
 'atr',
 'atrx']
```

In [26]:

```
print("train_gene_feature_onehotCoding is converted feature using one-hot encoding method.")
```

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding  
method. The shape of gene feature: (2124, 242)
```

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

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

In [27]:

```

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

# read more about SGDClassifier() at http://scikit-Learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal'
# 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 Desce
# 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.cla

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()

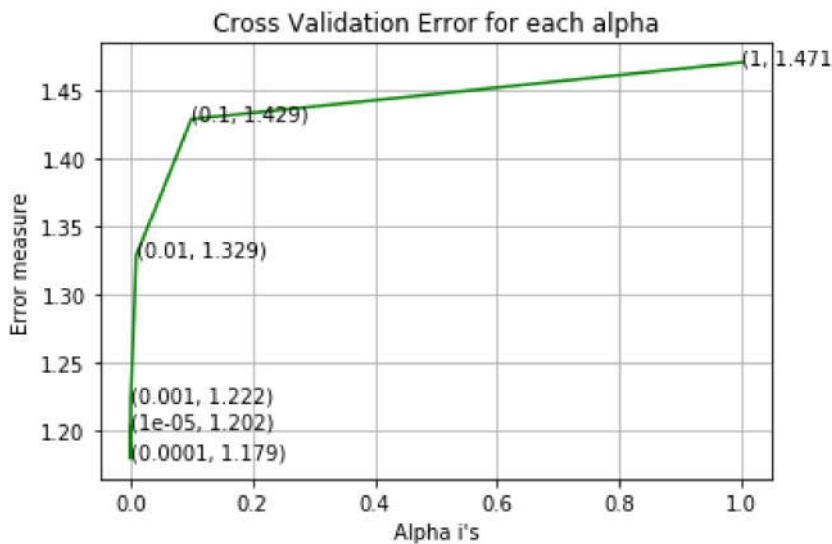
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_gene_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_gene_feature_onehotCoding, y_train)

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_cv, predict_y))
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:")
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_cv, predict_y))

```

For values of alpha = 1e-05 The log loss is: 1.201921313392381  
 For values of alpha = 0.0001 The log loss is: 1.1791330653831467  
 For values of alpha = 0.001 The log loss is: 1.2215086234804406  
 For values of alpha = 0.01 The log loss is: 1.328720183487397

For values of alpha = 0.1 The log loss is: 1.4291311026378817  
 For values of alpha = 1 The log loss is: 1.4706833951979281



For values of best alpha = 0.0001 The train log loss is: 0.9990278493040442  
 For values of best alpha = 0.0001 The cross validation log loss is: 1.1791330653831467  
 For values of best alpha = 0.0001 The test log loss is: 1.1724574298406922

**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 [28]:

```
print("Q6. How many data points in Test and CV datasets are covered by the ", unique_genes)

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 243 genes in train dataset?

Ans

1. In test data 651 out of 665 : 97.89473684210527
2. In cross validation data 518 out of 532 : 97.36842105263158

### 3.2.2 Univariate Analysis on Variation Feature

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

**Ans.** Variation is a categorical variable

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

In [29]:

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

Number of Unique Variations : 1922

|                      |    |
|----------------------|----|
| Truncating_Mutations | 60 |
| Amplification        | 54 |
| Deletion             | 47 |
| Fusions              | 23 |
| Overexpression       | 5  |
| Q61R                 | 3  |
| Q61L                 | 3  |
| TMPRSS2-ETV1_Fusion  | 2  |
| P130S                | 2  |
| R173C                | 2  |

Name: Variation, dtype: int64

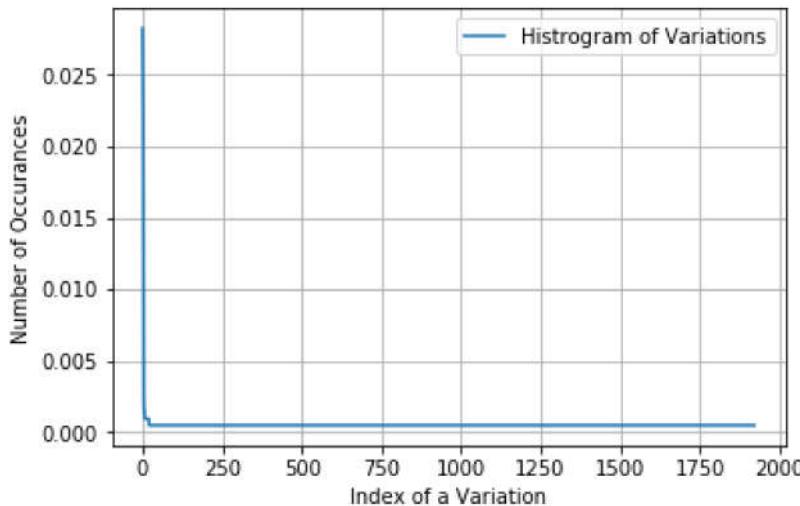
In [30]:

```
print("Ans: There are", unique_variations.shape[0] , "different categories of variations in")
```

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

In [31]:

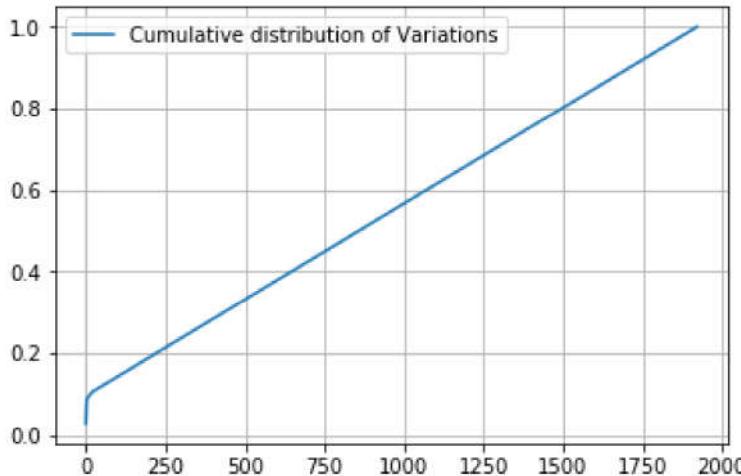
```
s = sum(unique_variations.values);
h = unique_variations.values/s;
plt.plot(h, label="Histogram of Variations")
plt.xlabel('Index of a Variation')
plt.ylabel('Number of Occurrences')
plt.legend()
plt.grid()
plt.show()
```



In [32]:

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

[0.02824859 0.05367232 0.07580038 ... 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/> (<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 [33]:

```
# alpha is used for Laplace smoothing
alpha = 1
# train gene feature
train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation", train_
# 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 [34]:

```
print("train_variation_feature_responseCoding is a converted feature using the response cod
```

train\_variation\_feature\_responseCoding is a converted feature using the response coding method. The shape of Variation feature: (2124, 9)

In [35]:

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

In [36]:

```
print("train_variation_feature_onehotEncoded is converted feature using the onne-hot encodi
```

```
train_variation_feature_onehotEncoded is converted feature using the onne-ho
t encoding method. The shape of Variation feature: (2124, 1958)
```

**Q10.** How good is this Variation feature in predicting y\_i?

Let's build a model just like the earlier!

In [37]:

```

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

# read more about SGDClassifier() at http://scikit-Learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal'
# 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 Desce
# 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=)

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_
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:"
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_

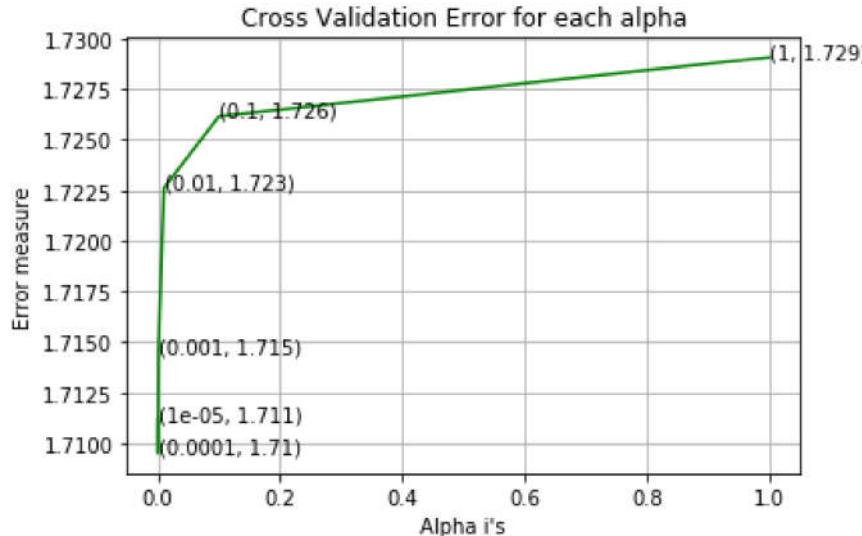
```

For values of alpha = 1e-05 The log loss is: 1.7111485336443009  
 For values of alpha = 0.0001 The log loss is: 1.7095044941194737  
 For values of alpha = 0.001 The log loss is: 1.7145224311737974

For values of alpha = 0.01 The log loss is: 1.7226140046074891

For values of alpha = 0.1 The log loss is: 1.7261573984424987

For values of alpha = 1 The log loss is: 1.7290792753568696



For values of best alpha = 0.0001 The train log loss is: 0.7825931880670998

For values of best alpha = 0.0001 The cross validation log loss is: 1.7095044941194737

For values of best alpha = 0.0001 The test log loss is: 1.714682948943186

**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 [38]:

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

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

Ans

1. In test data 66 out of 665 : 9.924812030075188
2. In cross validation data 51 out of 532 : 9.586466165413533

### 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 [39]:

```
# 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 [40]:

```
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)+len(total_dict)))+1)
            text_feature_responseCoding[row_index][i] = math.exp(sum_prob/len(row['TEXT'].split()))
            row_index += 1
    return text_feature_responseCoding
```

In [41]:

train\_df['TEXT'].head(4)

Out[41]:

```
1707 ppm1d wip1 negatively regulates dephosphorylat...
252 epidermal growth factor receptor egfr overexpr...
1638 abstract approach identify human oncogenes gen...
2371 protein tyrosine phosphatases ptps key positiv...
Name: TEXT, dtype: object
```

In [42]:

```

def top_tfidf_feats(row, features, top_n=25):
    ''' Get top n tfidf values in row and return them with their corresponding feature name
    topn_ids = np.argsort(row)[::-1][:top_n]
    top_feats = [(features[i], row[i]) for i in topn_ids]
    df = pd.DataFrame(top_feats)
    df.columns = ['feature', 'tfidf']
    return df

def top_mean_feats(Xtr, features, min_tfidf=0.1, grp_ids=None, top_n=25):
    ''' Return the top n features that on average are most important amongst documents in r
       indentified by indices in grp_ids. '''
    if grp_ids:
        D = Xtr[grp_ids].toarray()
    else:
        D = Xtr.toarray()

    D[D < min_tfidf] = 0
    tfidf_means = np.mean(D, axis=0)
    return top_tfidf_feats(tfidf_means, features, top_n)

```

In [43]:

```

#building a CountVectorizer with all the words that occurred minimum times in train data
text_vectorizer = TfidfVectorizer(min_df=3)
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 rows)
train_textfea_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_textfea_counts))
print("Total number of unique words in train data :", len(train_text_features))

```

Total number of unique words in train data : 53084

In [44]:

train\_textfea\_counts

Out[44]:

```
array([8.50542282, 8.86859573, 0.02160533, ..., 0.02864126, 0.03927349,
       0.07306987])
```

In [45]:

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

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

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

In [46]:

```
#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 [47]:

```
# 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 [48]:

```
# 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 [49]:

```
#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 [50]:

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

In [51]:

```
# 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,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal',
# 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_))

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)

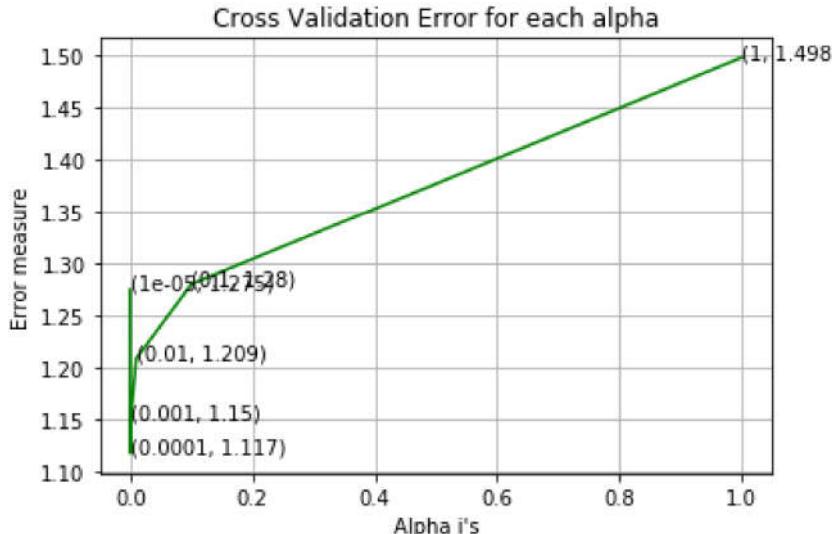
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_cv, predict_y))
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y))
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_cv, predict_y))
```

For values of alpha = 1e-05 The log loss is: 1.2750369752789361  
 For values of alpha = 0.0001 The log loss is: 1.117498533121457  
 For values of alpha = 0.001 The log loss is: 1.1498026734310902

For values of alpha = 0.01 The log loss is: 1.208641149955801

For values of alpha = 0.1 The log loss is: 1.2800110087797332

For values of alpha = 1 The log loss is: 1.4977782747371609



For values of best alpha = 0.0001 The train log loss is: 0.6517753508893039

For values of best alpha = 0.0001 The cross validation log loss is: 1.117498533121457

For values of best alpha = 0.0001 The test log loss is: 1.0599211094859267

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

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

In [52]:

```
def get_intersec_text(df):
    df_text_vec = TfidfVectorizer(min_df=3)
    df_textfea = df_text_vec.fit_transform(df['TEXT'])

    df_text_features = df_text_vec.get_feature_names()

    df_textfea_counts = df_textfea.sum(axis=0).A1
    df_textfea_dict = dict(zip(list(df_text_features),df_textfea_counts))
    len1 = len(set(df_text_features))
    len2 = len(set(train_text_features) & set(df_text_features))
    return len1,len2
```

In [53]:

```
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")
```

96.519 % of word of test data appeared in train data

98.549 % of word of Cross Validation appeared in train data

## 4. Machine Learning Models

In [54]:

```
#Data preparation for ML models.

#Misc. functionns for ML models

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

    # for calculating log_loss we will provide the array of probabilities belongs to each
    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.size)
    plot_confusion_matrix(test_y, pred_y)
```

In [55]:

```
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 [56]:

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

    gene_vec = gene_count_vec.fit(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,ye
        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(wc
        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,ye

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

## Stacking the three types of features

In [57]:

```
# merging gene, variance and text features

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

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

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

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

cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding)).toarray()
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 [58]:

```
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, 55284)
(number of data points * number of features) in test data = (665, 55284)
(number of data points * number of features) in cross validation data = (53
2, 55284)
```

In [59]:

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

Response encoding features :

```
(number of data points * number of features) in train data = (2124, 27)
(number of data points * number of features) in test data = (665, 27)
(number of data points * number of features) in cross validation data = (53
2, 27)
```

## 4.1. Base Line Model

### 4.1.1. Naive Bayes

#### 4.1.1.1. Hyper parameter tuning

In [60]:

```
# find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/
# -----
# 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
# -----


# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/gen
# -----
# 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
# -----


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

```

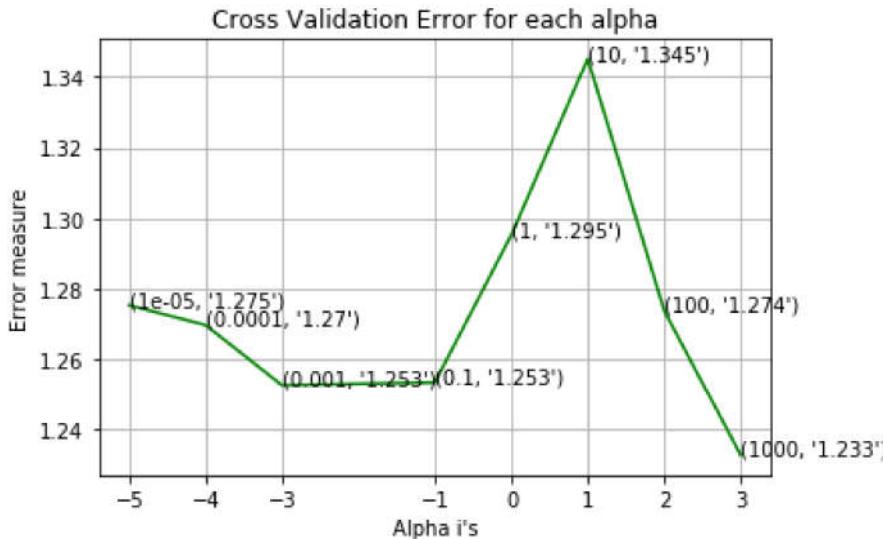
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:"
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_

```

```

for alpha = 1e-05
Log Loss : 1.2752622356177004
for alpha = 0.0001
Log Loss : 1.269664395462353
for alpha = 0.001
Log Loss : 1.2526392943859541
for alpha = 0.01
Log Loss : 1.253422477372519
for alpha = 0.1
Log Loss : 1.295172377428481
for alpha = 1
Log Loss : 1.3451704291963553
for alpha = 10
Log Loss : 1.2737632202886557
for alpha = 100
Log Loss : 1.232804010031772

```



```

For values of best alpha = 1000 The train log loss is: 0.9217499566643759
For values of best alpha = 1000 The cross validation log loss is: 1.2328040
10031772
For values of best alpha = 1000 The test log loss is: 1.194405811390313

```

#### 4.1.1.2. Testing the model with best hyper parameters

In [61]:

```
# find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules
# -----
# 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
# -----
```

```
# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/gen
# -----
# 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
# -----
```

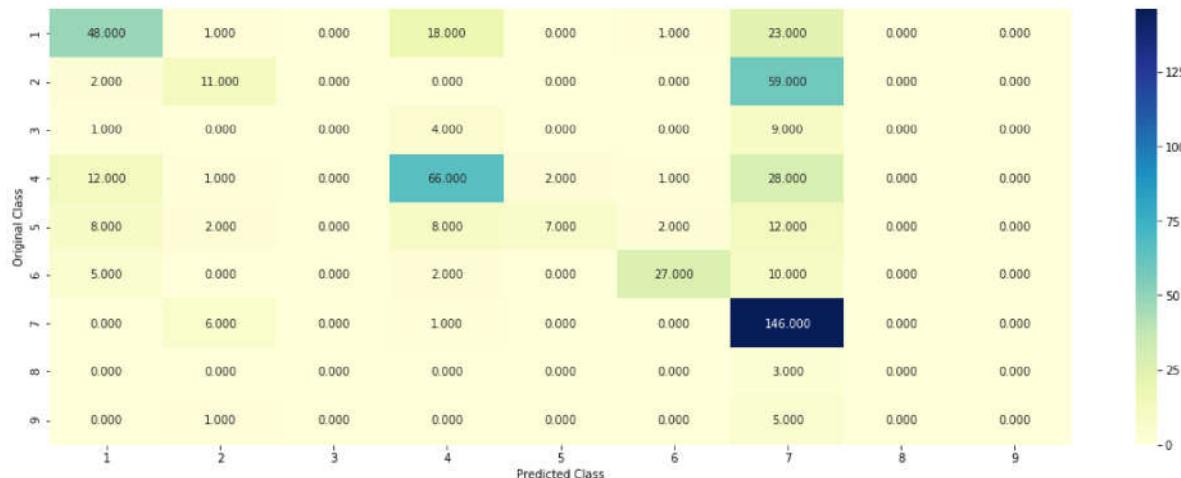
  

```
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
# to avoid rounding error while multiplying probabilités 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)))
plot_confusion_matrix(cv_y, sig_clf.predict(cv_x_onehotCoding.toarray()))
```

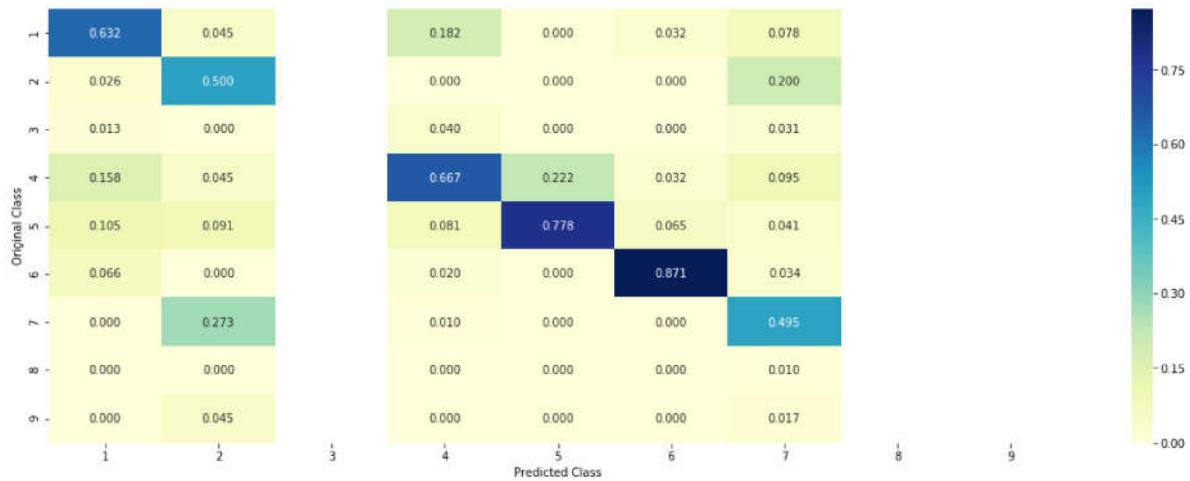
Log Loss : 1.232804010031772

Number of missclassified point : 0.4266917293233083

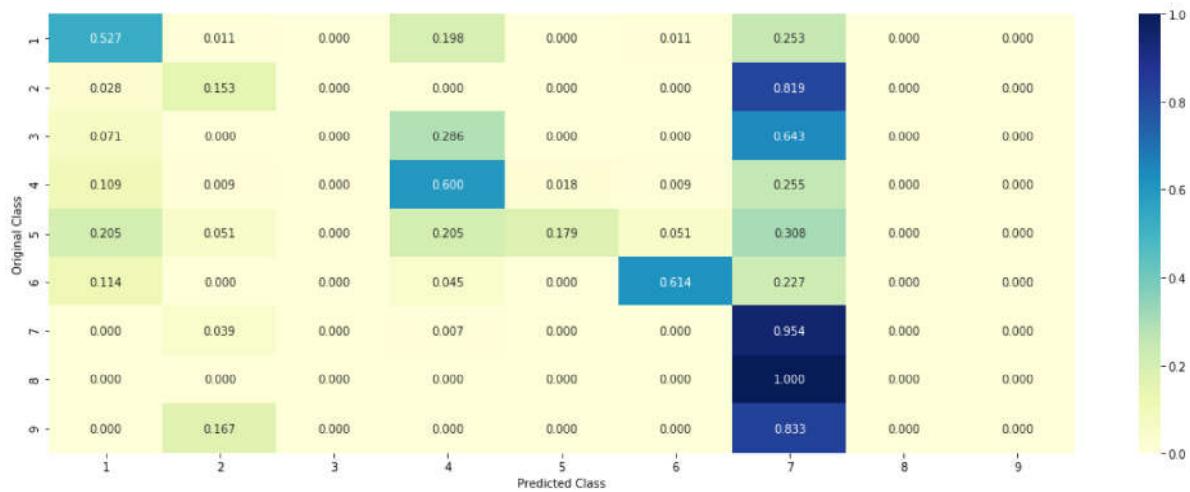
----- Confusion matrix -----



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



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



#### 4.1.1.3. Feature Importance, Correctly classified point

In [62]:

```
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[0]), 3))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imptfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[0])
```

```
Predicted Class : 7
Predicted Class Probabilities: [[5.690e-02 6.300e-02 1.900e-03 4.800e-02 9.1
50e-02 2.130e-02 7.164e-01
6.000e-04 4.000e-04]]
Actual Class : 5
```

```
-----  
14 Text feature [cells] present in test data point [True]  
15 Text feature [downstream] present in test data point [True]  
17 Text feature [kinase] present in test data point [True]  
18 Text feature [cell] present in test data point [True]  
19 Text feature [activation] present in test data point [True]  
23 Text feature [phosphorylation] present in test data point [True]  
25 Text feature [contrast] present in test data point [True]  
27 Text feature [expressing] present in test data point [True]  
28 Text feature [shown] present in test data point [True]  
29 Text feature [also] present in test data point [True]  
30 Text feature [inhibitor] present in test data point [True]  
31 Text feature [signaling] present in test data point [True]  
32 Text feature [growth] present in test data point [True]  
33 Text feature [however] present in test data point [True]  
34 Text feature [suggest] present in test data point [True]  
35 Text feature [independent] present in test data point [True]  
37 Text feature [previously] present in test data point [True]  
38 Text feature [increased] present in test data point [True]  
39 Text feature [compared] present in test data point [True]  
40 Text feature [10] present in test data point [True]  
41 Text feature [addition] present in test data point [True]  
42 Text feature [similar] present in test data point [True]  
43 Text feature [found] present in test data point [True]  
45 Text feature [1a] present in test data point [True]  
46 Text feature [higher] present in test data point [True]  
47 Text feature [mechanism] present in test data point [True]  
48 Text feature [potential] present in test data point [True]  
49 Text feature [enhanced] present in test data point [True]  
50 Text feature [figure] present in test data point [True]  
52 Text feature [well] present in test data point [True]  
53 Text feature [activating] present in test data point [True]  
54 Text feature [mutations] present in test data point [True]  
56 Text feature [mutant] present in test data point [True]  
57 Text feature [tyrosine] present in test data point [True]  
58 Text feature [3b] present in test data point [True]  
59 Text feature [consistent] present in test data point [True]  
61 Text feature [treatment] present in test data point [True]  
62 Text feature [described] present in test data point [True]  
66 Text feature [interestingly] present in test data point [True]  
67 Text feature [may] present in test data point [True]  
68 Text feature [using] present in test data point [True]  
69 Text feature [demonstrated] present in test data point [True]
```

72 Text feature [various] present in test data point [True]  
75 Text feature [inhibitors] present in test data point [True]  
76 Text feature [followed] present in test data point [True]  
77 Text feature [lines] present in test data point [True]  
78 Text feature [including] present in test data point [True]  
79 Text feature [proliferation] present in test data point [True]  
81 Text feature [3a] present in test data point [True]  
82 Text feature [furthermore] present in test data point [True]  
84 Text feature [report] present in test data point [True]  
86 Text feature [reported] present in test data point [True]  
87 Text feature [observed] present in test data point [True]  
90 Text feature [mutation] present in test data point [True]  
91 Text feature [expression] present in test data point [True]  
93 Text feature [two] present in test data point [True]  
94 Text feature [identified] present in test data point [True]  
95 Text feature [obtained] present in test data point [True]  
96 Text feature [increase] present in test data point [True]  
97 Text feature [examined] present in test data point [True]  
Out of the top 100 features 60 are present in query point

#### 4.1.1.4. Feature Importance, Incorrectly classified point

In [63]:

```
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[0]), 3))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imptfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[0])
```

```
Predicted Class : 4
Predicted Class Probabilities: [[1.136e-01 9.000e-03 3.790e-02 6.822e-01 5.9
80e-02 4.940e-02 4.750e-02
 4.000e-04 2.000e-04]]
Actual Class : 6
```

```
-----  
10 Text feature [proteins] present in test data point [True]  
11 Text feature [activity] present in test data point [True]  
12 Text feature [protein] present in test data point [True]  
14 Text feature [experiments] present in test data point [True]  
15 Text feature [loss] present in test data point [True]  
16 Text feature [acid] present in test data point [True]  
18 Text feature [determined] present in test data point [True]  
19 Text feature [whereas] present in test data point [True]  
20 Text feature [whether] present in test data point [True]  
21 Text feature [shown] present in test data point [True]  
23 Text feature [indicated] present in test data point [True]  
25 Text feature [described] present in test data point [True]  
26 Text feature [function] present in test data point [True]  
27 Text feature [missense] present in test data point [True]  
28 Text feature [two] present in test data point [True]  
29 Text feature [important] present in test data point [True]  
30 Text feature [results] present in test data point [True]  
35 Text feature [type] present in test data point [True]  
36 Text feature [amino] present in test data point [True]  
38 Text feature [bind] present in test data point [True]  
39 Text feature [also] present in test data point [True]  
40 Text feature [mutations] present in test data point [True]  
41 Text feature [functions] present in test data point [True]  
42 Text feature [expressed] present in test data point [True]  
45 Text feature [vitro] present in test data point [True]  
48 Text feature [wild] present in test data point [True]  
50 Text feature [thus] present in test data point [True]  
52 Text feature [indicate] present in test data point [True]  
53 Text feature [levels] present in test data point [True]  
54 Text feature [purified] present in test data point [True]  
55 Text feature [partially] present in test data point [True]  
56 Text feature [reduced] present in test data point [True]  
57 Text feature [suppressor] present in test data point [True]  
58 Text feature [related] present in test data point [True]  
59 Text feature [transfected] present in test data point [True]  
60 Text feature [determine] present in test data point [True]  
61 Text feature [containing] present in test data point [True]  
62 Text feature [tagged] present in test data point [True]  
63 Text feature [may] present in test data point [True]  
64 Text feature [either] present in test data point [True]  
65 Text feature [see] present in test data point [True]  
66 Text feature [lower] present in test data point [True]
```

```
67 Text feature [three] present in test data point [True]
68 Text feature [analyzed] present in test data point [True]
70 Text feature [although] present in test data point [True]
71 Text feature [30] present in test data point [True]
73 Text feature [standard] present in test data point [True]
74 Text feature [mm] present in test data point [True]
75 Text feature [effects] present in test data point [True]
83 Text feature [except] present in test data point [True]
86 Text feature [made] present in test data point [True]
87 Text feature [critical] present in test data point [True]
90 Text feature [result] present in test data point [True]
92 Text feature [using] present in test data point [True]
93 Text feature [suggest] present in test data point [True]
Out of the top 100 features 55 are present in query point
```

## 4.2. K Nearest Neighbour Classification

### 4.2.1. Hyper parameter tuning

In [64]:

```
# find more about KNeighborsClassifier() here http://scikit-Learn.org/stable/modules/general.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
#-----


# find more about CalibratedClassifierCV here at http://scikit-Learn.org/stable/modules/general.html
# -----
# default parameters
# 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 probabilites we use log-probability estimation
    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)
```

```

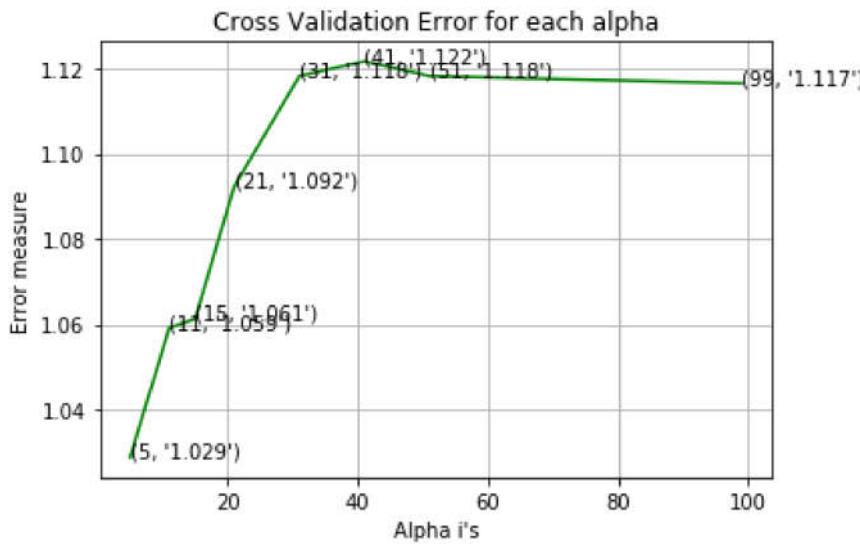
predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:"
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_

```

```

for alpha = 5
Log Loss : 1.0287617628155055
for alpha = 11
Log Loss : 1.0591593928899152
for alpha = 15
Log Loss : 1.0612819158060707
for alpha = 21
Log Loss : 1.0923449041197113
for alpha = 31
Log Loss : 1.1182542439970928
for alpha = 41
Log Loss : 1.121808887531119
for alpha = 51
Log Loss : 1.1182624339938598
for alpha = 99
Log Loss : 1.1165451019710568

```



For values of best alpha = 5 The train log loss is: 0.47188162335748696  
 For values of best alpha = 5 The cross validation log loss is: 1.0287617628  
 155055  
 For values of best alpha = 5 The test log loss is: 1.0301971131023804

#### 4.2.2. Testing the model with best hyper parameters

In [65]:

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

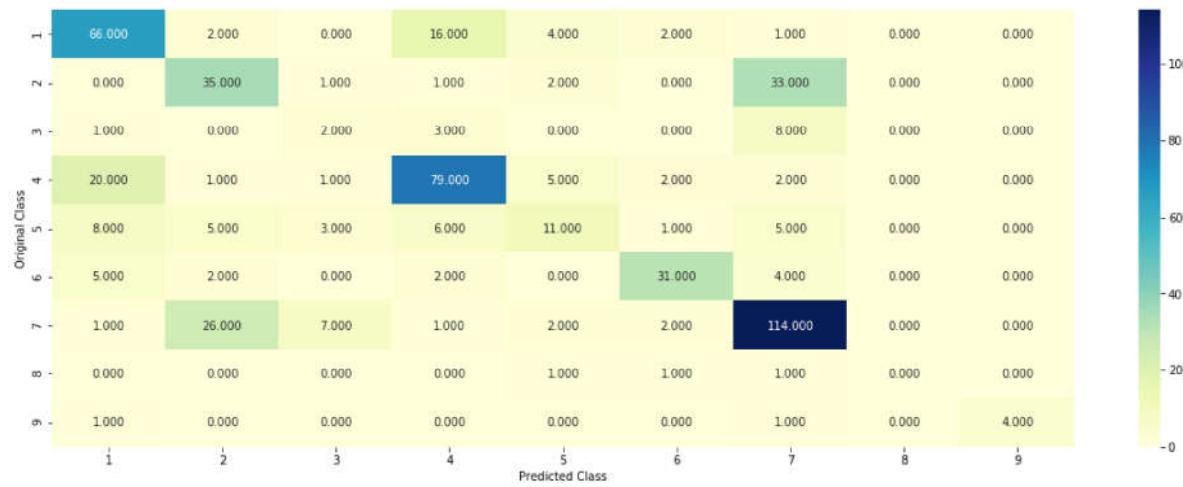
# methods of
# fit(X, y) : Fit the model using X as training data and y as target values
# predict(X):Predict the class labels for the provided data
# predict_proba(X):Return probability estimates for the test data X.
# -----
# video Link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nea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

```

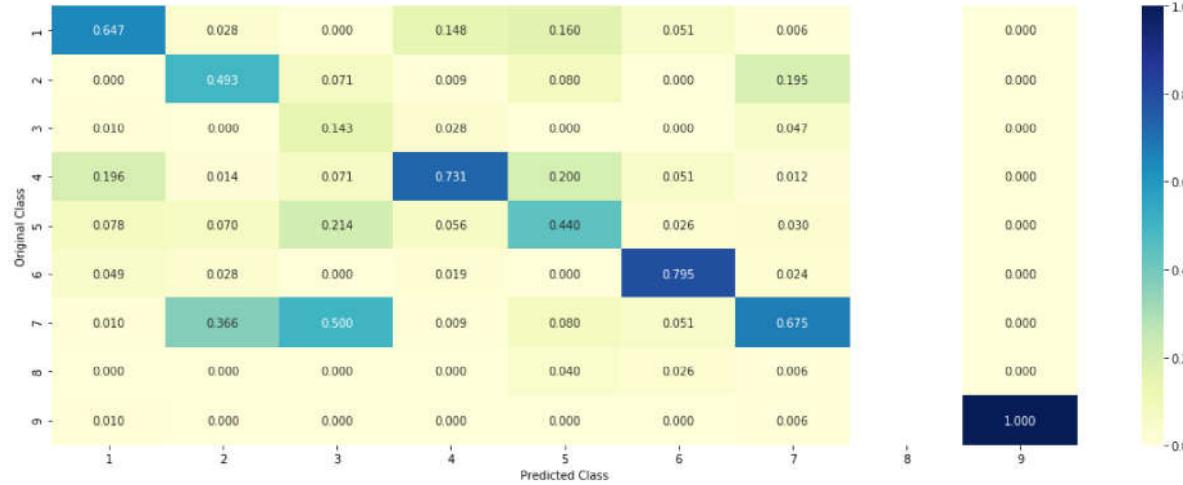
Log loss : 1.0287617628155055

Number of mis-classified points : 0.35714285714285715

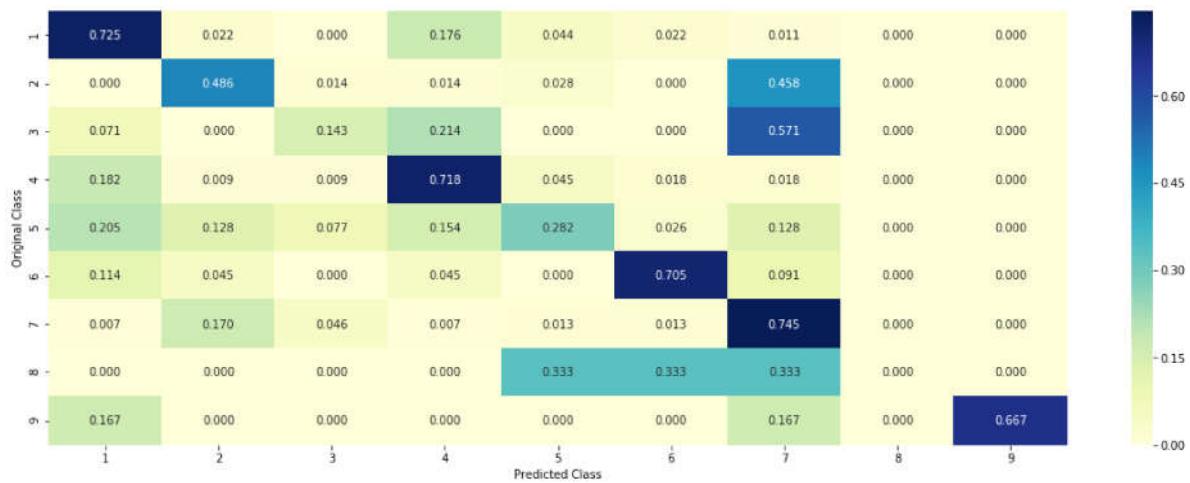
----- Confusion matrix -----



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



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



#### 4.2.3. Sample Query point -1

In [66]:

```
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[be
print("The ",alpha[best_alpha]," nearest neighbours of the test points belongs to classes",
print("Frequency of nearest points : ",Counter(train_y[neighbors[1][0]])))
```

Predicted Class : 7

Actual Class : 5

The 5 nearest neighbours of the test points belongs to classes [1 2 7 5 5]

Frequency of nearest points : Counter({5: 2, 1: 1, 2: 1, 7: 1})

#### 4.2.4. Sample Query Point-2

In [67]:

```
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[be
print("the k value for knn is",alpha[best_alpha],"and the nearest neighbours of the test po
print("Frequency of nearest points :",Counter(train_y[neighbors[1][0]])))
```

```
Predicted Class : 4
Actual Class : 6
the k value for knn is 5 and the nearest neighbours of the test points belon
gs to classes [3 4 3 4 4]
Frequency of nearest points : Counter({4: 3, 3: 2})
```

## 4.3. Logistic Regression

### 4.3.1. With Class balancing

#### 4.3.1.1. Hyper parameter tuning

In [68]:

```
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal'
# 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/geom
#-----

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/gen
# -----
# 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_
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use Log-probability estimat
    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='l
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
```

```

sig_clf.fit(train_x_onehotCoding, train_y)

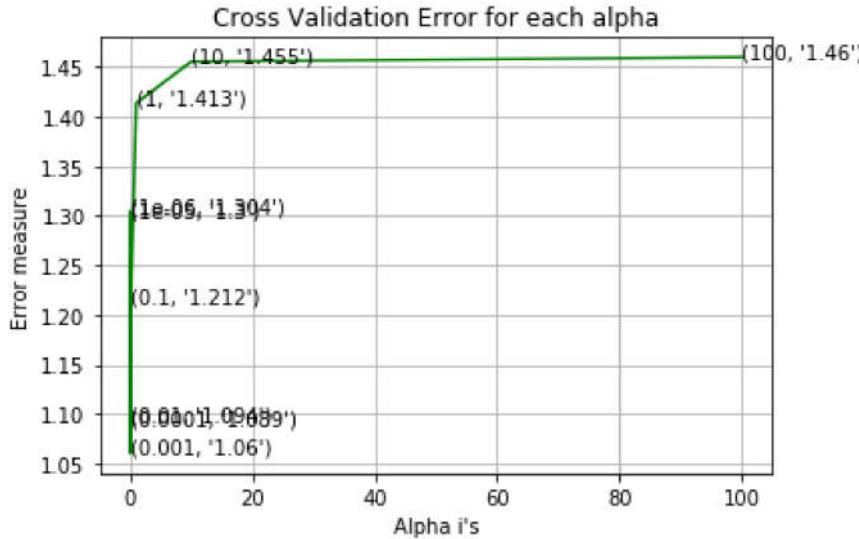
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:"
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_

```

```

for alpha = 1e-06
Log Loss : 1.3042815707526412
for alpha = 1e-05
Log Loss : 1.2998182904486122
for alpha = 0.0001
Log Loss : 1.0894281825665182
for alpha = 0.001
Log Loss : 1.0601647505037206
for alpha = 0.01
Log Loss : 1.093890891832136
for alpha = 0.1
Log Loss : 1.2115862539200877
for alpha = 1
Log Loss : 1.413384028448876
for alpha = 10
Log Loss : 1.4552120861811846
for alpha = 100
Log Loss : 1.460117203415138

```



For values of best alpha = 0.001 The train log loss is: 0.5215763162424594  
 For values of best alpha = 0.001 The cross validation log loss is: 1.060164  
 7505037206  
 For values of best alpha = 0.001 The test log loss is: 1.0071251506052734

#### 4.3.1.2. Testing the model with best hyper parameters

In [69]:

```
# read more about SGDClassifier() at http://scikit-Learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal',
# class_weight=None, warm_start=False, average=False, n_iter=None)

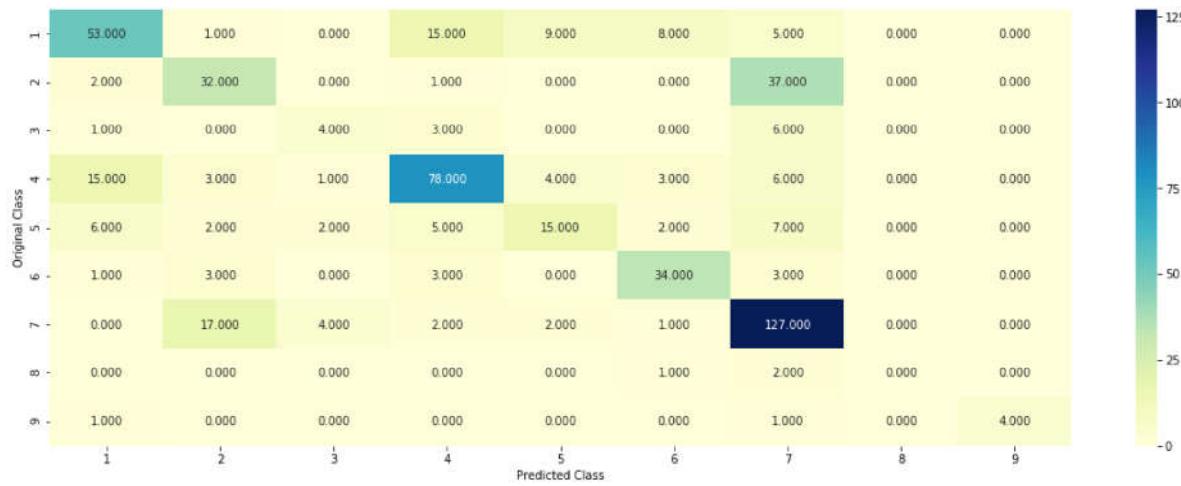
# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit Linear model with Stochastic Gradient Desce
# predict(X) Predict class labels for samples in X.

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/geom
# -----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='l
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, c
```

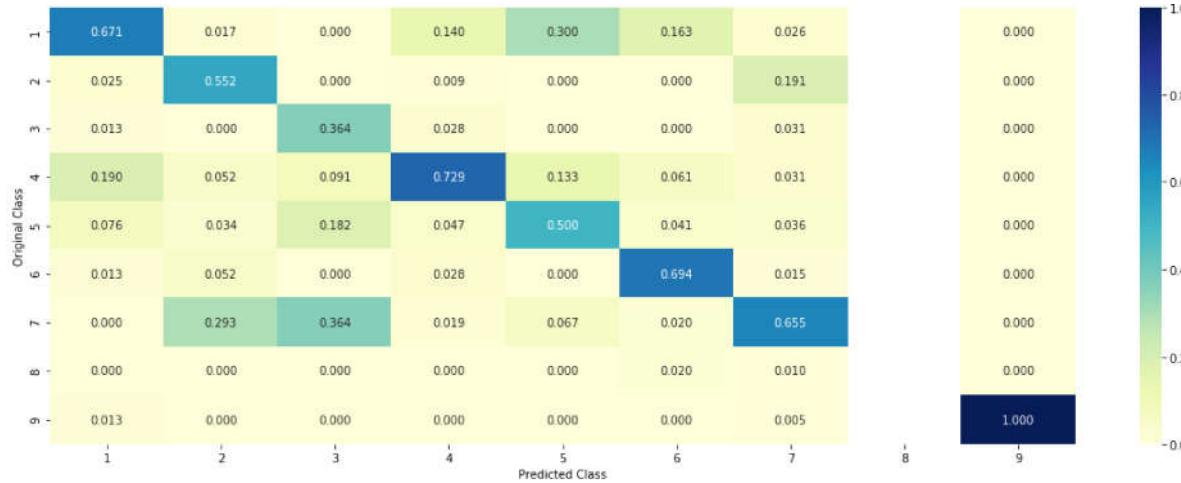
Log loss : 1.0601647505037206

Number of mis-classified points : 0.34774436090225563

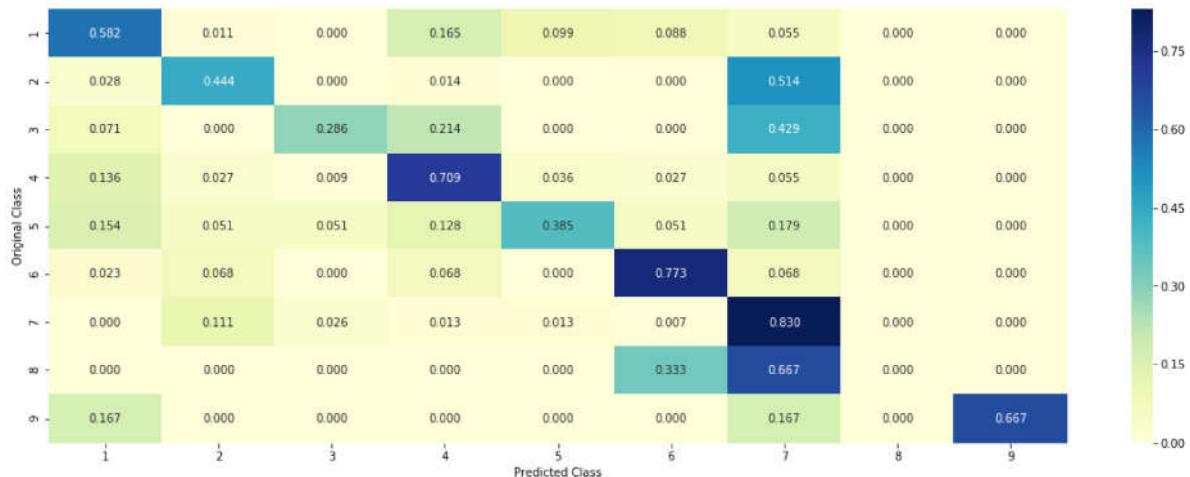
----- Confusion matrix -----



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



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



#### 4.3.1.3. Feature Importance

In [70]:

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

##### 4.3.1.3.1. Correctly Classified point

In [71]:

```
# from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='')
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[0]), 3))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imptfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[0])
```

```
Predicted Class : 5
Predicted Class Probabilities: [[0.0753 0.1092 0.0033 0.0083 0.5285 0.0122
0.2596 0.0015 0.002 ]]
Actual Class : 5
-----
270 Text feature [figure22] present in test data point [True]
Out of the top 500 features 1 are present in query point
```

#### 4.3.1.3.2. Incorrectly Classified point

In [72]:

```
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[
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imptfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc
```

```
Predicted Class : 4
Predicted Class Probabilities: [[2.400e-03 2.400e-03 4.854e-01 5.001e-01 3.3
00e-03 5.000e-03 8.000e-04
3.000e-04 3.000e-04]]
Actual Class : 6
-----
44 Text feature [cys61gly] present in test data point [True]
126 Text feature [fluorescing] present in test data point [True]
129 Text feature [gcuccucucacucuuagu] present in test data point [True]
130 Text feature [jdp] present in test data point [True]
131 Text feature [t37k] present in test data point [True]
132 Text feature [uaaaauuggacauaaggaguccucc] present in test data point [True]
133 Text feature [dr13] present in test data point [True]
134 Text feature [asterisked] present in test data point [True]
135 Text feature [ay273801] present in test data point [True]
136 Text feature [c27a] present in test data point [True]
137 Text feature [cbas] present in test data point [True]
152 Text feature [pdr] present in test data point [True]
155 Text feature [osu] present in test data point [True]
164 Text feature [dilemma] present in test data point [True]
168 Text feature [decrements] present in test data point [True]
208 Text feature [fourths] present in test data point [True]
230 Text feature [suppressor] present in test data point [True]
232 Text feature [immunoaffinity] present in test data point [True]
300 Text feature [adapted] present in test data point [True]
301 Text feature [nonsense] present in test data point [True]
307 Text feature [magnitudes] present in test data point [True]
335 Text feature [claes] present in test data point [True]
338 Text feature [truncate] present in test data point [True]
363 Text feature [ionizing] present in test data point [True]
377 Text feature [monolayers] present in test data point [True]
392 Text feature [deficiencies] present in test data point [True]
397 Text feature [germline] present in test data point [True]
428 Text feature [302] present in test data point [True]
495 Text feature [neoplasias] present in test data point [True]
Out of the top 500 features 29 are present in query point
```

### 4.3.2. Without Class balancing

#### 4.3.2.1. Hyper parameter tuning

In [73]:

```
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal',
# 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 Desce
# predict(X) Predict class labels for samples in X.

#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/geom
#-----
```

```
# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/gen
# -----
# 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)
```

```

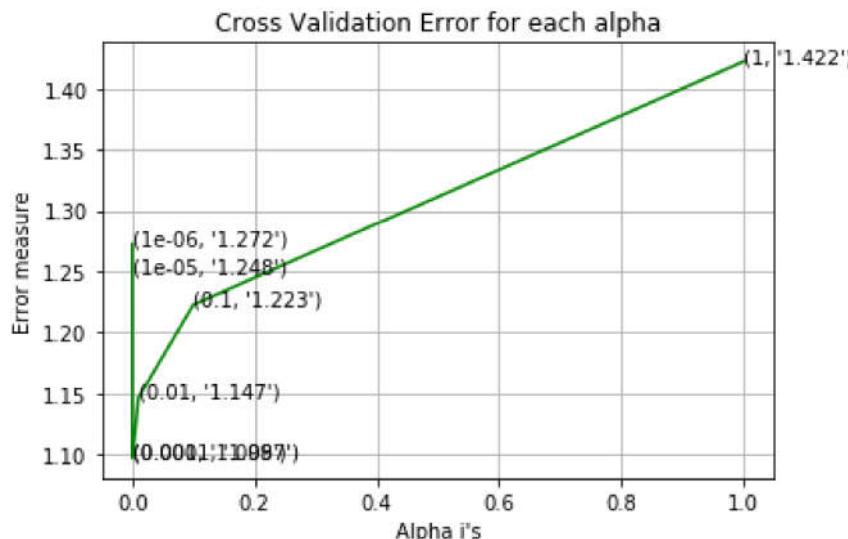
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:"
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_

```

```

for alpha = 1e-06
Log Loss : 1.2721582907725741
for alpha = 1e-05
Log Loss : 1.2482753265185487
for alpha = 0.0001
Log Loss : 1.0973509608941066
for alpha = 0.001
Log Loss : 1.0980745627243476
for alpha = 0.01
Log Loss : 1.146519582987
for alpha = 0.1
Log Loss : 1.2226600891460393
for alpha = 1
Log Loss : 1.422323509881528

```



For values of best alpha = 0.0001 The train log loss is: 0.516129615400089  
 For values of best alpha = 0.0001 The cross validation log loss is: 1.09735  
 09608941066  
 For values of best alpha = 0.0001 The test log loss is: 1.0187249180666924

#### 4.3.2.2. Testing model with best hyper parameters

In [74]:

```
# read more about SGDClassifier() at http://scikit-Learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal',
# 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 Desce
# predict(X) Predict class labels for samples in X.

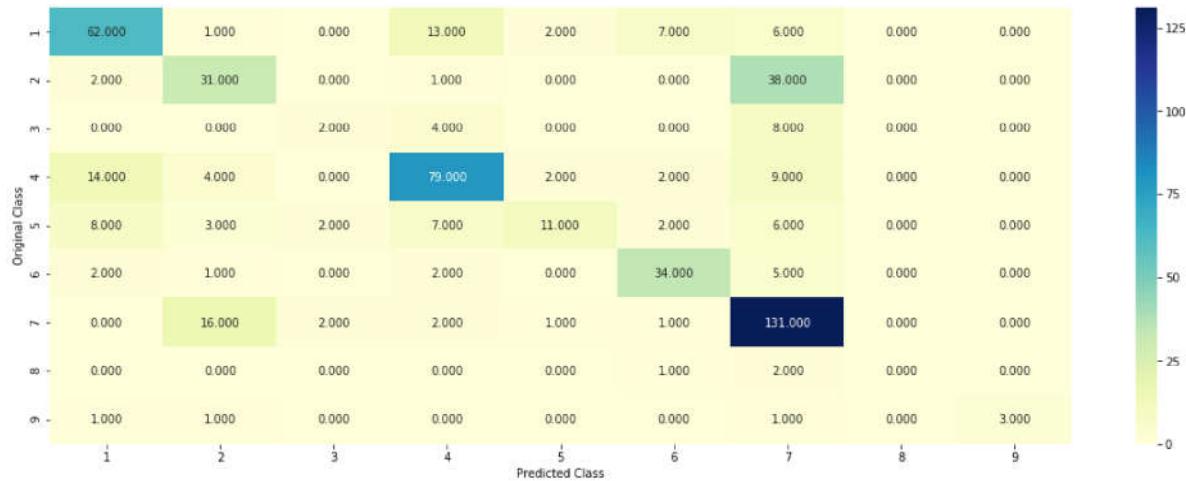
#-----
# video link:
#-----
```

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, c

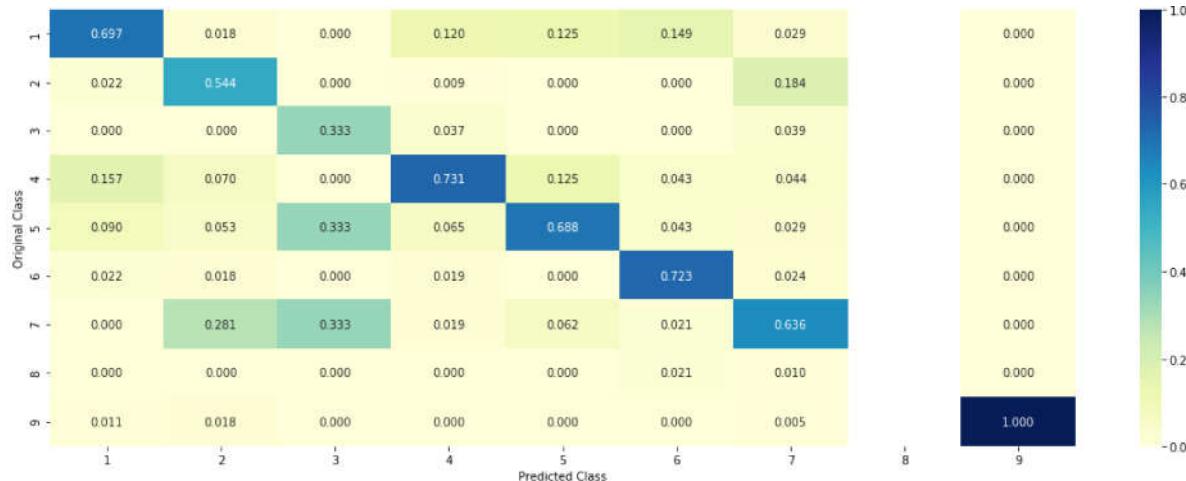
Log loss : 1.0973509608941066

Number of mis-classified points : 0.33646616541353386

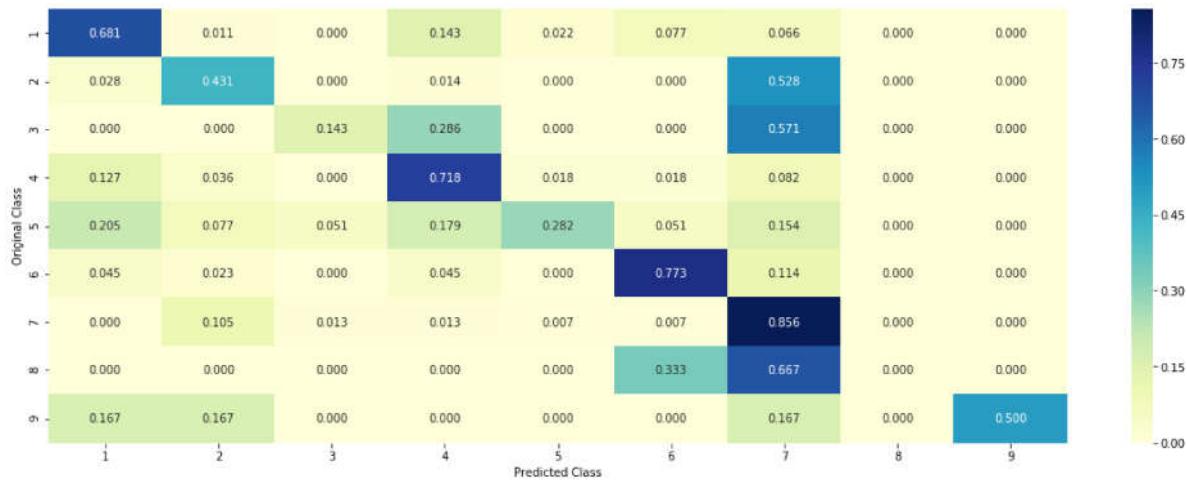
----- Confusion matrix -----



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



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



#### 4.3.2.3. Feature Importance, Correctly Classified point

In [75]:

```
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[
print("Actual Class : ", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc
```

```

Predicted Class : 7
Predicted Class Probabilities: [[0.1321 0.1186 0.0024 0.0142 0.2964 0.0123
0.4204 0.0023 0.0014]]
Actual Class : 5
-----
239 Text feature [transformation] present in test data point [True]
256 Text feature [activation] present in test data point [True]
274 Text feature [phosphorylation] present in test data point [True]
282 Text feature [downstream] present in test data point [True]
283 Text feature [overexpression] present in test data point [True]
390 Text feature [phospho] present in test data point [True]
473 Text feature [responding] present in test data point [True]
Out of the top 500 features 7 are present in query point
```

#### 4.3.2.4. Feature Importance, Incorrectly Classified point

In [76]:

```
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[0]), 3))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imptfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc[test_point_index])

```

Predicted Class : 4  
Predicted Class Probabilities: [[4.700e-03 1.040e-02 3.682e-01 5.998e-01  
3.100e-03 3.300e-03 9.300e-03  
9.000e-04 4.000e-04]]  
Actual Class : 6

---

157 Text feature [cys61gly] present in test data point [True]  
164 Text feature [cbas] present in test data point [True]  
165 Text feature [c27a] present in test data point [True]  
166 Text feature [asterisked] present in test data point [True]  
167 Text feature [dr13] present in test data point [True]  
168 Text feature [gcuccucucacucuucagu] present in test data point [True]  
169 Text feature [t37k] present in test data point [True]  
170 Text feature [jdp] present in test data point [True]  
171 Text feature [uaaaauuggacauaaggagucccc] present in test data point [True]  
172 Text feature [ay273801] present in test data point [True]  
179 Text feature [pdr] present in test data point [True]  
186 Text feature [fluorescing] present in test data point [True]

## 4.4. Linear Support Vector Machines

### 4.4.1. Hyper parameter tuning

In [77]:

```
# read more about support vector machines with Linear kernels here http://scikit-learn.org/
# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovo')
# -----
# 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/mathematics-for-machine-learning-linear-algebra
# -----
```

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

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

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

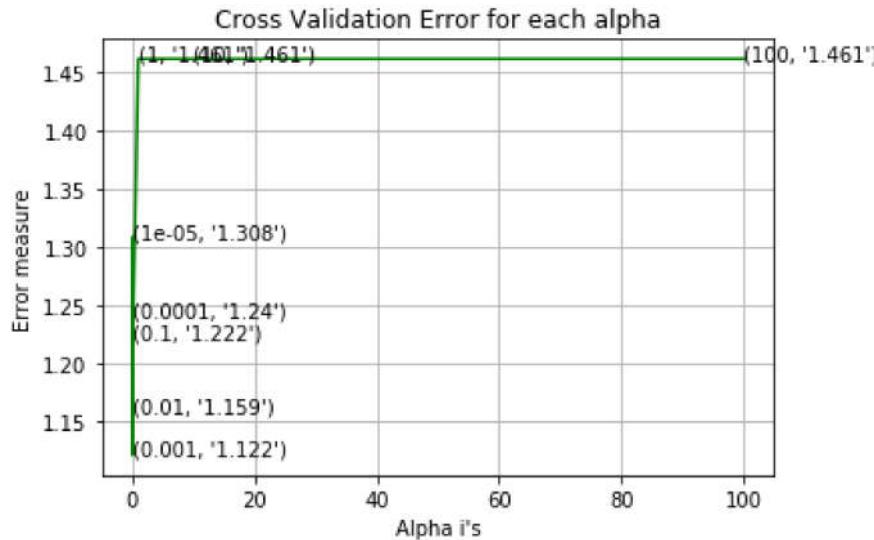
```
best_alpha = np.argmin(cv_log_error_array)
# clf = SVC(C=i, kernel='linear', probability=True, class_weight='balanced')
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='hinge')
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")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

```

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_
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:"
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_
```

for C = 1e-05  
Log Loss : 1.308425202952882  
for C = 0.0001  
Log Loss : 1.2401831076206464  
for C = 0.001  
Log Loss : 1.1217251801897135  
for C = 0.01  
Log Loss : 1.1585454249452924  
for C = 0.1  
Log Loss : 1.2221690495249855  
for C = 1  
Log Loss : 1.4611226682070682  
for C = 10  
Log Loss : 1.4610616839040127  
for C = 100  
Log Loss : 1.4610609215510388



For values of best alpha = 0.001 The train log loss is: 0.5671580307946241  
For values of best alpha = 0.001 The cross validation log loss is: 1.121725  
1801897135  
For values of best alpha = 0.001 The test log loss is: 1.034955350080306

#### 4.4.2. Testing model with best hyper parameters

In [78]:

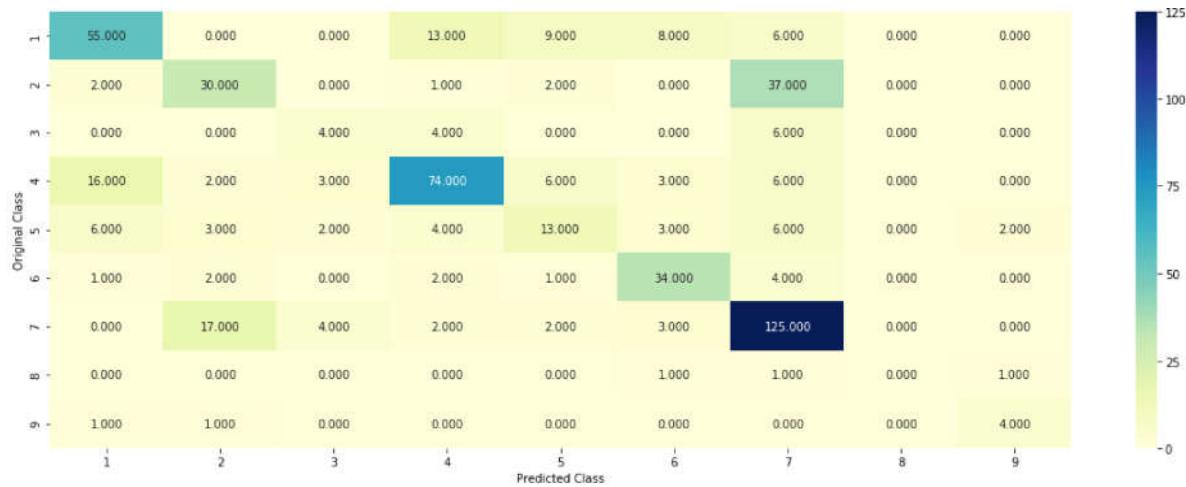
```
# read more about support vector machines with Linear kernels here http://scikit-learn.org/
# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovo')
# -----
# 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/mathematics-for-machine-learning-svm
# -----
```

# `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)`

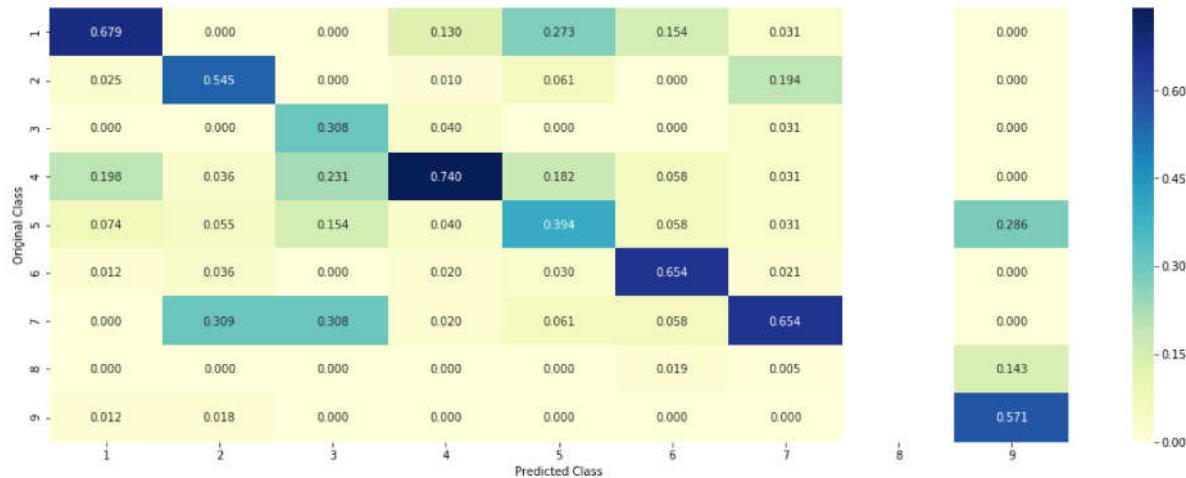
Log loss : 1.1217251801897135

Number of mis-classified points : 0.36278195488721804

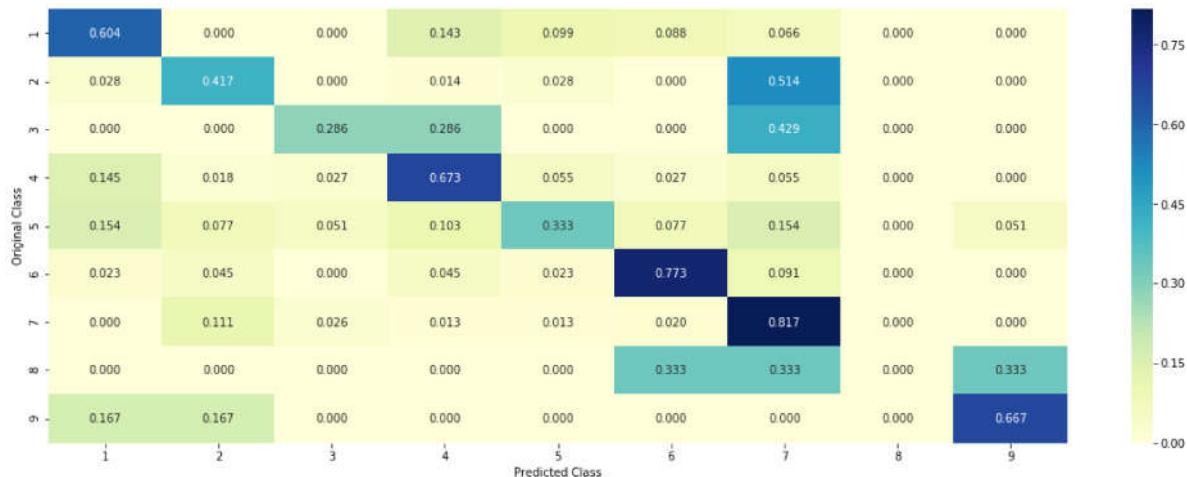
----- Confusion matrix -----



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



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



### 4.3.3. Feature Importance

#### 4.3.3.1. For Correctly classified point

In [79]:

```
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[
print("Actual Class : ", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gene'].iloc
```

```
Predicted Class : 5
Predicted Class Probabilities: [[0.099  0.1166  0.0127  0.078  0.5022  0.0277
0.1522 0.003  0.0086]]
Actual Class : 5
-----
297 Text feature [figure22] present in test data point [True]
307 Text feature [outcomes] present in test data point [True]
442 Text feature [limitations] present in test data point [True]
Out of the top 500 features 3 are present in query point
```

#### 4.3.3.2. For Incorrectly classified point

In [80]:

```
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[
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imptfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0156 0.02    0.4597 0.4739 0.0022 0.0033
0.0201 0.0015 0.0036]]
Actual Class : 6
```

```
-----
134 Text feature [c27a] present in test data point [True]
135 Text feature [cbas] present in test data point [True]
136 Text feature [t37k] present in test data point [True]
137 Text feature [ay273801] present in test data point [True]
138 Text feature [gcuccucucacucuucagu] present in test data point [True]
139 Text feature [uaaaauuugggacauaaggaguccucc] present in test data point [True]
e]
140 Text feature [asterisked] present in test data point [True]
141 Text feature [jdp] present in test data point [True]
142 Text feature [dr13] present in test data point [True]
144 Text feature [pdr] present in test data point [True]
156 Text feature [dilemma] present in test data point [True]
160 Text feature [fluorescing] present in test data point [True]
162 Text feature [decrements] present in test data point [True]
168 Text feature [fourths] present in test data point [True]
209 Text feature [osu] present in test data point [True]
211 Text feature [magnitudes] present in test data point [True]
218 Text feature [cys61gly] present in test data point [True]
257 Text feature [adapted] present in test data point [True]
298 Text feature [immunoaffinity] present in test data point [True]
311 Text feature [i21v] present in test data point [True]
399 Text feature [recombined] present in test data point [True]
410 Text feature [302] present in test data point [True]
443 Text feature [ionizing] present in test data point [True]
444 Text feature [1292] present in test data point [True]
452 Text feature [monolayers] present in test data point [True]
453 Text feature [guiding] present in test data point [True]
484 Text feature [798] present in test data point [True]
Out of the top 500 features 27 are present in query point
```

## 4.5 Random Forest Classifier

### 4.5.1. Hyper parameter tuning (With One hot Encoding)

In [81]:

```

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


# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/gen
# -----
# 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_
        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-9))
        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_a
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_c  
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:")  
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)  
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The cross validation lo  
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:",  
  
for n_estimators = 100 and max depth = 5  
Log Loss : 1.1927150731927265  
for n_estimators = 100 and max depth = 10  
Log Loss : 1.167893323525949  
for n_estimators = 200 and max depth = 5  
Log Loss : 1.1809069441947004  
for n_estimators = 200 and max depth = 10  
Log Loss : 1.1506606079557182  
for n_estimators = 500 and max depth = 5  
Log Loss : 1.1664947074027054  
for n_estimators = 500 and max depth = 10  
Log Loss : 1.1446997931033098  
for n_estimators = 1000 and max depth = 5  
Log Loss : 1.171247356853659  
for n_estimators = 1000 and max depth = 10  
Log Loss : 1.1451307725487907  
for n_estimators = 2000 and max depth = 5  
Log Loss : 1.169172266154061  
for n_estimators = 2000 and max depth = 10  
Log Loss : 1.1455971669668354  
For values of best estimator = 500 The train log loss is: 0.650748846844593  
7  
For values of best estimator = 500 The cross validation log loss is: 1.1446  
997931033096  
For values of best estimator = 500 The test log loss is: 1.154552185659132
```

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

In [82]:

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

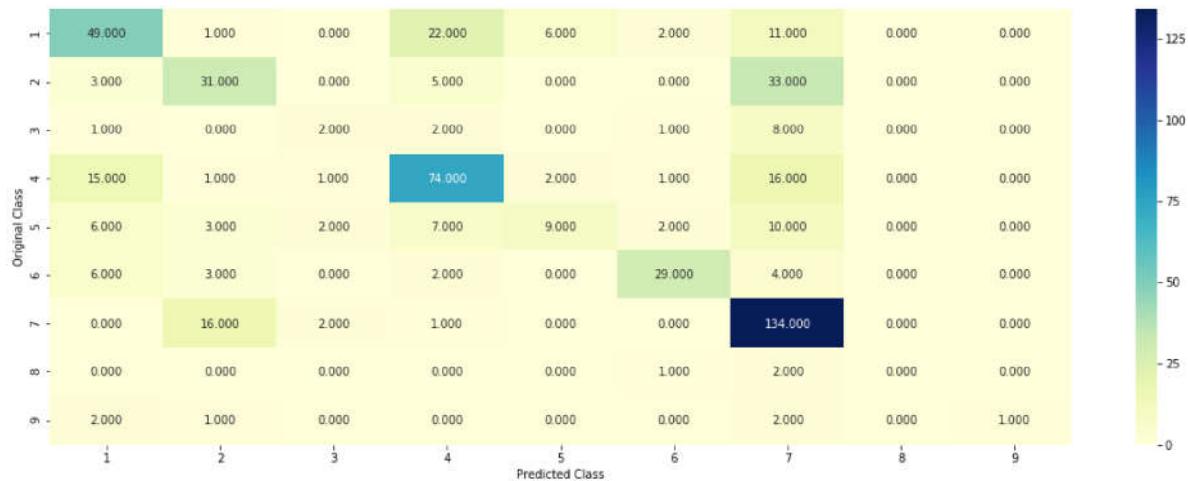
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/random-forest-classifier
# -----
```

```
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_c
```

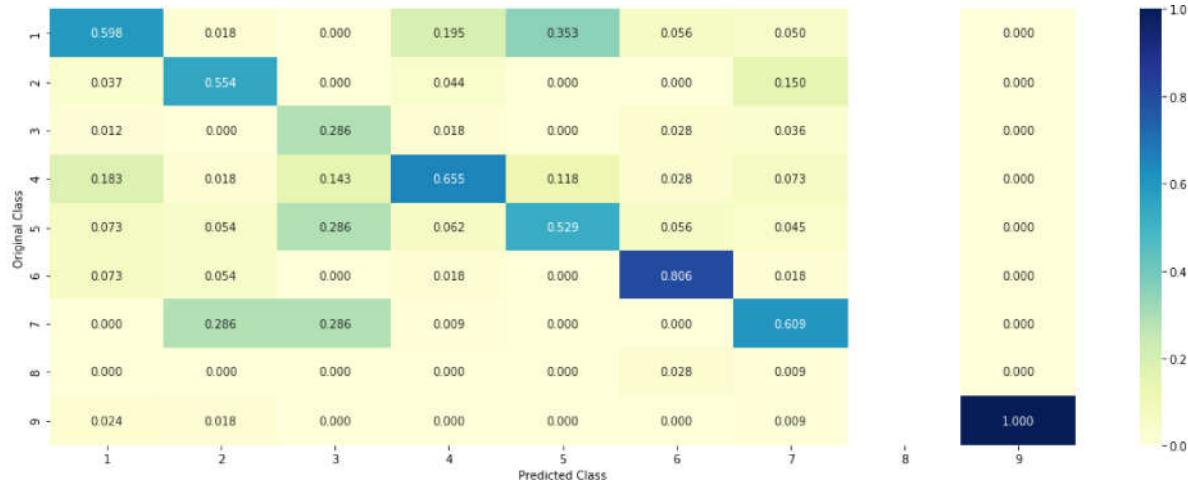
Log loss : 1.1446997931033098

Number of mis-classified points : 0.3815789473684211

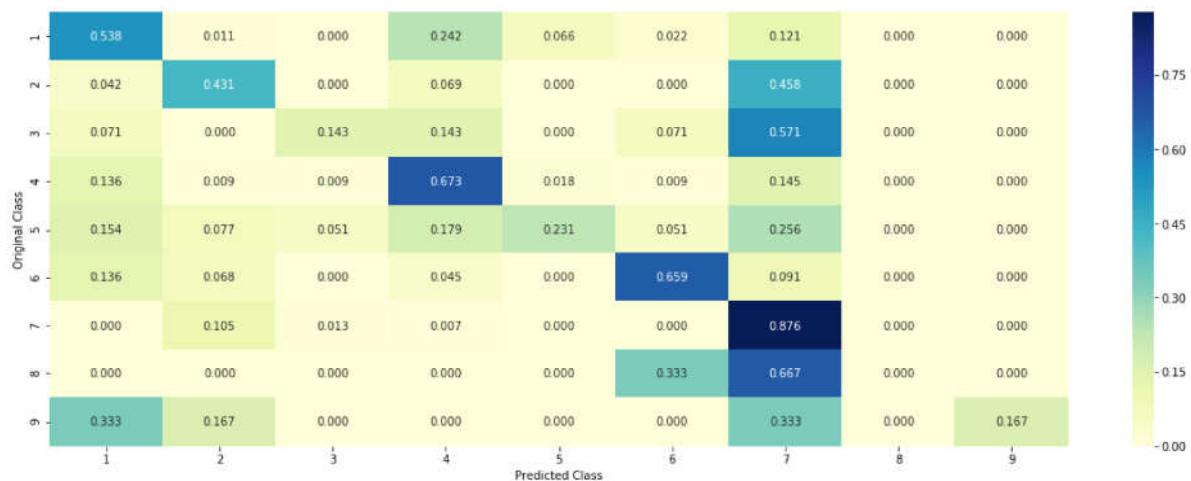
----- Confusion matrix -----



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



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



### 4.5.3. Feature Importance

#### 4.5.3.1. Correctly Classified point

In [83]:

```
# test_point_index = 10
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_c
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[
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_imptfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index], test_df['
```

```
Predicted Class : 7
Predicted Class Probabilities: [[0.1159 0.1735 0.0195 0.0716 0.1946 0.0582
0.3529 0.0051 0.0088]]
Actual Class : 5
```

```
-----
0 Text feature [kinase] present in test data point [True]
1 Text feature [activation] present in test data point [True]
2 Text feature [inhibitor] present in test data point [True]
3 Text feature [inhibitors] present in test data point [True]
6 Text feature [phosphorylation] present in test data point [True]
7 Text feature [tyrosine] present in test data point [True]
8 Text feature [treatment] present in test data point [True]
9 Text feature [activating] present in test data point [True]
10 Text feature [missense] present in test data point [True]
11 Text feature [oncogenic] present in test data point [True]
15 Text feature [resistance] present in test data point [True]
17 Text feature [signaling] present in test data point [True]
18 Text feature [function] present in test data point [True]
20 Text feature [growth] present in test data point [True]
21 Text feature [therapy] present in test data point [True]
22 Text feature [extracellular] present in test data point [True]
23 Text feature [drug] present in test data point [True]
27 Text feature [phospho] present in test data point [True]
28 Text feature [transforming] present in test data point [True]
29 Text feature [loss] present in test data point [True]
32 Text feature [receptor] present in test data point [True]
34 Text feature [cells] present in test data point [True]
35 Text feature [therapeutic] present in test data point [True]
36 Text feature [advanced] present in test data point [True]
37 Text feature [efficacy] present in test data point [True]
38 Text feature [downstream] present in test data point [True]
39 Text feature [cell] present in test data point [True]
40 Text feature [trials] present in test data point [True]
46 Text feature [pathogenic] present in test data point [True]
48 Text feature [variants] present in test data point [True]
51 Text feature [functional] present in test data point [True]
52 Text feature [effective] present in test data point [True]
53 Text feature [stability] present in test data point [True]
55 Text feature [lines] present in test data point [True]
57 Text feature [harboring] present in test data point [True]
58 Text feature [kinases] present in test data point [True]
59 Text feature [clinical] present in test data point [True]
```

```
61 Text feature [phosphorylated] present in test data point [True]
64 Text feature [protein] present in test data point [True]
65 Text feature [inactivation] present in test data point [True]
66 Text feature [variant] present in test data point [True]
69 Text feature [tkis] present in test data point [True]
71 Text feature [proliferation] present in test data point [True]
72 Text feature [expressing] present in test data point [True]
73 Text feature [amplification] present in test data point [True]
77 Text feature [likely] present in test data point [True]
82 Text feature [responses] present in test data point [True]
84 Text feature [assay] present in test data point [True]
86 Text feature [potential] present in test data point [True]
87 Text feature [sensitivity] present in test data point [True]
95 Text feature [ring] present in test data point [True]
97 Text feature [functions] present in test data point [True]
Out of the top 100 features 52 are present in query point
```

#### 4.5.3.2. Incorrectly Classified point

In [84]:

```
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[0]), 3))
print("Actual Class : ", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_imptfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index], test_df['Cancer Type'])
```

Predicted Class : 4  
Predicted Class Probabilities: [[0.1136 0.0387 0.2864 0.3797 0.0754 0.0685  
0.026 0.0046 0.0071]]

Actual Class : 6

---

1 Text feature [activation] present in test data point [True]  
4 Text feature [suppressor] present in test data point [True]  
10 Text feature [missense] present in test data point [True]  
14 Text feature [nonsense] present in test data point [True]  
15 Text feature [resistance] present in test data point [True]  
18 Text feature [function] present in test data point [True]  
24 Text feature [treated] present in test data point [True]  
28 Text feature [transforming] present in test data point [True]  
29 Text feature [loss] present in test data point [True]  
32 Text feature [receptor] present in test data point [True]  
33 Text feature [brca1] present in test data point [True]  
34 Text feature [cells] present in test data point [True]  
38 Text feature [downstream] present in test data point [True]  
39 Text feature [cell] present in test data point [True]  
46 Text feature [pathogenic] present in test data point [True]  
48 Text feature [variants] present in test data point [True]  
49 Text feature [classify] present in test data point [True]  
51 Text feature [functional] present in test data point [True]  
55 Text feature [lines] present in test data point [True]  
59 Text feature [clinical] present in test data point [True]  
64 Text feature [protein] present in test data point [True]  
66 Text feature [variant] present in test data point [True]  
70 Text feature [frameshift] present in test data point [True]  
72 Text feature [expressing] present in test data point [True]  
74 Text feature [mutants] present in test data point [True]  
75 Text feature [patients] present in test data point [True]  
76 Text feature [factor] present in test data point [True]  
77 Text feature [likely] present in test data point [True]  
79 Text feature [membranes] present in test data point [True]  
80 Text feature [deleterious] present in test data point [True]  
81 Text feature [defective] present in test data point [True]  
84 Text feature [assay] present in test data point [True]  
86 Text feature [potential] present in test data point [True]  
89 Text feature [database] present in test data point [True]  
95 Text feature [ring] present in test data point [True]  
96 Text feature [patient] present in test data point [True]  
97 Text feature [functions] present in test data point [True]  
98 Text feature [resistant] present in test data point [True]

Out of the top 100 features 38 are present in query point

#### 4.5.3. Hyper parameter tuning (With Response Coding)

In [85]:

```

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


# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/gen
# -----
# 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_
        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_a
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_c  
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  
predict_y = sig_clf.predict_proba(cv_x_responseCoding)  
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log lo  
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_
```

```
< [REDACTED] >  
for n_estimators = 10 and max depth = 2  
Log Loss : 1.9486622409397882  
for n_estimators = 10 and max depth = 3  
Log Loss : 1.4907936910889927  
for n_estimators = 10 and max depth = 5  
Log Loss : 1.5852770855785963  
for n_estimators = 10 and max depth = 10  
Log Loss : 1.6793600896699592  
for n_estimators = 50 and max depth = 2  
Log Loss : 1.7281234592499115  
for n_estimators = 50 and max depth = 3  
Log Loss : 1.3754404025704692  
for n_estimators = 50 and max depth = 5  
Log Loss : 1.3450240748584454  
for n_estimators = 50 and max depth = 10  
Log Loss : 1.5806884017329255  
for n_estimators = 100 and max depth = 2  
Log Loss : 1.6188344722355248  
for n_estimators = 100 and max depth = 3  
Log Loss : 1.42605906849819  
for n_estimators = 100 and max depth = 5  
Log Loss : 1.3071470065512318  
for n_estimators = 100 and max depth = 10  
Log Loss : 1.5962872636991463  
for n_estimators = 200 and max depth = 2  
Log Loss : 1.6397260550043664  
for n_estimators = 200 and max depth = 3  
Log Loss : 1.4200257625767485  
for n_estimators = 200 and max depth = 5  
Log Loss : 1.308526547861796  
for n_estimators = 200 and max depth = 10  
Log Loss : 1.6541913031563706  
for n_estimators = 500 and max depth = 2  
Log Loss : 1.6591453540682177  
for n_estimators = 500 and max depth = 3  
Log Loss : 1.5050668971902854  
for n_estimators = 500 and max depth = 5  
Log Loss : 1.3306255586784295  
for n_estimators = 500 and max depth = 10  
Log Loss : 1.680413522370788  
for n_estimators = 1000 and max depth = 2  
Log Loss : 1.610037876760534  
for n_estimators = 1000 and max depth = 3  
Log Loss : 1.4734393642535262  
for n_estimators = 1000 and max depth = 5  
Log Loss : 1.3218422287061997
```

```
for n_estimators = 1000 and max depth = 10
Log Loss : 1.648580122707394
For values of best alpha = 100 The train log loss is: 0.06425748890616
For values of best alpha = 100 The cross validation log loss is: 1.30714700
65512318
For values of best alpha = 100 The test log loss is: 1.317529213672818
```

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

In [86]:

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

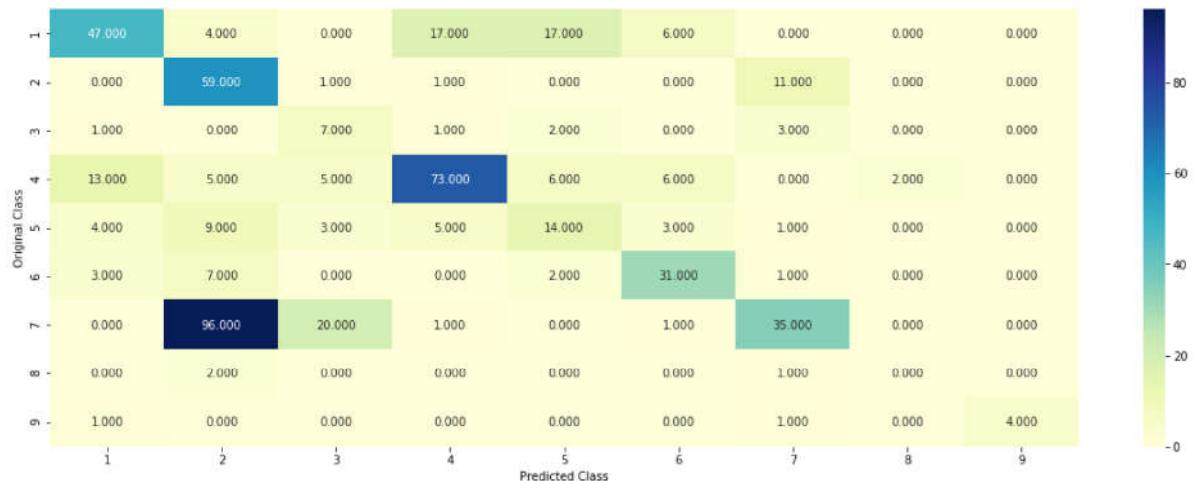
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/random-forest-classifier
# -----
```

```
clf = RandomForestClassifier(max_depth=max_depth[int(best_alpha%4)], n_estimators=alpha[int(best_alpha%4)])
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseCoding, cv_y,
```

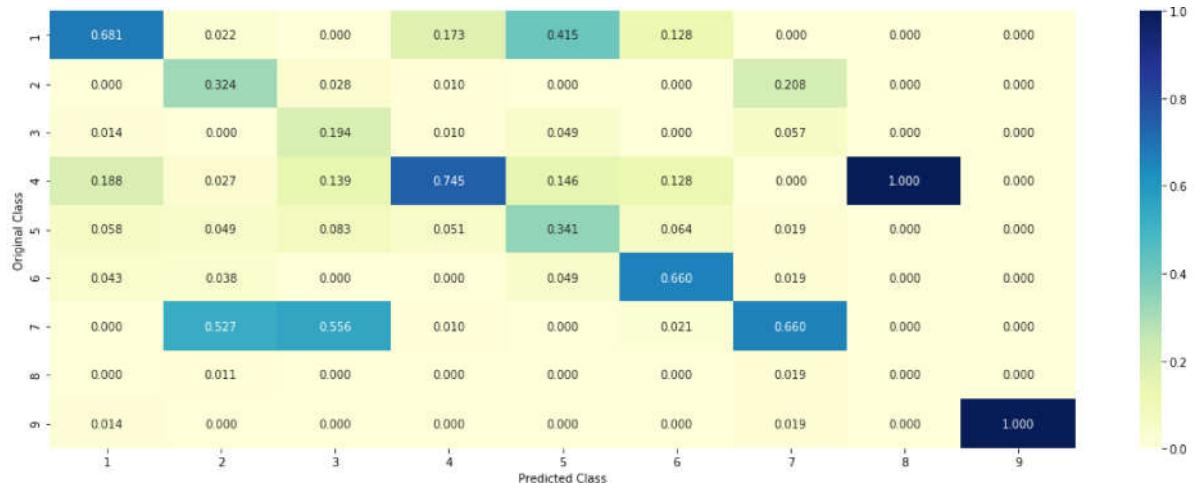
Log loss : 1.3071470065512318

Number of mis-classified points : 0.4924812030075188

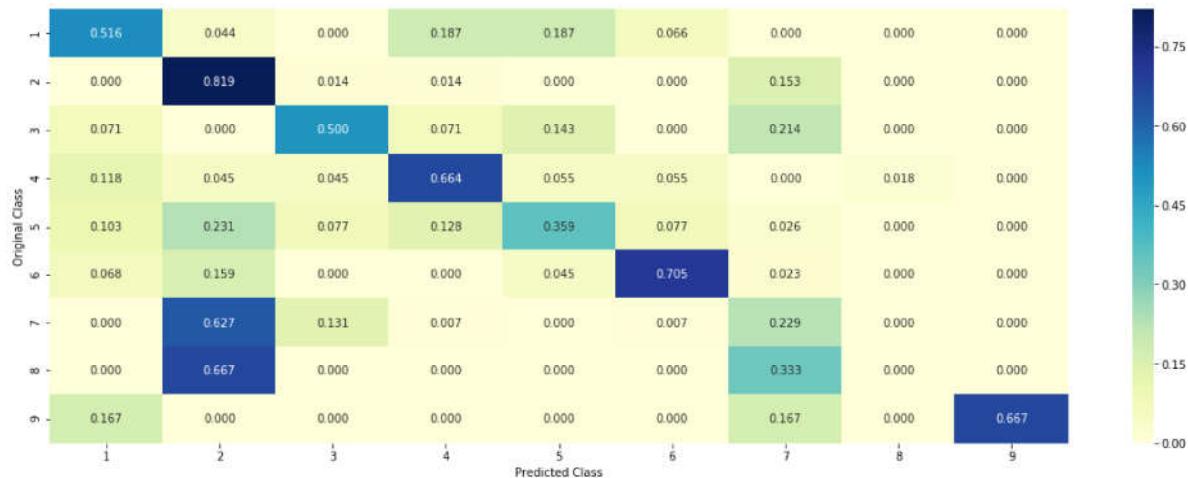
----- Confusion matrix -----



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



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



## 4.5.5. Feature Importance

### 4.5.5.1. Correctly Classified point

In [87]:

```
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini', max_c
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_responseCodir
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
for i in indices:
    if i<9:
        print("Gene is important feature")
    elif i<18:
        print("Variation is important feature")
    else:
        print("Text is important feature")
```

```
Predicted Class : 2
Predicted Class Probabilities: [[0.0163 0.4061 0.1024 0.0239 0.0863 0.0516
0.2856 0.0137 0.014 ]]
Actual Class : 5
```

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

#### 4.5.5.2. Incorrectly Classified point

In [88]:

```
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]), 3))
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 : 5  
Predicted Class Probabilities: [[0.0449 0.0065 0.2397 0.0818 0.4828 0.1291  
0.0041 0.0046 0.0065]]  
Actual Class : 6

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

## 4.7 Stack the models

### 4.7.1 testing with hyper parameter tuning

In [89]:

```
# read more about SGDClassifier() at http://scikit-Learn.org/stable/modules/generated/sklear
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='L2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal'
# 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 Desce
# predict(X) Predict class labels for samples in X.

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/geom
# -----


# read more about support vector machines with linear kernels here http://scikit-Learn.org/
# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=F
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='o

# 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/mathem
# -----


# read more about support vector machines with linear kernels here http://scikit-Learn.org/
# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, ve
# 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
# -----


clf1 = SGDClassifier(alpha=0.001, penalty='l2', loss='log', class_weight='balanced', random_s
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_s
clf2.fit(train_x_onehotCoding, train_y)
sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")
```

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

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

```
Logistic Regression : Log Loss: 1.06
Support vector machines : Log Loss: 1.46
Naive Bayes : Log Loss: 1.25
```

```
-----
Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.177
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.030
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.489
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.119
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.235
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.566
```

#### 4.7.2 testing the model with the best hyper parameters

In [90]:

```

lr = LogisticRegression(C=0.1)
sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr, n_jobs=-1)
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)))
plot_confusion_matrix(test_y=test_y, predict_y=sclf.predict(test_x_onehotCoding))

```

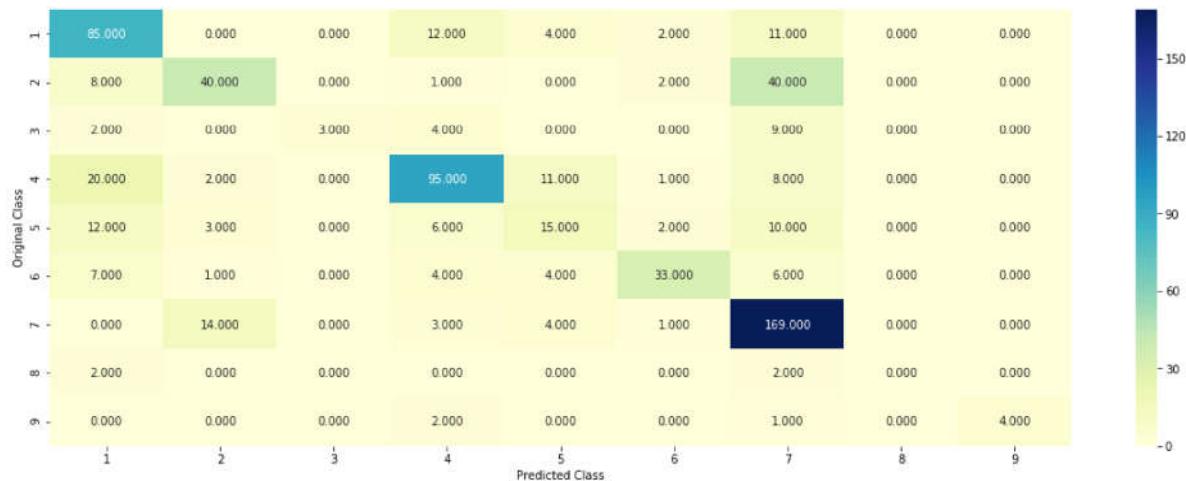
Log loss (train) on the stacking classifier : 0.610040141357214

Log loss (CV) on the stacking classifier : 1.1191471782009077

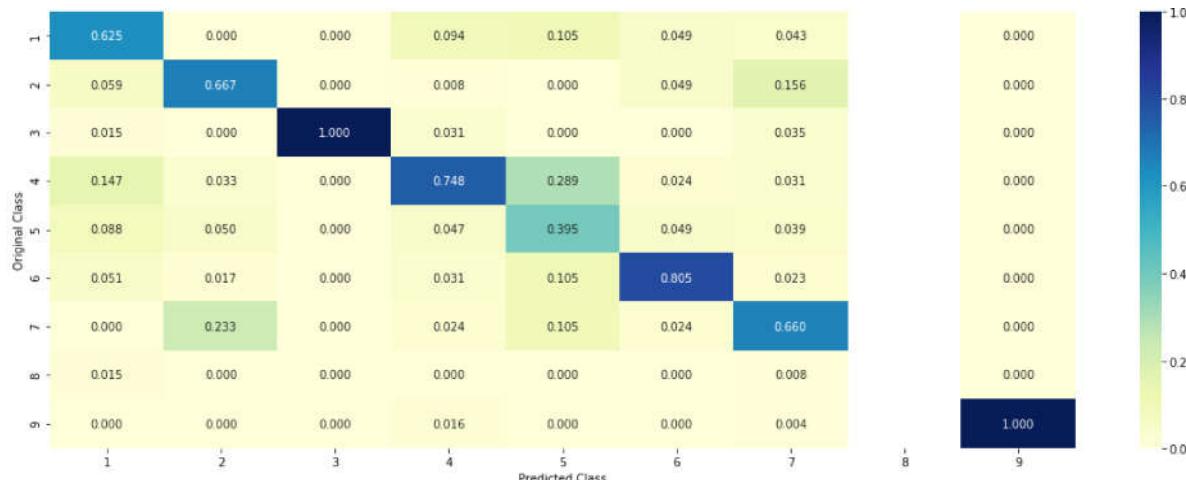
Log loss (test) on the stacking classifier : 1.084215456193397

Number of missclassified point : 0.3323308270676692

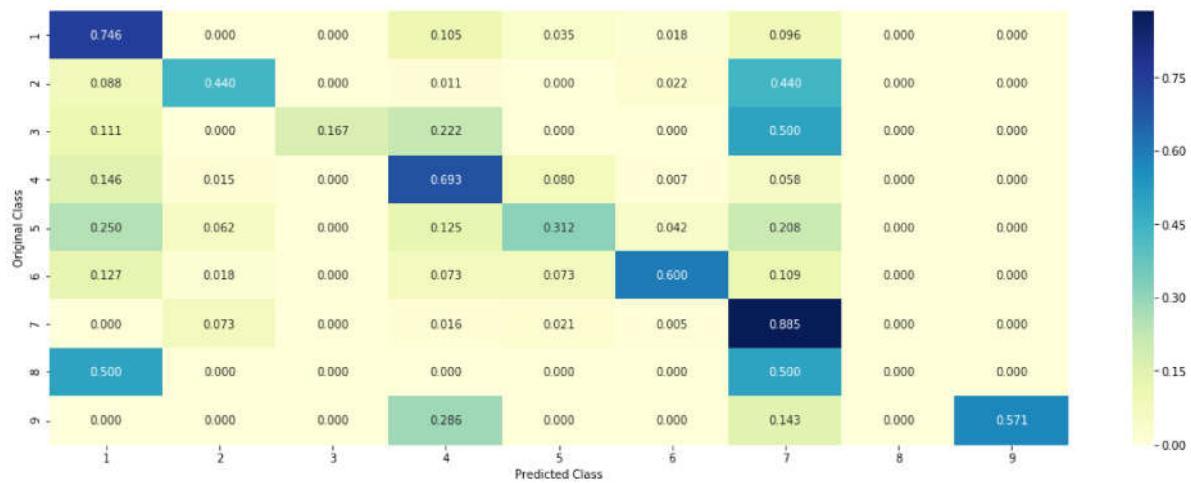
----- Confusion matrix -----



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



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



#### 4.7.3 Maximum Voting classifier

In [91]:

```
#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)], n_jobs=-1)
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) != 0))
plot_confusion_matrix(test_y=test_y, predict_y=vclf.predict(test_x_onehotCoding))
```

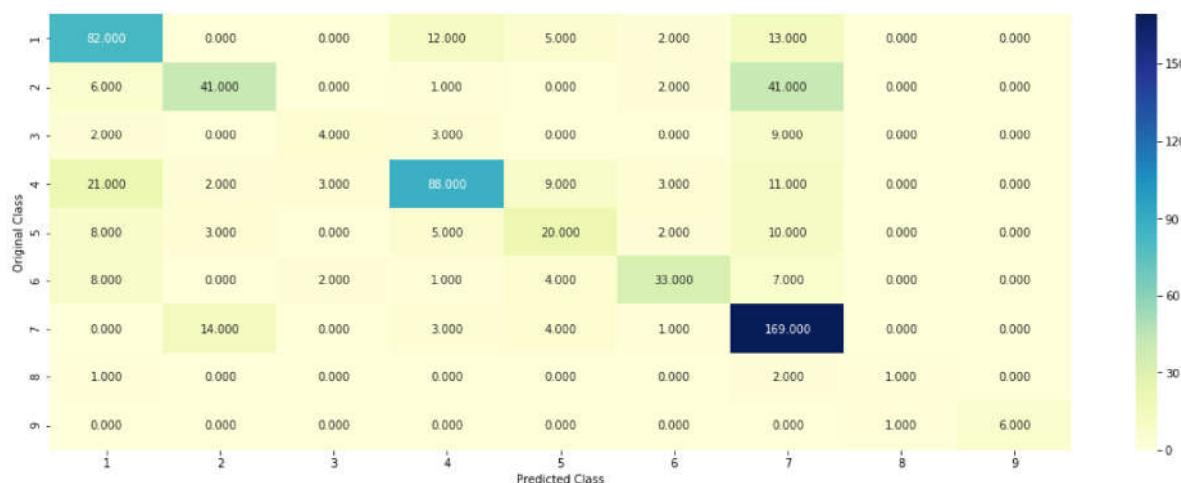
Log loss (train) on the VotingClassifier : 0.8282366920273875

Log loss (CV) on the VotingClassifier : 1.1263915231839183

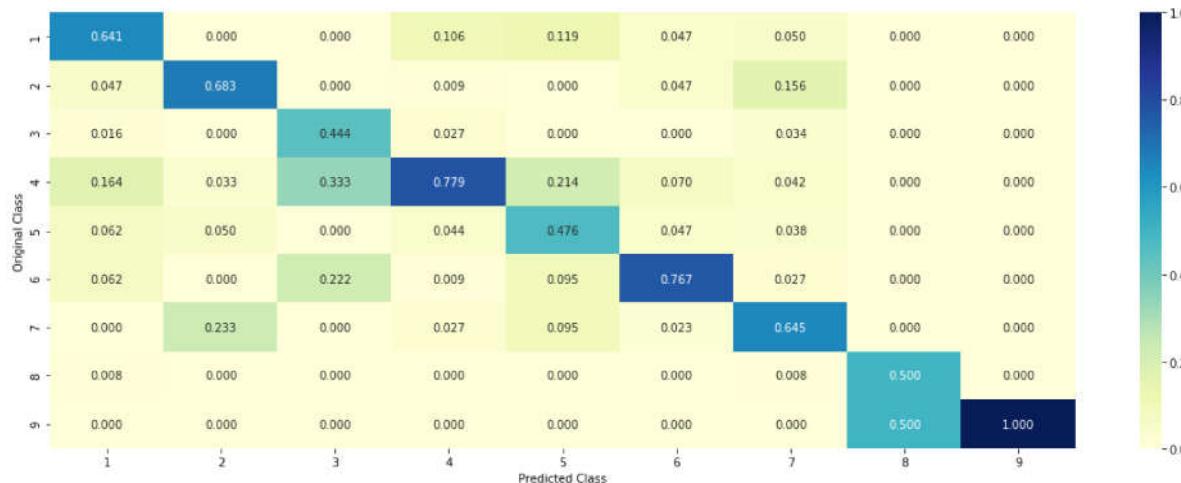
Log loss (test) on the VotingClassifier : 1.1143193824417803

Number of missclassified point : 0.3323308270676692

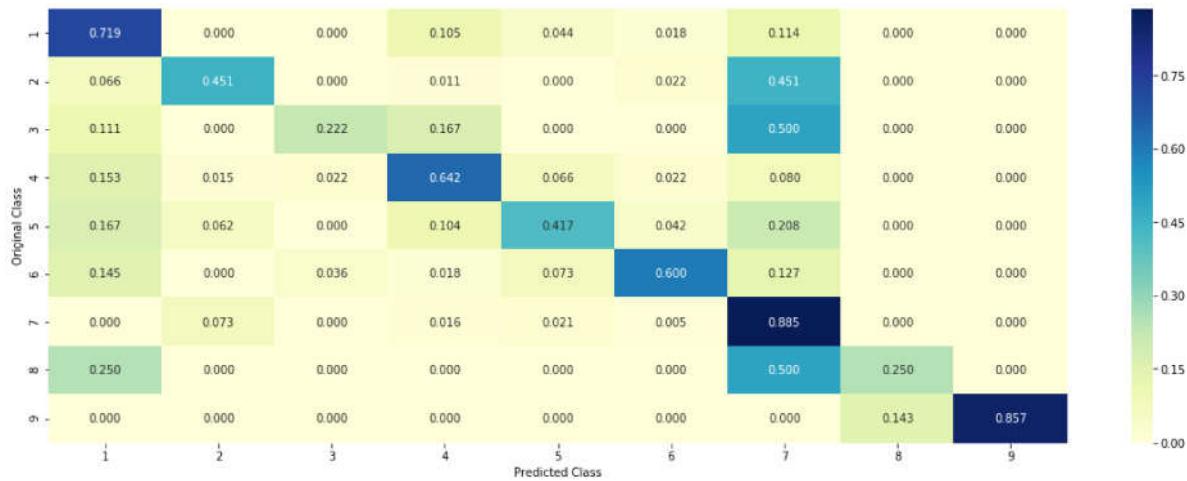
----- Confusion matrix -----



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



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



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

In [92]:

```
#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 [93]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer(ngram_range=(1, 2))
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'])

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

## Variation Feature

In [94]:

```
# 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 [95]:

```
# one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer(ngram_range=(1, 2))
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'])

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

## Text Feature

In [96]:

```
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = CountVectorizer(min_df=3, ngram_range=(1, 2))
train_text_feature_onehotCoding = text_vectorizer.fit_transform(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 words)
train_textfea_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
textfea_dict = dict(zip(list(train_text_features), train_textfea_counts))

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

Total number of unique words in train data : 776987

In [97]:

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

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

In [98]:

```
# 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)
```

## Stack above three features

In [99]:

```
# merging gene, variance and text features

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

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

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

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

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

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

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

In [100]:

```
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, 779288)
(number of data points * number of features) in test data = (665, 779288)
(number of data points * number of features) in cross validation data = (53, 779288)
```

In [101]:

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

Response encoding features :

```
(number of data points * number of features) in train data = (2124, 27)
(number of data points * number of features) in test data = (665, 27)
(number of data points * number of features) in cross validation data = (53, 27)
```

# Logistic Regression

In [102]:

```

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_
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability estimat
    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='l
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))

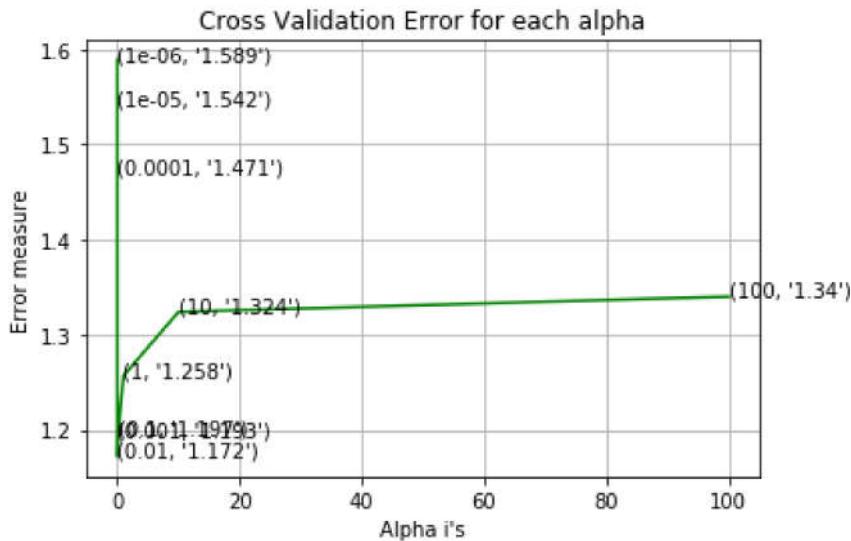
```

```

for alpha = 1e-06
Log Loss : 1.5888878065123857
for alpha = 1e-05
Log Loss : 1.5423532119912036
for alpha = 0.0001
Log Loss : 1.4709323256653155
for alpha = 0.001
Log Loss : 1.1933233102008804
for alpha = 0.01
Log Loss : 1.1719492128354452
for alpha = 0.1
Log Loss : 1.1966895193122504
for alpha = 1
Log Loss : 1.2576540959113942

```

```
for alpha = 10
Log Loss : 1.3242734239822405
for alpha = 100
Log Loss : 1.3402567752899845
```



For values of best alpha = 0.01 The train log loss is: 0.7162500668498496

For values of best alpha = 0.01 The cross validation log loss is: 1.1719492  
128354452

For values of best alpha = 0.01 The test log loss is: 1.1339179938706174

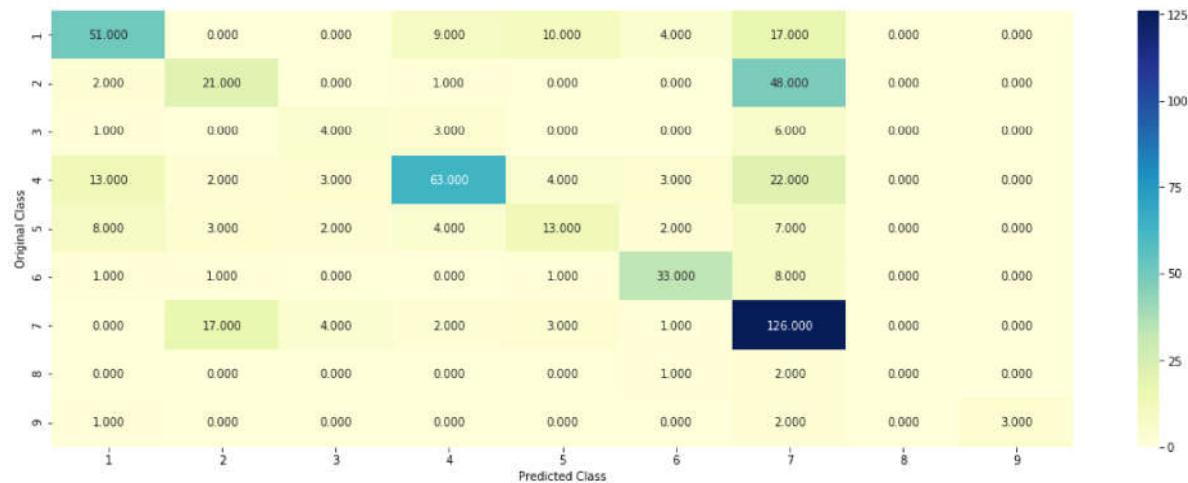
In [103]:

```
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, c)
```

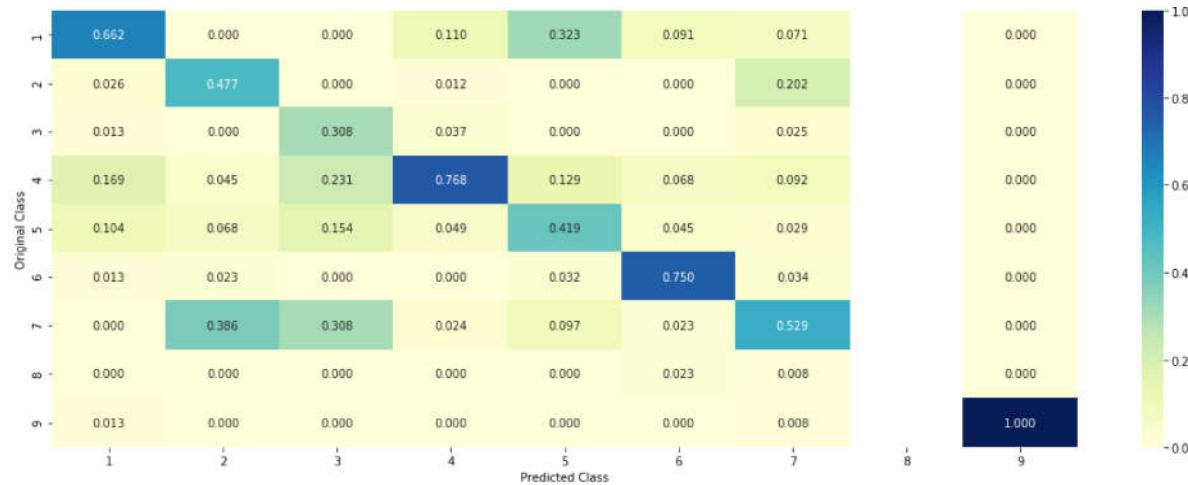
Log loss : 1.1719492128354452

Number of mis-classified points : 0.40977443609022557

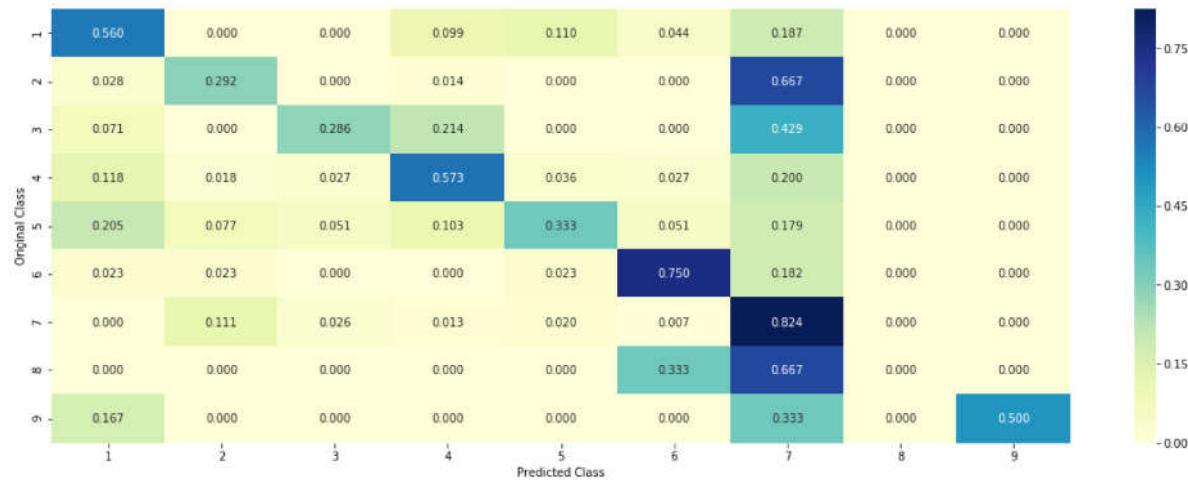
----- Confusion matrix -----



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



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



Still model does not decreases log loss values after using unigram and bigram features

# Apply Logistic regression with CountVectorizer Features, including both unigrams and Quatergrams

## Gene features

In [108]:

```
#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 [109]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer(ngram_range=(1, 4))
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'])

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

## Variation Feature

In [110]:

```
# 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 [111]:

```
# one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer(ngram_range=(1, 4))
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'])

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

## Text Feature

In [114]:

```
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = CountVectorizer(min_df=3,ngram_range=(1, 4))
train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])

# getting all the feature names (words)
# getting top 2000 feature names (words)
#train_text_features = top_mean_feats(train_text_feature_onehotCoding,
#                                      #text_vectorizer.get_feature_names(),
#                                      #top_n=2000)['feature'].tolist()
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) array
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 : 2968930

In [115]:

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

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

In [116]:

```
# 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)
```

## Stack above three features

In [118]:

```
# merging gene, variance and text features

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

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

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

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

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

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

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

In [119]:

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

One hot encoding features :  
(number of data points \* number of features) in train data = (2124, 297121  
7)  
(number of data points \* number of features) in test data = (665, 2971217)  
(number of data points \* number of features) in cross validation data = (53  
2, 2971217)

In [120]:

```
print(" Response encoding features :")
print("(number of data points * number of features) in train data = ", train_x_responseCodi
print("(number of data points * number of features) in test data = ", test_x_responseCoding
print("(number of data points * number of features) in cross validation data =", cv_x_respo
```

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 = (53  
2, 27)

## Logistic Regression

In [121]:

```

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_
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability estimat
    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='l
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 : 1.8066316539448797
for alpha = 1e-05
Log Loss : 1.6294875917407785
for alpha = 0.0001
Log Loss : 1.5972869998570707
for alpha = 0.001
Log Loss : 1.3986384646613301
for alpha = 0.01
Log Loss : 1.2497237169616435
for alpha = 0.1
Log Loss : 1.276333987473048
for alpha = 1

```

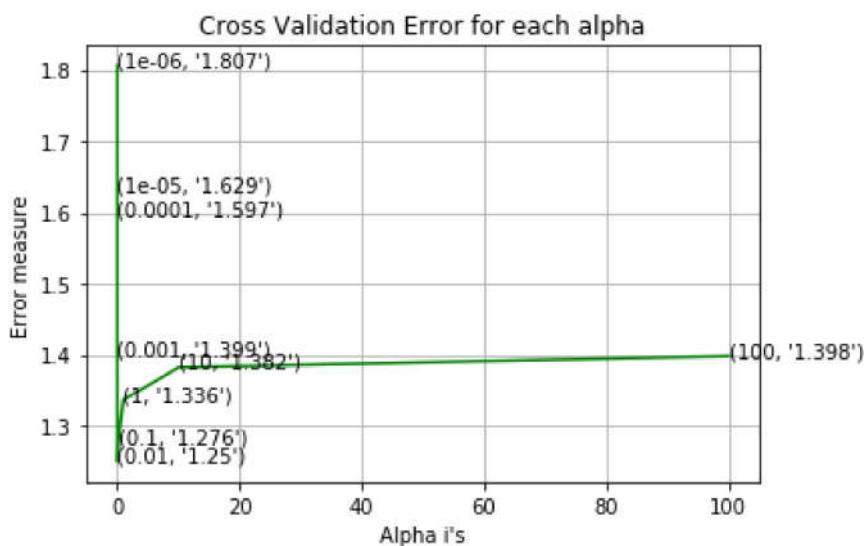
```
Log Loss : 1.3362297460847055
```

```
for alpha = 10
```

```
Log Loss : 1.3823747150482502
```

```
for alpha = 100
```

```
Log Loss : 1.3981015542348065
```



For values of best alpha = 0.01 The train log loss is: 0.7771611922961507

For values of best alpha = 0.01 The cross validation log loss is: 1.2497237  
169616435

For values of best alpha = 0.01 The test log loss is: 1.2784499657581352

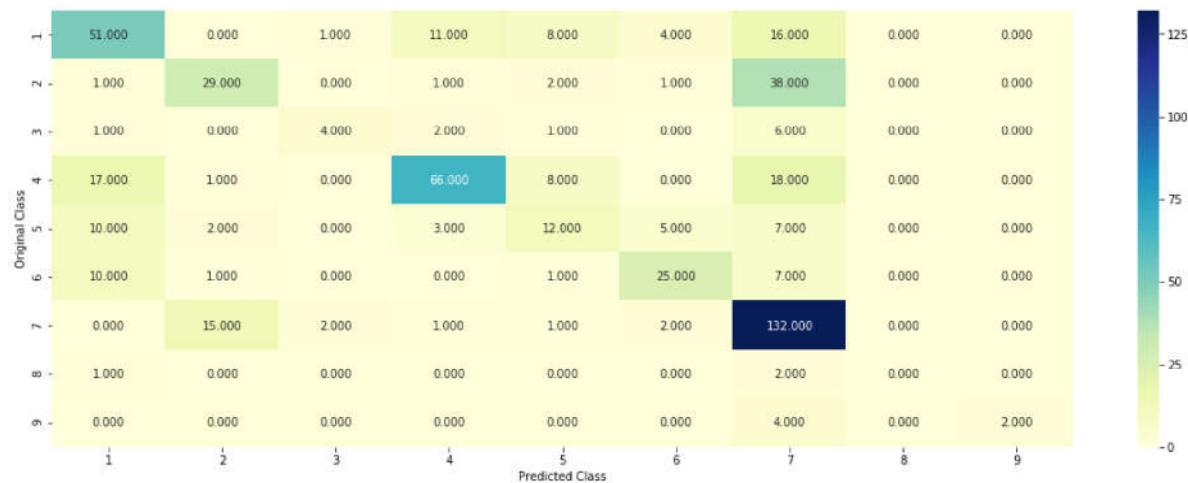
In [122]:

```
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, c)
```

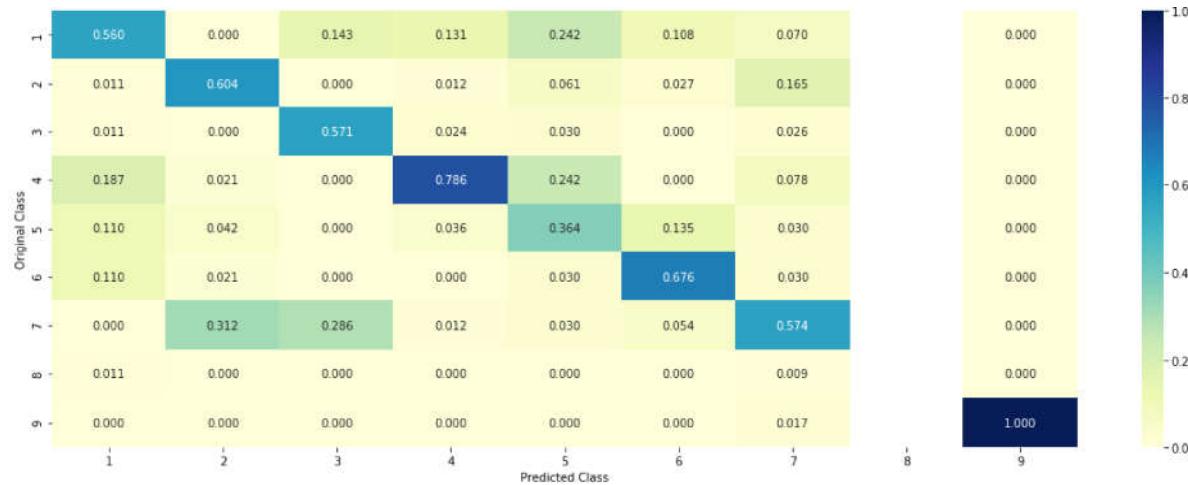
Log loss : 1.2497237169616435

Number of mis-classified points : 0.3966165413533835

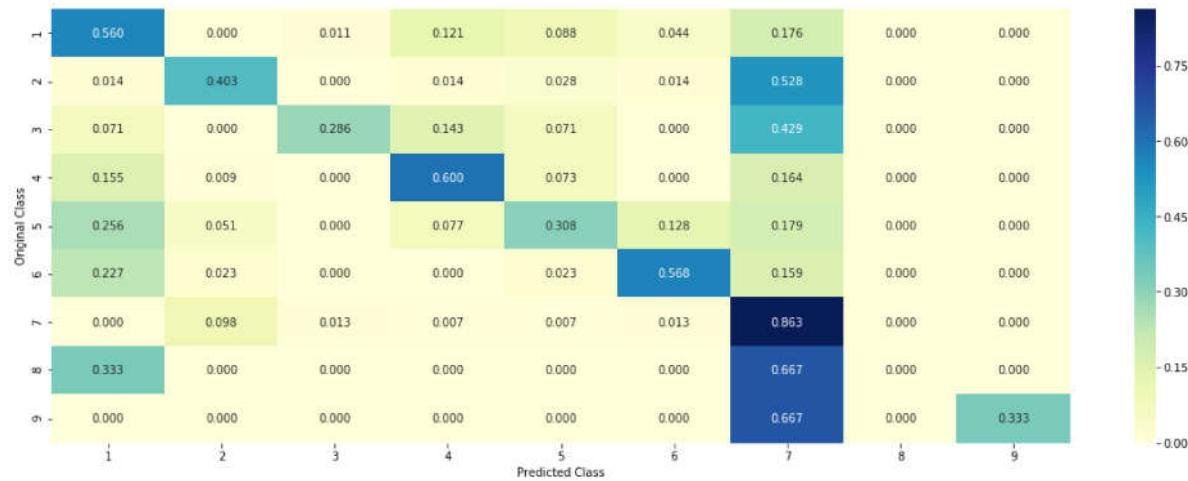
----- Confusion matrix -----



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



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



## 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**

**Lets merge gene and variation data into one list and apply TFidVectorizer on top of it.**

## Gene Feature

In [123]:

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

x_train, test_df, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_
train_df, cv_df, y_train, y_cv = train_test_split(x_train, y_train, stratify=y_train, test_
```

In [124]:

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

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

In [125]:

```
#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 [126]:

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

## Variation Feature

In [127]:

```
# 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 [128]:

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

## Text Feature

In [129]:

```

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

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

```

In [130]:

```

# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer()
train_text_feature_onehotCoding = text_vectorizer.fit_transform(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 words)
train_textfea_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_textfea_counts))

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

```

Total number of unique words in train data : 125275

In [131]:

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

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

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

In [132]:

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

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

In [133]:

```
test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
cv_text_feature_onehotCoding = text_vectorizer.transform(cv_df['TEXT'])
```

## Features after feature engineering

In [134]:

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

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

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

In [135]:

```
tfidfVectorizer = TfidfVectorizer(max_features=1000)
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'])
```

## Stack above three features

In [136]:

```
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(train_df['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(test_df['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(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 [137]:

```
print("One hot encoding features :")
print("(number of data points * number of features) in train data = ", train_x_onehotCoding.shape[1])
print("(number of data points * number of features) in test data = ", test_x_onehotCoding.shape[1])
print("(number of data points * number of features) in cross validation data = ", cv_x_onehotCoding.shape[1])
```

One hot encoding features :

```
(number of data points * number of features) in train data = (2124, 128461)
(number of data points * number of features) in test data = (665, 128461)
(number of data points * number of features) in cross validation data = (53
2, 128461)
```

In [138]:

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

```
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 = (53
2, 27)
```

In [139]:

```

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_
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability estimat
    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='l
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 : 1.105527446324402
for alpha = 1e-05
Log Loss : 1.0477104228840888
for alpha = 0.0001
Log Loss : 1.0011576044702446
for alpha = 0.001
Log Loss : 1.0573274038173124
for alpha = 0.01
Log Loss : 1.2699878622635565
for alpha = 0.1
Log Loss : 1.62245317287363
for alpha = 1

```

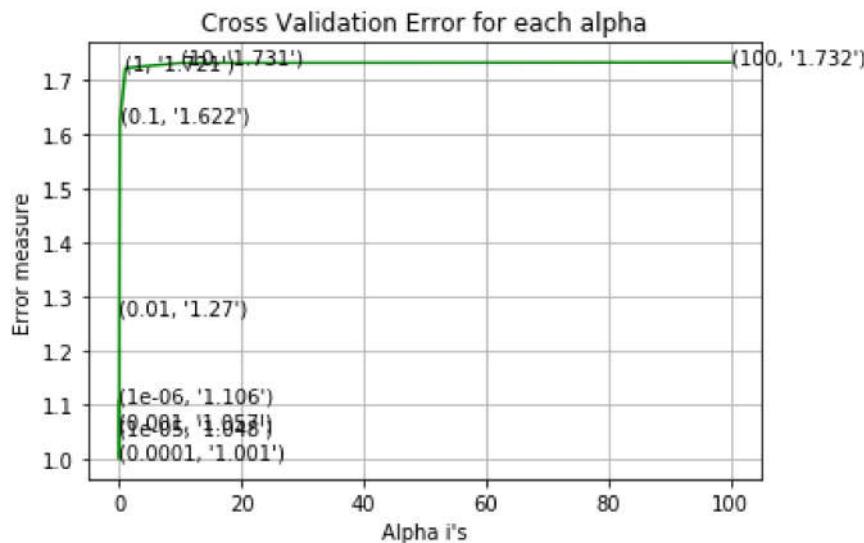
Log Loss : 1.7205552001995568

for alpha = 10

Log Loss : 1.7307834952612375

for alpha = 100

Log Loss : 1.7318319240228155



For values of best alpha = 0.0001 The train log loss is: 0.4375377376630504

3

For values of best alpha = 0.0001 The cross validation log loss is: 1.00115  
76044702446

For values of best alpha = 0.0001 The test log loss is: 1.002581386661197

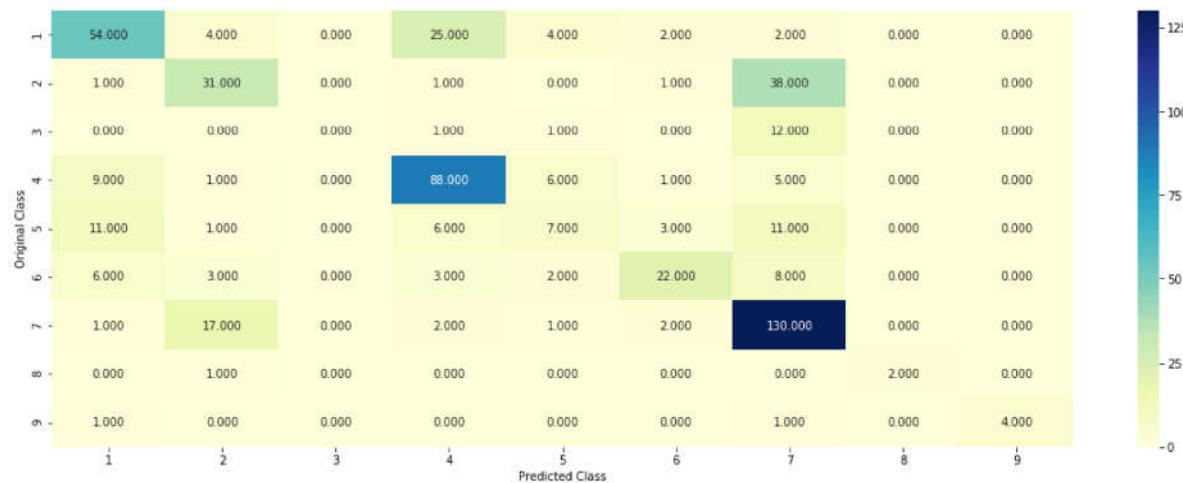
In [140]:

```
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, c)
```

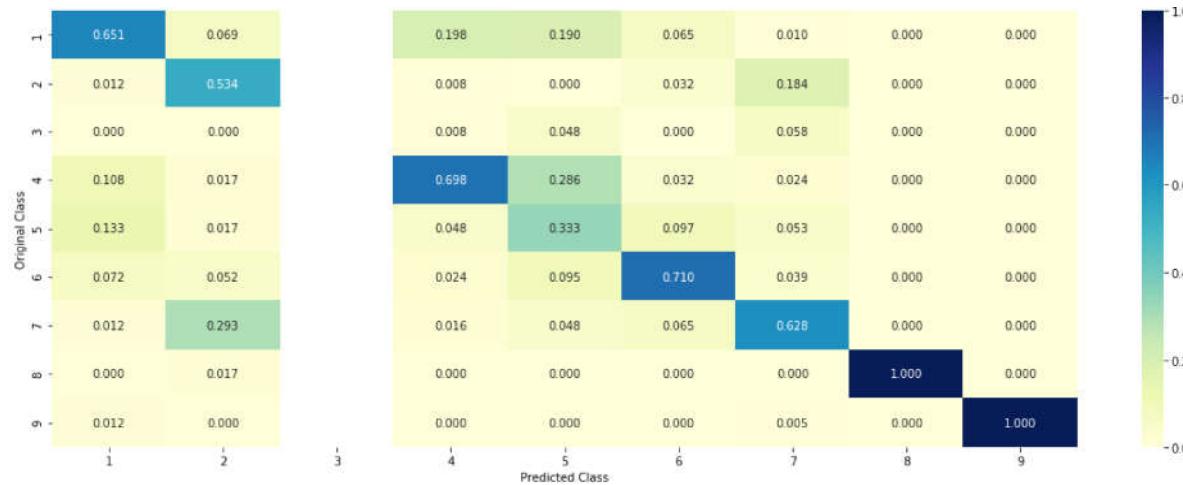
Log loss : 1.0011576044702446

Number of mis-classified points : 0.36466165413533835

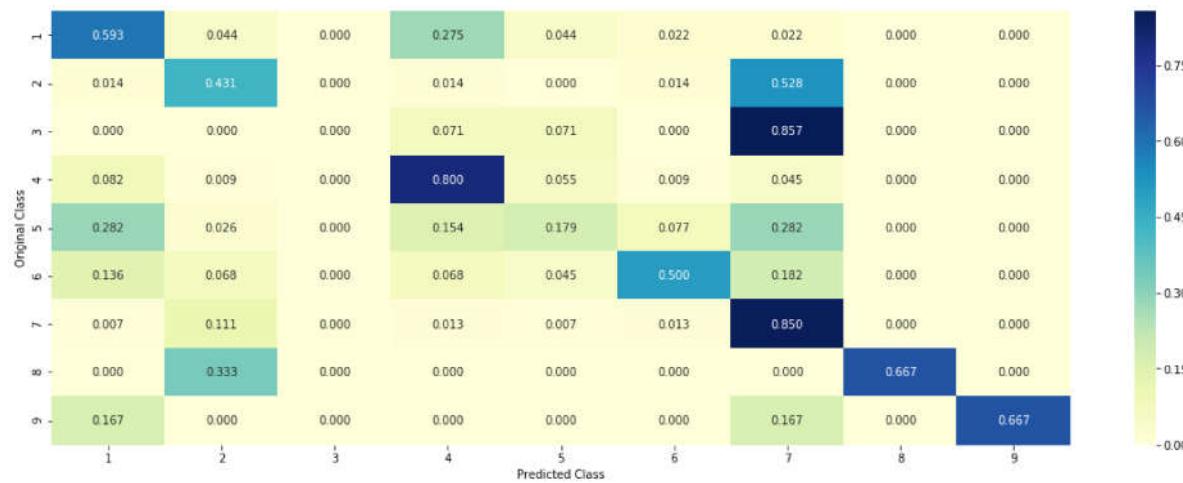
----- Confusion matrix -----



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



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



## Logistic regression + feature engineering+unigrams

# and bigrams

**merging gene and variation data into one list and apply TFidVectorizer on top of it**

## Gene Feature

In [144]:

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

x_train, test_df, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_
train_df, cv_df, y_train, y_cv = train_test_split(x_train, y_train, stratify=y_train, test_
```

In [146]:

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

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

In [147]:

```
#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 [149]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer(ngram_range=(1, 2))
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 [150]:

```
print(train_gene_feature_onehotCoding.shape)
print(test_gene_feature_onehotCoding.shape)
print(cv_gene_feature_onehotCoding.shape)
```

```
(2124, 226)
(665, 226)
(532, 226)
```

## Variation Feature

In [152]:

```
# 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 [153]:

```
# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer(ngram_range=(1, 2))
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 [154]:

```
print(train_variation_feature_onehotCoding.shape)
print(test_variation_feature_onehotCoding.shape)
print(cv_variation_feature_onehotCoding.shape)
```

```
(2124, 2048)
(665, 2048)
(532, 2048)
```

## Text Feature

In [155]:

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

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

In [156]:

```
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer(ngram_range=(1, 2))
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 words)
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 : 2353896

In [157]:

```
print(train_text_feature_onehotCoding.shape)
(2124, 2353896)
```

In [158]:

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

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

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

In [159]:

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

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

## Features after feature engineering

In [161]:

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

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

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

In [162]:

```
tfidfVectorizer = TfidfVectorizer(max_features=1000)
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 [163]:

```
print(train_text.shape)
print(test_text.shape)
print(cv_text.shape)
```

```
(2124, 1000)
(665, 1000)
(532, 1000)
```

In [164]:

```
test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
cv_text_feature_onehotCoding = text_vectorizer.transform(cv_df['TEXT'])
```

In [165]:

```
print(train_text_feature_onehotCoding.shape )
print(test_text_feature_onehotCoding.shape )
print(cv_text_feature_onehotCoding.shape )
```

```
(2124, 2353896)
(665, 2353896)
(532, 2353896)
```

## Stack above three features

In [166]:

```

train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding, train_variation_featu
test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding, test_variation_feature_
cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding, cv_variation_feature_onehot

# 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)).tocs
train_y = np.array(list(train_df['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(test_df['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(cv_df['Class']))

train_gene_var_responseCoding = np.hstack((train_gene_feature_responseCoding, train_variation
test_gene_var_responseCoding = np.hstack((test_gene_feature_responseCoding, test_variation_f
cv_gene_var_responseCoding = np.hstack((cv_gene_feature_responseCoding, cv_variation_feature

train_x_responseCoding = np.hstack((train_gene_var_responseCoding, train_text_feature_respo
test_x_responseCoding = np.hstack((test_gene_var_responseCoding, test_text_feature_response
cv_x_responseCoding = np.hstack((cv_gene_var_responseCoding, cv_text_feature_responseCoding

```

In [167]:

```

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

```

One hot encoding features :

```

(number of data points * number of features) in train data = (2124, 235717
0)
(number of data points * number of features) in test data = (665, 2357170)
(number of data points * number of features) in cross validation data = (53
2, 2357170)

```

In [168]:

```

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

```

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 = (53
2, 27)

```

In [169]:

```
cv_x_onehotCoding.shape
```

Out[169]:

```
(532, 2357170)
```

In [170]:

```
print(train_x_onehotCoding.shape)
test_x_onehotCoding.shape
```

```
(2124, 2357170)
```

Out[170]:

```
(665, 2357170)
```

## LR + FE + Response coding

In [171]:

```

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_
    clf.fit(train_x_responseCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_responseCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability estimat
    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='l
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))

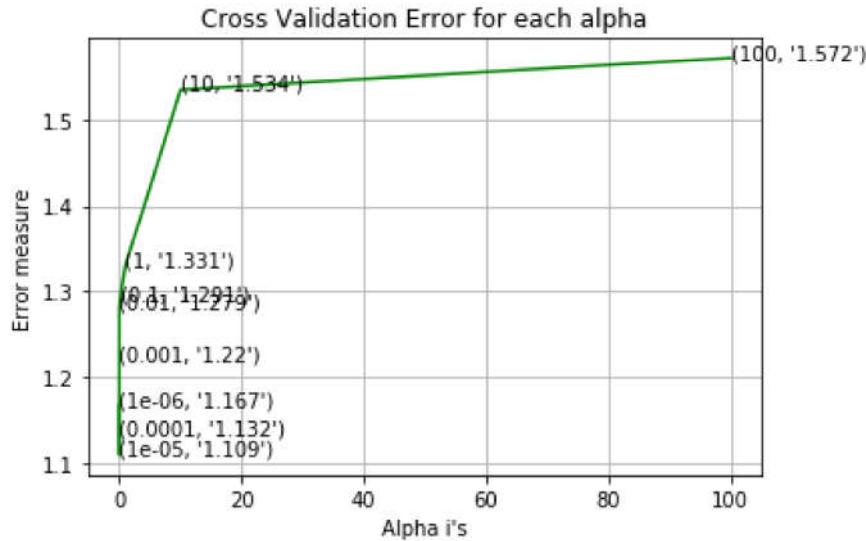
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ',
      alpha[best_alpha],
      "The cross validation log loss is:",
      log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))

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 = 1e-06
Log Loss : 1.166615377941956
for alpha = 1e-05
Log Loss : 1.108949698481365
for alpha = 0.0001
Log Loss : 1.1319978949685288
for alpha = 0.001
Log Loss : 1.2200980363200777
for alpha = 0.01
Log Loss : 1.279240571838944
for alpha = 0.1
Log Loss : 1.2908105447756197

```

```
for alpha = 1
Log Loss : 1.3305357199221168
for alpha = 10
Log Loss : 1.5344923112657416
for alpha = 100
Log Loss : 1.571581129875967
```



For values of best alpha = 1e-05 The train log loss is: 0.6159240484249302

For values of best alpha = 1e-05 The cross validation log loss is: 1.108949  
698481365

For values of best alpha = 1e-05 The test log loss is: 1.1037071899186286

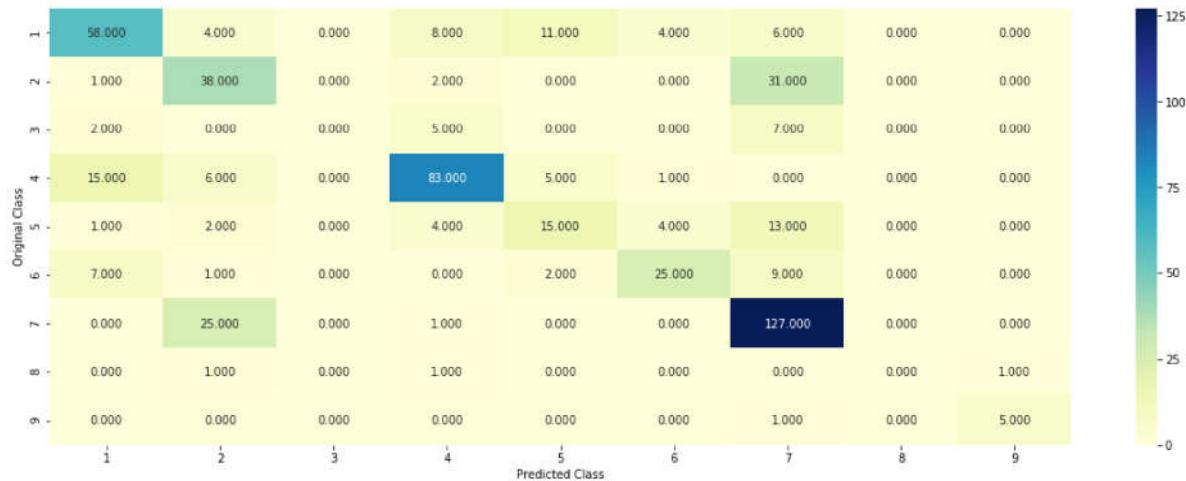
In [172]:

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

Log loss : 1.108949698481365

Number of mis-classified points : 0.34022556390977443

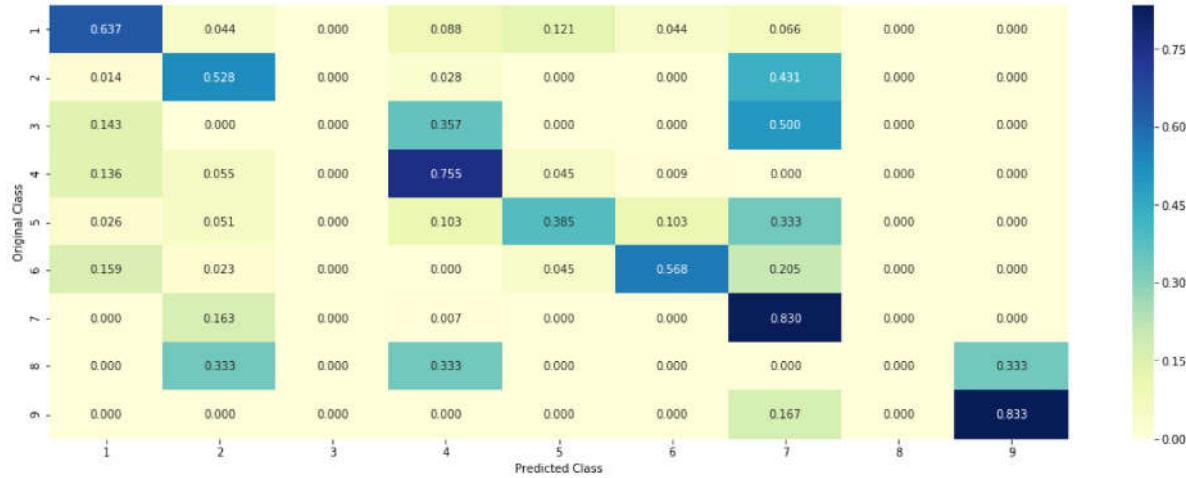
----- Confusion matrix -----



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



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

**LR + FE + TfIdf**

In [173]:

```

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_
    clf.fit(train_x_onehotCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability estimat
    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='l
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 : 1.1075287259590694
for alpha = 1e-05
Log Loss : 1.0350873487835524
for alpha = 0.0001
Log Loss : 0.9969212115686542
for alpha = 0.001
Log Loss : 1.053306564221487
for alpha = 0.01
Log Loss : 1.2536572910934247
for alpha = 0.1
Log Loss : 1.5936442310422554
for alpha = 1

```

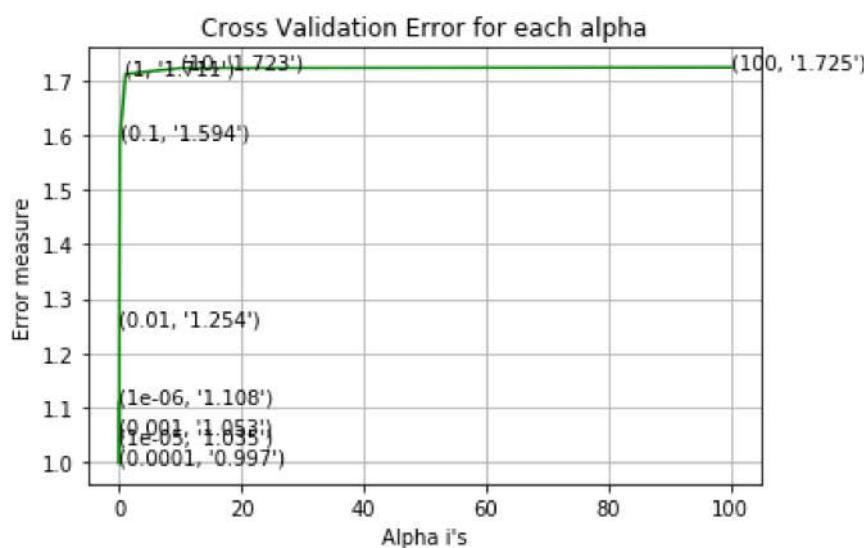
```
Log Loss : 1.7113381715545082
```

```
for alpha = 10
```

```
Log Loss : 1.723327353998306
```

```
for alpha = 100
```

```
Log Loss : 1.7246019205915866
```



```
For values of best alpha = 0.0001 The train log loss is: 0.4066870708103957
```

```
3
```

```
For values of best alpha = 0.0001 The cross validation log loss is: 0.99692  
12115686542
```

```
For values of best alpha = 0.0001 The test log loss is: 1.0556962301711788
```

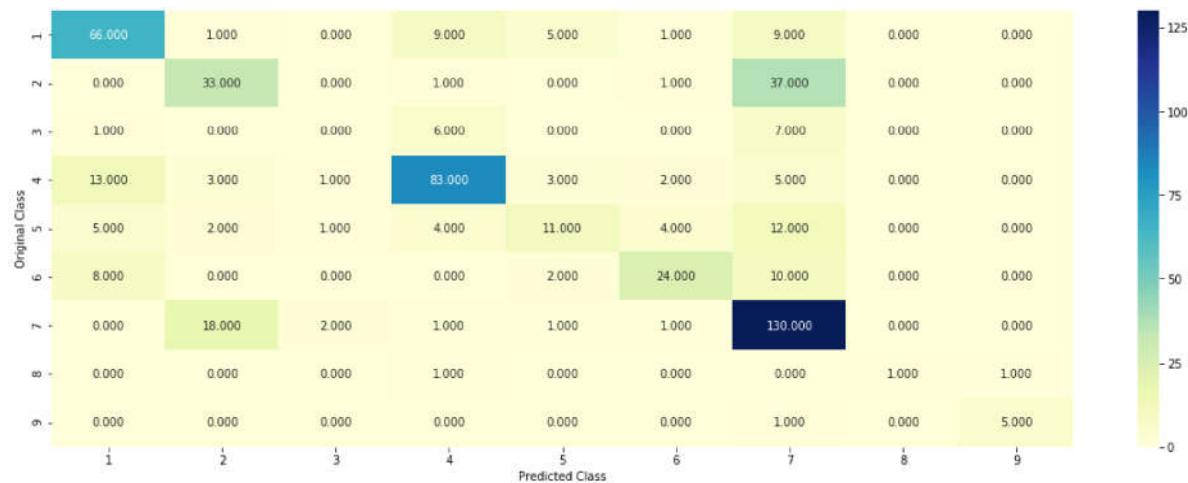
In [174]:

```
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, c)
```

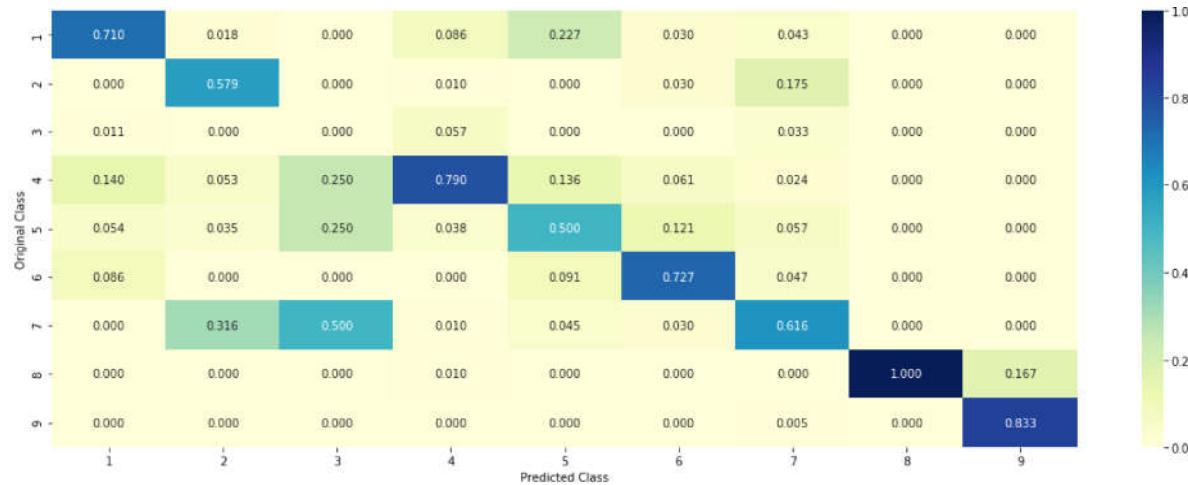
Log loss : 0.9969212115686542

Number of mis-classified points : 0.33646616541353386

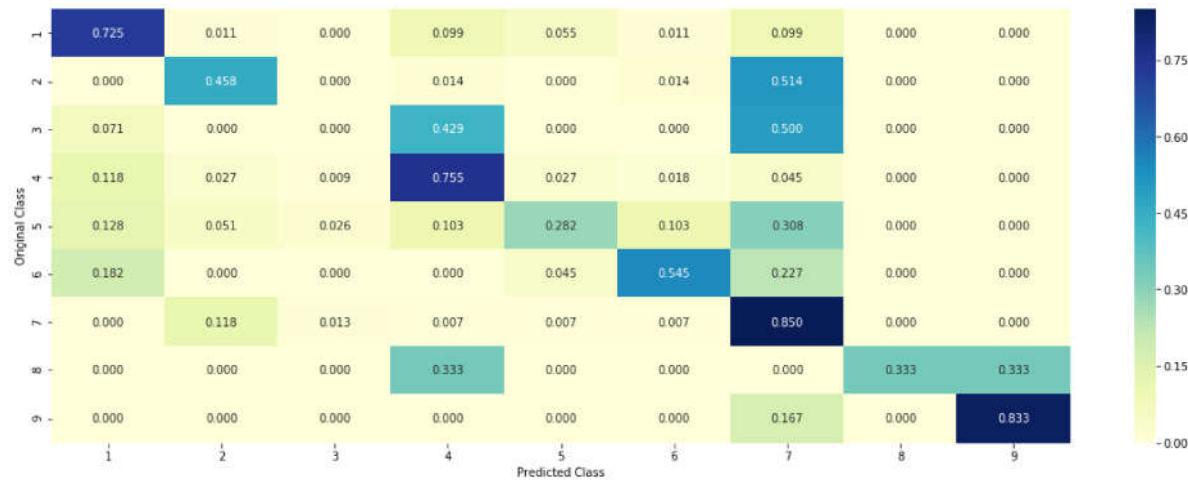
----- Confusion matrix -----



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



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



## Conclusion

# without Feature Engineering

In [176]:

```
print()
from prettytable import PrettyTable
ptable = PrettyTable()
ptable.title = "*** Model Summary *** [Performance Metric: Log-Loss]"
ptable.field_names=["Model Name","Train","CV","Test","% Misclassified Points"]
ptable.add_row(["Naive Bayes","0.92","1.23","1.19","42"])
ptable.add_row(["KNN","0.47","1.02","1.03","35"])
ptable.add_row(["Logistic Regression With Class balancing","0.52","1.06","1.00","34"])
ptable.add_row(["Logistic Regression Without Class balancing","0.51","1.09","1.01","33"])
ptable.add_row(["Linear SVM","0.56","1.12","1.03","36"])
ptable.add_row(["Random Forest Classifier With One hot Encoding","0.65","1.14","1.15","38"])
ptable.add_row(["Random Forest Classifier With Response Coding","0.064","1.30","1.31","49"])
ptable.add_row(["Stack Models:LR+NB+SVM","1.06","1.46","1.25","33"])
ptable.add_row(["Maximum Voting classifier","0.82","1.12","1.11","33"])
print(ptable)
print()
```

|    | Model Name                                     | Train | CV   | Test | % M<br>isclassified Points |
|----|--|-------|------|------|----------------------------|
| 42 | Naive Bayes                                    | 0.92  | 1.23 | 1.19 |                            |
| 35 | KNN  | 0.47  | 1.02 | 1.03 |                            |
| 34 | Logistic Regression With Class balancing       | 0.52  | 1.06 | 1.00 |                            |
| 33 | Logistic Regression Without Class balancing    | 0.51  | 1.09 | 1.01 |                            |
| 36 | Linear SVM                                     | 0.56  | 1.12 | 1.03 |                            |
| 38 | Random Forest Classifier With One hot Encoding | 0.65  | 1.14 | 1.15 |                            |
| 49 | Random Forest Classifier With Response Coding  | 0.064 | 1.30 | 1.31 |                            |
| 33 | Stack Models:LR+NB+SVM                         | 1.06  | 1.46 | 1.25 |                            |
| 33 | Maximum Voting classifier                      | 0.82  | 1.12 | 1.11 |                            |

# With Feature Engineering

In [177]:

```
from prettytable import PrettyTable
ptable1 = PrettyTable()
ptable1.field_names=["Model Name","Train Log-Loss","CV Log-Loss","Test Log-Loss","% Misclas"]
ptable1.add_row(["LR with Response Coding(unigrams and bigrams)","0.71","1.17","1.13","40"])
ptable1.add_row(["LR with Response Coding(unigrams and quatergram)","0.77","1.24","1.27","39"])
ptable1.add_row(["LR with one hot encoding + feature engineering","0.43","1.00","1.0","36"])
ptable1.add_row(["LR with one hot encoding+unigram-quatergram+feature engineering","0.61","1.1","1.1","36"])
ptable1.add_row(["LR with Response coding+unigram-quatergram+feature engineering","0.40","0.99","1.0","33"])
print(ptable1)
```

| Model Name  | Train Log   |               |                        |  |
|---|-------------|---------------|------------------------|--|
| -Loss   | CV Log-Loss | Test Log-Loss | % Misclassified Points |  |
| LR with Response Coding(unigrams and bigrams)                   | 0.71        |               |                        |  |
| 1.17   1.13   1.13   40   |             |               |                        |  |
| LR with Response Coding(unigrams and quatergram)                | 0.77        |               |                        |  |
| 1.24   1.27   1.27   39   |             |               |                        |  |
| LR with one hot encoding + feature engineering                  | 0.43        |               |                        |  |
| 1.00   1.0   1.0   36   |             |               |                        |  |
| LR with one hot encoding+unigram-quatergram+feature engineering | 0.61        |               |                        |  |
| 1.1   1.1   1.1   36  |             |               |                        |  |
| LR with Response coding+unigram-quatergram+feature engineering  | 0.40        |               |                        |  |
| 0.99   1.0   1.0   33   |             |               |                        |  |

- We have applied some feature engineering like merging gene and variation and length of text as a new feature but it's difficult to decrease the test logloss to less than 1.
- In our logistic Regression(without class balancing) model the train log loss is great but the difference between train log loss and test/cv log loss is high here so, model may overfits.
- In terms of train log loss our logistic regression(without class balancing) model is good but overall logistic regression(with class balancing) is out best model because test/cv logloss are comparatively low as compared to others and the difference between train and test/cv logloss is also low.