

Personalized Cancer Diagnosis

1. Business Problem

1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

Context:

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

Problem statement :

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

1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

1. <https://www.forbes.com/sites/matthewherper/2017/06/03/a-new-cancer-drug-helped-almost-everyone-who-took-it-almost-heres-what-it-teaches-us/#2a44ee2f6b25> (<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 probabilities => Metric is Log-loss.
- No Latency constraints.

2.3. Train, CV and Test Datasets

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

In [2]:

```
import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import numpy as np
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import SGDClassifier
from 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
from tqdm import tqdm
```

Reading Data

Reading Gene and Variation Data

In [3]:

```
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[3]:

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

Reading Text Data

In [4]:

```
# 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[4]:

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...

Preprocessing of text

In [5]:

```
# 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 [6]:

```
#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 : 58.23205780818062 seconds

In [7]:

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

Out[7]:

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

In [8]:

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

Out[8]:

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 [9]:

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

In [10]:

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

Out[10]:

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

Test, Train and Cross Validation Split

In [11]:

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

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

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

In [12]:

```
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
```

Prediction using a 'Random' Model

In [13]:

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

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

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

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

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

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

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

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

In [31]:

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

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

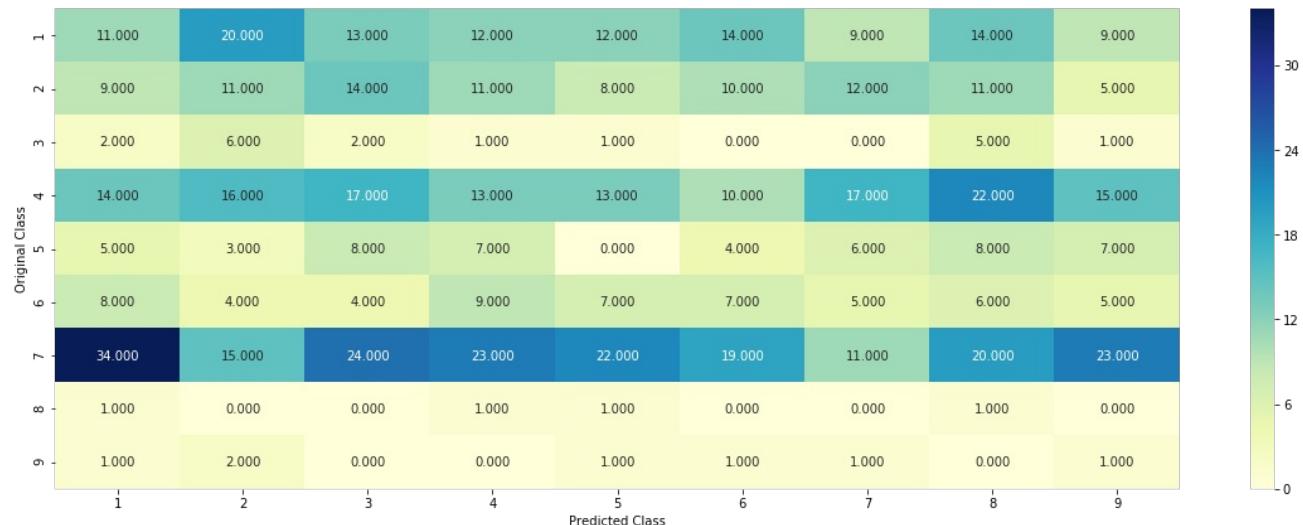
# Test-Set error.
#we create a output array that has exactly same as the test data
test_predicted_y = np.zeros((test_data_len,9))
for i in range(test_data_len):
    rand_probs = np.random.rand(1,9)
    test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random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

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

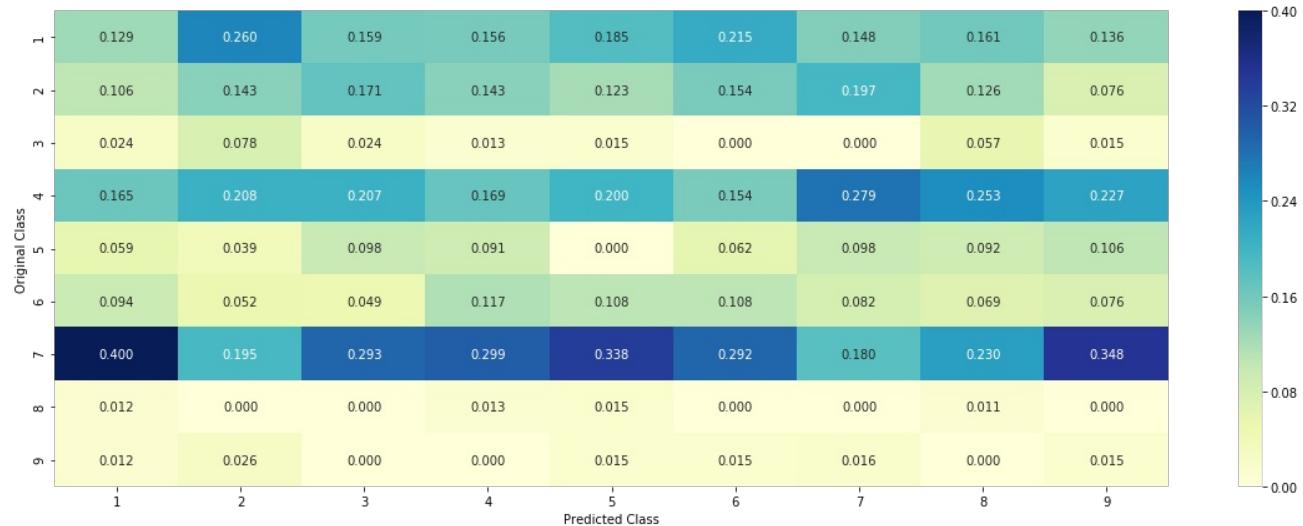
Log loss on Cross Validation Data using Random Model 2.506267624986588

Log loss on Test Data using Random Model 2.510451429078548

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



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



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



1.Applying all the models with tfdf feature

Univariate Analysis

In [16]:

```
# code for response coding with Laplace smoothing.
# alpha : used for laplace smoothing
# feature: ['gene', 'variation']
```

```

# df: ['train_df', 'test_df', 'cv_df']

# algorithm
# -----
# Consider all unique values and the number of occurrences of given feature in train data dataframe
# build a vector (1*9) , the first element = (number of times it occurred in class1 + 10*alpha / number of time it
occurred in total data+90*alpha)
# gv_dict is like a look up table, for every gene it store a (1*9) representation of it
# for a value of feature in df:
# if it is in train data:
# we add the vector that was stored in 'gv_dict' look up table to 'gv_fea'
# if it is not there is train:
# we add [1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9] to 'gv_fea'
# return 'gv_fea'
# ----

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

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

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

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    # { 'BRCA1': [0.200757575757575, 0.03787878787878788, 0.0681818181818177, 0.13636363636363635, 0.25, 0
.193181818181818, 0.03787878787878788, 0.03787878787878788, 0.03787878787878788],
    #   'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.27040816326530615, 0.06122
489795918366, 0.066326530612244902, 0.051020408163265307, 0.051020408163265307, 0.056122448979591837],
    #   ...
    # }

```

```

#      'EGFR': [0.0568181818181816, 0.21590909090909091, 0.0625, 0.0681818181818177, 0.0681818181818177
, 0.0625, 0.34659090909090912, 0.0625, 0.0568181818181816],
#      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.0787878787878782, 0.139
3939393939394, 0.34545454545454546, 0.060606060606060608, 0.060606060606060608, 0.060606060606060608],
#      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.46540880503144655, 0.0754
71698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081761006289, 0.062893081761006289],
#      'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295, 0.06622
5165562913912, 0.066225165562913912, 0.27152317880794702, 0.066225165562913912, 0.066225165562913912],
#      'BRAF': [0.066666666666666666, 0.17999999999999999, 0.07333333333333334, 0.07333333333333334, 0.0933
3333333333338, 0.08000000000000002, 0.29999999999999999, 0.066666666666666666, 0.066666666666666666],
#      ...
#      ]
gv_dict = get_gvfea_dict(alpha, feature, df)
# value_count is similar in get_gvfea_dict
value_count = train_df[feature].value_counts()

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

```

Gene feature categories

In [14]:

```

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 : 237
BRCA1      164
TP53       105
EGFR        84
BRCA2       82
PTEN        77
BRAF        60
KIT          59
PDGFRA     46
ERBB2       46
ALK          42
Name: Gene, dtype: int64

```

In [33]:

```

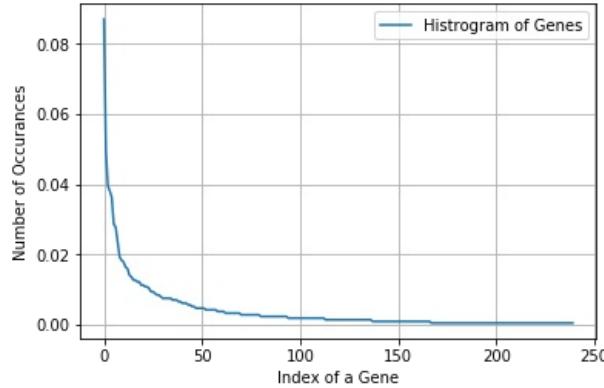
print("Ans: There are", unique_genes.shape[0] , "different categories of genes in the train data, and they are dis
tributed as follows",)

```

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

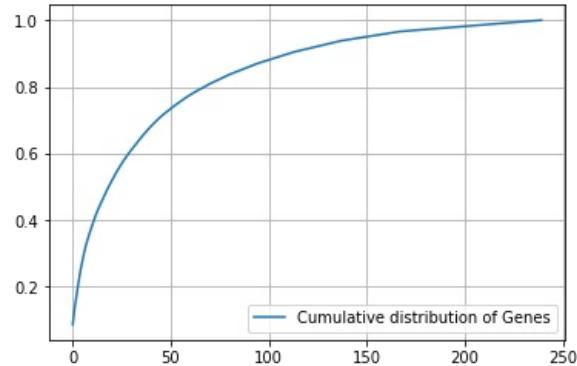
In [35]:

```
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 [36]:

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



response coding of gene feature

In [17]:

```
#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 [18]:

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

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

One hot encoding with gene feature

In [39]:

```
# 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 [40]:

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

Out[40]:

```
681      CDKN2A
2017     MAP2K1
1606      VHL
2213     PTEN
2761     BRAF
Name: Gene, dtype: object
```

In [41]:

```
gene_vectorizer.get_feature_names()
```

Out[41]:

```
['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ari',
 'araf',
 'arid1a',
 'asxl1',
 'atm',
 'atr',
 'atrx',
 'aurka',
 'aurkb',
 'axin1',
 'b2m',
 'bap1',
 'bccl10',
 'bccl2',
 'bccl2l11',
 'bcor',
 'braf',
 'brcal',
 'brc2a',
 'brd4',
 'brip1',
 'btik',
 'card11',
 'carm1',
 'casp8',
 'cb1',
 'ccnd1',
 'ccnd2',
 'ccnd3',
 'ccne1',
 'cdh1',
 'cdk12',
 'cdk4',
 'cdk6',
 'cdk8',
 'cdkn1a',
 'cdkn1b',
 'cdkn2a',
 'cdkn2b',
 'cdkn2c',
 'chek2',
 'cic',
 'crebbp',
 'ctcf',
 'ctla4',
 'ctnnb1',
 'ddr2',
 'dicer1',
 'dnmt3a',
```

'dnmt3b',
'egfr',
'eif1ax',
'elf3',
'ep300',
'epas1',
'epcam',
'erb2bb',
'erb3bb',
'erb4bb',
'ercc2',
'ercc3',
'ercc4',
'erg',
'errfi1',
'esr1',
'etv6',
'ewsrl1',
'ezh2',
'fam58a',
'fanca',
'fancc',
'fat1',
'fbxw7',
'fgf19',
'fgf3',
'fgf4',
'fgfr1',
'fgfr2',
'fgfr3',
'fgfr4',
'flt1',
'flt3',
'foxal1',
'foxl2',
'foxo1',
'foxp1',
'fubp1',
'gata3',
'gli1',
'gnas',
'h3f3a',
'hist1h1c',
'hnf1a',
'hras',
'idh1',
'idh2',
'igfir',
'ikbbk',
'il7r',
'jak1',
'jak2',
'kdm5a',
'kdm5c',
'kdm6a',
'kdr',
'keap1',
'kit',
'klf4',
'kmt2a',
'kmt2c',
'kmt2d',
'knstrn',
'kras',
'lats1',
'lats2',
'map2k1',
'map2k2',
'map2k4',
'map3kl',
'mapk1',
'mdm4',
'med12',
'mef2b',
'men1',
'met',
'mga',
'mlh1',
'mpl',
'msh2',
'msh6',
'mtor',
'myc',

'mycn',
'myd88',
'myod1',
'ncor1',
'nf1',
'nf2',
'nfe2l2',
'nkbia',
'nkx2',
'notch1',
'notch2',
'npm1',
'nras',
'nsd1',
'ntrk1',
'ntrk2',
'ntrk3',
'pak1',
'pbprm1',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pik3r3',
'pim1',
'pms2',
'pole',
'ppp2r1a',
'ppp6c',
'prdm1',
'ptch1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rab35',
'rac1',
'rad21',
'rad50',
'rad51b',
'rad51c',
'rad51d',
'rad541',
'rafi',
'rara',
'rasal',
'rb1',
'rbm10',
'ret',
'rheb',
'rhoa',
'rictor',
'rit1',
'rnf43',
'ros1',
'runx1',
'rxra',
'rybp',
'sdhb',
'sdhc',
'setd2',
'sf3b1',
'shoc2',
'smad2',
'smad3',
'smad4',
'smarca4',
'smarcb1',
'smo',
'sos1',
'sox9',
'spop',
'src',
'srpf2',
'stag2',
'stat3',
'stk11',
'tcf3',
'tcf7l2',
'tert',

```
'tet1',  
'tet2',  
'tgfb1',  
'tgfb2',  
'tmpRSS2',  
'tp53',  
'tp53bp1',  
'tsc1',  
'tsc2',  
'u2af1',  
'vegfa',  
'vh1',  
'whsc1',  
'whsc1l1',  
'xpol',  
'xrcc2',  
'yap1']
```

In [42]:

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

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

Applying LR model on Gene feature

In [43]:

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

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

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

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


cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_gene_feature_onehotCoding, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_gene_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-1
5))

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

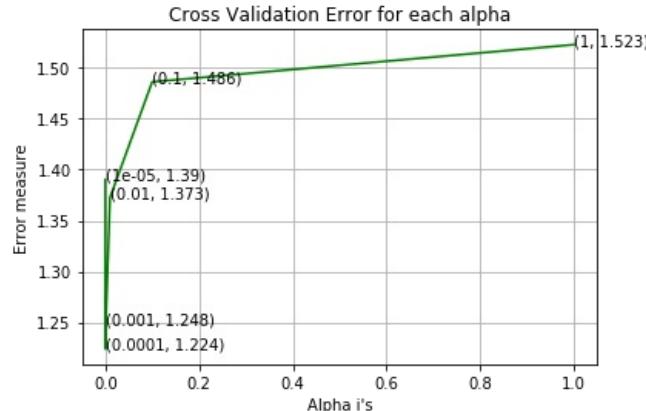
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_gene_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_gene_feature_onehotCoding, y_train)

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

```

For values of alpha = 1e-05 The log loss is: 1.390240860545501
For values of alpha = 0.0001 The log loss is: 1.2236828539239384
For values of alpha = 0.001 The log loss is: 1.2483652490081318
For values of alpha = 0.01 The log loss is: 1.372569880430752
For values of alpha = 0.1 The log loss is: 1.4864186423031736
For values of alpha = 1 The log loss is: 1.5230512871565993

```



```

For values of best alpha = 0.0001 The train log loss is: 1.04459591155956
For values of best alpha = 0.0001 The cross validation log loss is: 1.2236828539239384
For values of best alpha = 0.0001 The test log loss is: 1.1888593393387916

```

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

yes it is stable

In [27]:

```

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

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

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

```

Q6. How many data points in Test and CV datasets are covered by the 235 genes in train dataset?

Ans

1. In test data 649 out of 665 : 97.59398496240601
2. In cross validation data 516 out of 532 : 96.99248120300751

Univariate Analysis on Variation Feature

How many categories are there

In [19]:

```

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 : 1924
Truncating_Mutations      55
Deletion                  53
Amplification              53
Fusions                   18
T58I                      3
Overexpression              3
T286A                     2
Q61L                      2
TMPRSS2-ETV1_Fusion        2
EWSR1-ETV1_Fusion          2
Name: Variation, dtype: int64

```

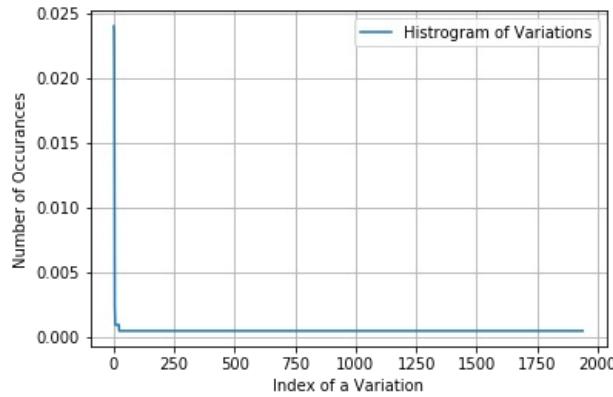
In [16]:

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

Ans: There are 1926 different categories of variations in the train data, and they are distributed as follows

In [46]:

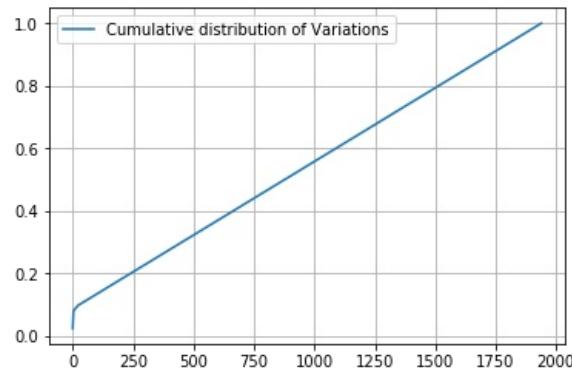
```
s = sum(unique_variations.values);
h = unique_variations.values/s;
plt.plot(h, label="Histogram of Variations")
plt.xlabel('Index of a Variation')
plt.ylabel('Number of Occurrences')
plt.legend()
plt.grid()
plt.show()
```



In [47]:

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

```
[0.0240113  0.04755179  0.06826742 ... 0.99905838  0.99952919  1.] ]
```



response coding of variation feature

In [20]:

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

In [21]:

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

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

Tfidf One hot encoding On variation feature

In [37]:

```
# 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 [38]:

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

train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The s hape of Variation feature: (2124, 1963)

Applying LR on Variation feature

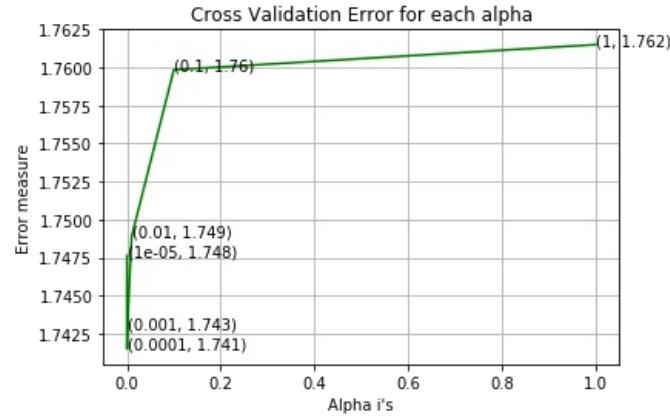
In [36]:

```
alpha = [10 ** x for x in range(-5, 1)]  
  
# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html  
# -----  
# default parameters  
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,  
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,  
# class_weight=None, warm_start=False, average=False, n_iter=None)  
  
# some of methods  
# fit(X, y[, coef_init, intercept_init, ...])      Fit linear model with Stochastic Gradient Descent.  
# predict(X)    Predict class labels for samples in X.  
  
#-----  
# video link:  
#-----  
  
cv_log_error_array=[]  
for i in alpha:  
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)  
    clf.fit(train_variation_feature_onehotCoding, y_train)  
  
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")  
    sig_clf.fit(train_variation_feature_onehotCoding, y_train)  
    predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)  
  
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))  
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))  
  
fig, ax = plt.subplots()  
ax.plot(alpha, cv_log_error_array,c='g')  
for i, txt in enumerate(np.round(cv_log_error_array,3)):  
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))  
plt.grid()  
plt.title("Cross Validation Error for each alpha")  
plt.xlabel("Alpha i's")  
plt.ylabel("Error measure")  
plt.show()  
  
best_alpha = np.argmin(cv_log_error_array)  
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)  
clf.fit(train_variation_feature_onehotCoding, y_train)  
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")  
sig_clf.fit(train_variation_feature_onehotCoding, y_train)  
  
predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)  
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))  
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)  
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))  
predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding)  
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```

For values of alpha = 1e-05 The log loss is: 1.7476281847848731
For values of alpha = 0.0001 The log loss is: 1.7414887600625417
For values of alpha = 0.001 The log loss is: 1.7428512778438146
For values of alpha = 0.01 The log loss is: 1.7488820642158998
For values of alpha = 0.1 The log loss is: 1.759847833865149
For values of alpha = 1 The log loss is: 1.761504344175262

```



```

For values of best alpha = 0.0001 The train log loss is: 0.751664368028053
For values of best alpha = 0.0001 The cross validation log loss is: 1.7414887600625417
For values of best alpha = 0.0001 The test log loss is: 1.7062888067314097

```

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

No it is not stable

In [37]:

```

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

```

Q12. How many data points are covered by total 1911 genes in test and cross validation data sets?
Ans

1. In test data 63 out of 665 : 9.473684210526317
2. In cross validation data 47 out of 532 : 8.834586466165414

Univariate Analysis on Text Feature

How many unique words are present in train data?

In [38]:

```

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

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

```

In [39]:

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

Tfidf One hot encoding On text feature

In [40]:

```
# building a CountVectorizer with all the words that occurred minimum 3 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 features) vector
train_text_fea_counts = train_text_feature_onehotCoding.sum(axis=0).A1

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

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

Total number of unique words in train data : 53206

In [41]:

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

# dict_list[i] is build on i'th class text data
# total_dict is 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 [42]:

```
#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 [43]:

```
# 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 [44]:

```
# 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 [45]:

```
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()))
```

Applying LR on Text feature

In [47]:

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

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

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

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


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

    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_text_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-1
5))

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

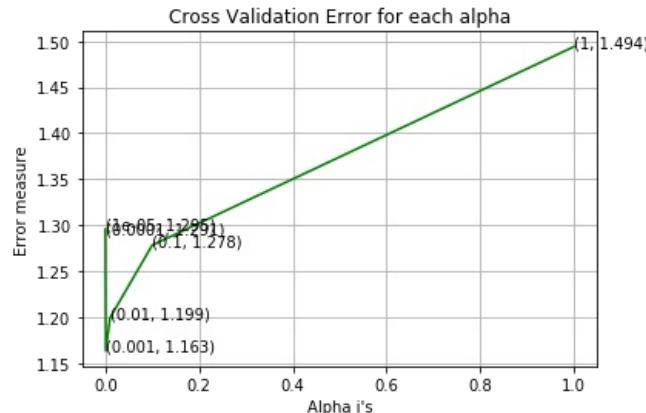
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_text_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_text_feature_onehotCoding, y_train)

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

```

For values of alpha = 1e-05 The log loss is: 1.2954800884598383
For values of alpha = 0.0001 The log loss is: 1.290754484121766
For values of alpha = 0.001 The log loss is: 1.1630255521069561
For values of alpha = 0.01 The log loss is: 1.1990275247747033
For values of alpha = 0.1 The log loss is: 1.2780559091099124
For values of alpha = 1 The log loss is: 1.494174023531577

```



```

For values of best alpha = 0.001 The train log loss is: 0.6598558801595399
For values of best alpha = 0.001 The cross validation log loss is: 1.1630255521069561
For values of best alpha = 0.001 The test log loss is: 1.1141572593109064

```

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

Ans. it is not stable

In [48]:

```

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

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

```

In [49]:

```

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

```

95.651 % of word of test data appeared in train data
98.769 % of word of Cross Validation appeared in train data

Machine Learning Models

In [50]:

```

#Data preparation for ML models.

#Misc. functionns for ML models

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

    # for calculating log_loss we will provide the array of probabilities belongs to each class
    print("Log loss :",log_loss(test_y, sig_clf.predict_proba(test_x)))
    # calculating the number of data points that are misclassified
    print("Number of mis-classified points :", np.count_nonzero((pred_y- test_y))/test_y.shape[0])
    plot_confusion_matrix(test_y, pred_y)

```

In [51]:

```
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 [52]:

```
# 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_imptfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = CountVectorizer()
    var_count_vec = CountVectorizer()
    text_count_vec = CountVectorizer(min_df=3)

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

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

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

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

Stacking all Three features

In [53]:

```
# merging gene, variance and text features

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

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

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

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

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

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

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

In [54]:

```
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, 55380)
(number of data points * number of features) in test data = (665, 55380)
(number of data points * number of features) in cross validation data = (532, 55380)
```

In [55]:

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

Response encoding features :

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

Base Line Model

Applying Naive Bayes on personalized data

Hyper parameter tuning

In [56]:

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

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


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


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

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

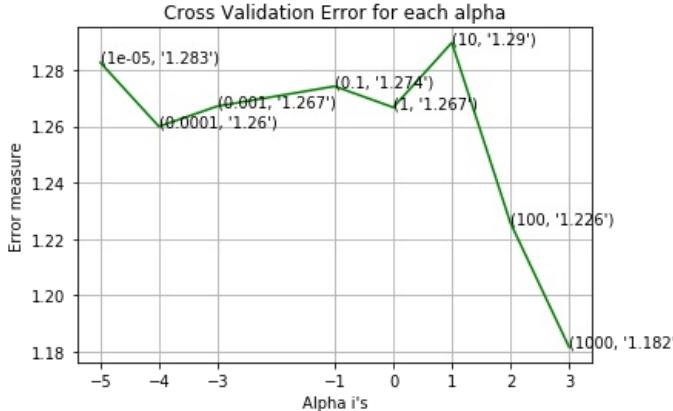
best_alpha = np.argmin(cv_log_error_array)
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_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-05
Log Loss : 1.282646492074082
for alpha = 0.0001
Log Loss : 1.2599386018193945
for alpha = 0.001
Log Loss : 1.2671508180447075
for alpha = 0.1
Log Loss : 1.2742154451226286
for alpha = 1
Log Loss : 1.2667298444310577
for alpha = 10
Log Loss : 1.2897043078779504
for alpha = 100
Log Loss : 1.2257175833775553
for alpha = 1000
Log Loss : 1.1816823252796311

```



```

For values of best alpha = 1000 The train log loss is: 0.920660992331587
For values of best alpha = 1000 The cross validation log loss is: 1.1816823252796311
For values of best alpha = 1000 The test log loss is: 1.1818599109742163

```

Testing the model with best hyper parameters

In [57]:

```

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

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


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

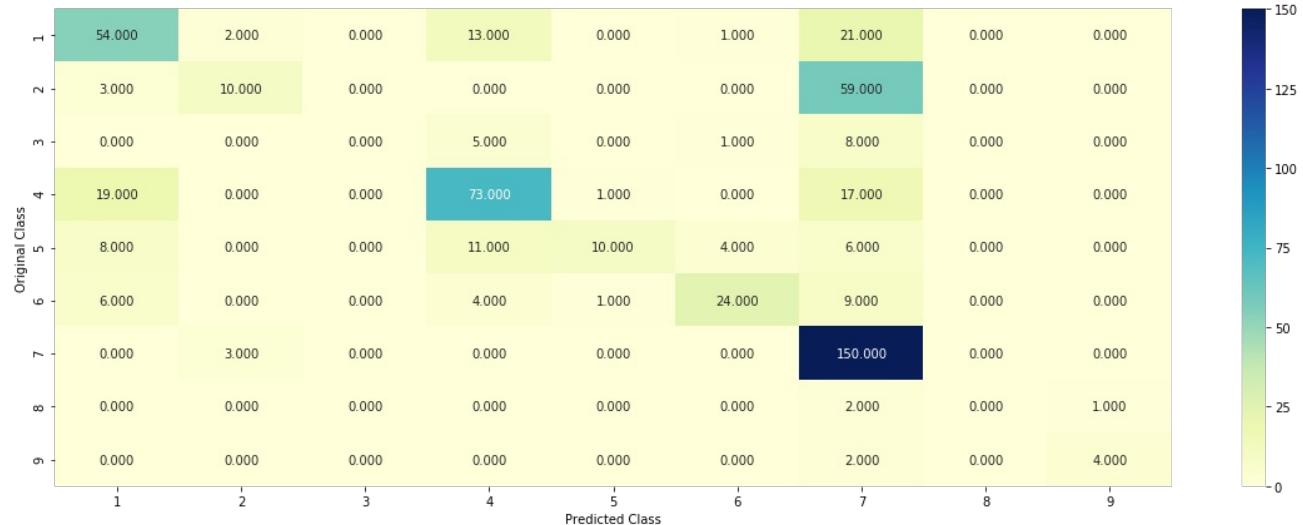

clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
# to avoid rounding error while multiplying probabilites we use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

```

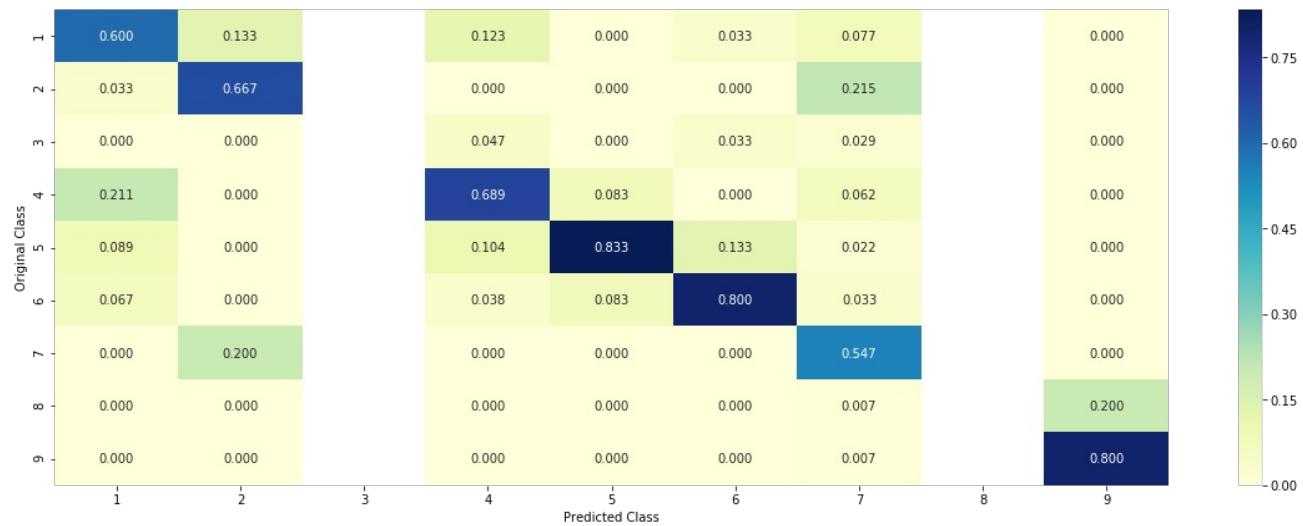
Log Loss : 1.1816823252796311

Number of missclassified point : 0.3890977443609023

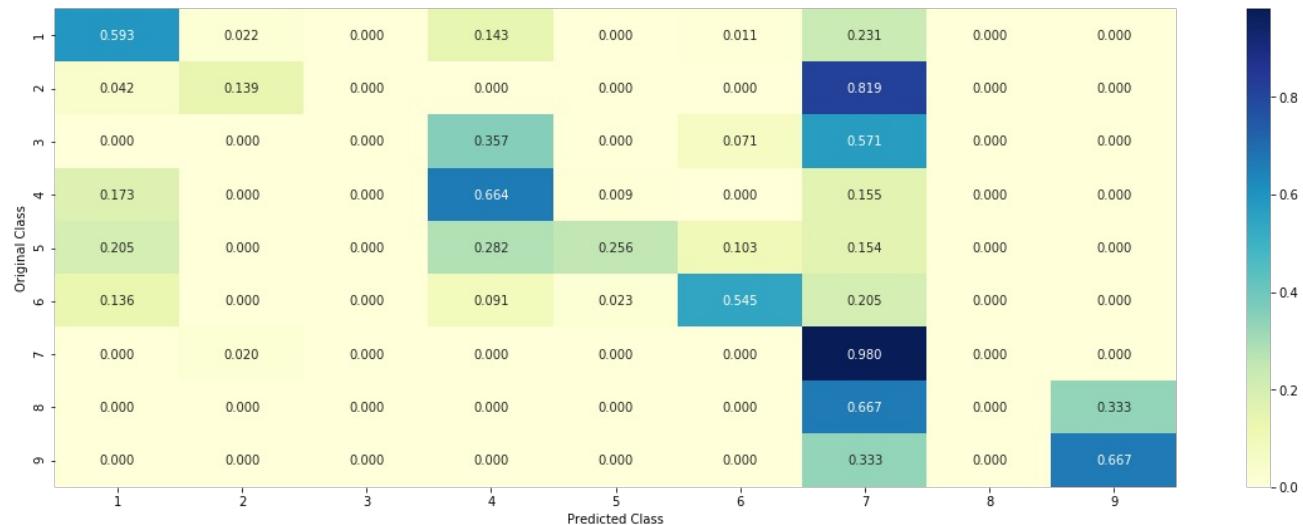
----- Confusion matrix -----



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



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



Feature Importance, Correctly classified point

In [58]:

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

Predicted Class : 7
Predicted Class Probabilities: [[0.0134 0.0299 0.0509 0.0222 0.0632 0.0161 0.8023 0.0018 0.]]
Actual Class : 3

16 Text feature [cells] present in test data point [True]
17 Text feature [kinase] present in test data point [True]
18 Text feature [contrast] present in test data point [True]
19 Text feature [activated] present in test data point [True]
20 Text feature [downstream] present in test data point [True]
21 Text feature [presence] present in test data point [True]
22 Text feature [phosphorylation] present in test data point [True]
23 Text feature [activation] present in test data point [True]
24 Text feature [cell] present in test data point [True]
27 Text feature [shown] present in test data point [True]
28 Text feature [growth] present in test data point [True]
29 Text feature [compared] present in test data point [True]
30 Text feature [inhibitor] present in test data point [True]
31 Text feature [10] present in test data point [True]
34 Text feature [signaling] present in test data point [True]
35 Text feature [also] present in test data point [True]
36 Text feature [however] present in test data point [True]
37 Text feature [addition] present in test data point [True]
38 Text feature [independent] present in test data point [True]
39 Text feature [recently] present in test data point [True]
40 Text feature [suggest] present in test data point [True]
41 Text feature [found] present in test data point [True]
42 Text feature [treated] present in test data point [True]
43 Text feature [similar] present in test data point [True]
44 Text feature [higher] present in test data point [True]
45 Text feature [showed] present in test data point [True]
46 Text feature [mutations] present in test data point [True]
47 Text feature [3b] present in test data point [True]
48 Text feature [increased] present in test data point [True]
49 Text feature [previously] present in test data point [True]
50 Text feature [well] present in test data point [True]
51 Text feature [tyrosine] present in test data point [True]
53 Text feature [potential] present in test data point [True]
55 Text feature [respectively] present in test data point [True]
56 Text feature [constitutive] present in test data point [True]
57 Text feature [consistent] present in test data point [True]
58 Text feature [described] present in test data point [True]
59 Text feature [mutant] present in test data point [True]
60 Text feature [constitutively] present in test data point [True]
61 Text feature [using] present in test data point [True]
62 Text feature [treatment] present in test data point [True]
63 Text feature [mutation] present in test data point [True]
64 Text feature [figure] present in test data point [True]
65 Text feature [absence] present in test data point [True]
66 Text feature [fig] present in test data point [True]
67 Text feature [sensitive] present in test data point [True]
68 Text feature [may] present in test data point [True]
69 Text feature [mechanism] present in test data point [True]
72 Text feature [reported] present in test data point [True]
74 Text feature [including] present in test data point [True]
75 Text feature [without] present in test data point [True]
77 Text feature [inhibitors] present in test data point [True]
78 Text feature [inhibition] present in test data point [True]
80 Text feature [activating] present in test data point [True]
81 Text feature [approximately] present in test data point [True]
82 Text feature [demonstrated] present in test data point [True]
84 Text feature [pathways] present in test data point [True]
85 Text feature [3a] present in test data point [True]
86 Text feature [followed] present in test data point [True]
87 Text feature [observed] present in test data point [True]
88 Text feature [expression] present in test data point [True]
89 Text feature [confirmed] present in test data point [True]
90 Text feature [activate] present in test data point [True]
93 Text feature [increase] present in test data point [True]
94 Text feature [expressed] present in test data point [True]
97 Text feature [4a] present in test data point [True]

Out of the top 100 features 66 are present in query point

Feature Importance, Incorrectly classified point

In [59]:

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

Predicted Class : 7
Predicted Class Probabilities: [[0.0047 0.0383 0.0036 0.0056 0.0171 0.012 0.9165 0.0021 0.]]
Actual Class : 7

16 Text feature [cells] present in test data point [True]
17 Text feature [kinase] present in test data point [True]
18 Text feature [contrast] present in test data point [True]
19 Text feature [activated] present in test data point [True]
20 Text feature [downstream] present in test data point [True]
21 Text feature [presence] present in test data point [True]
22 Text feature [phosphorylation] present in test data point [True]
23 Text feature [activation] present in test data point [True]
24 Text feature [cell] present in test data point [True]
25 Text feature [expressing] present in test data point [True]
26 Text feature [factor] present in test data point [True]
27 Text feature [shown] present in test data point [True]
28 Text feature [growth] present in test data point [True]
29 Text feature [compared] present in test data point [True]
30 Text feature [inhibitor] present in test data point [True]
31 Text feature [10] present in test data point [True]
34 Text feature [signaling] present in test data point [True]
35 Text feature [also] present in test data point [True]
36 Text feature [however] present in test data point [True]
37 Text feature [addition] present in test data point [True]
38 Text feature [independent] present in test data point [True]
39 Text feature [recently] present in test data point [True]
40 Text feature [suggest] present in test data point [True]
41 Text feature [found] present in test data point [True]
42 Text feature [treated] present in test data point [True]
43 Text feature [similar] present in test data point [True]
44 Text feature [higher] present in test data point [True]
45 Text feature [showed] present in test data point [True]
46 Text feature [mutations] present in test data point [True]
47 Text feature [3b] present in test data point [True]
48 Text feature [increased] present in test data point [True]
49 Text feature [previously] present in test data point [True]
50 Text feature [well] present in test data point [True]
51 Text feature [tyrosine] present in test data point [True]
52 Text feature [1a] present in test data point [True]
53 Text feature [potential] present in test data point [True]
55 Text feature [respectively] present in test data point [True]
56 Text feature [constitutive] present in test data point [True]
57 Text feature [consistent] present in test data point [True]
58 Text feature [described] present in test data point [True]
59 Text feature [mutant] present in test data point [True]
60 Text feature [constitutively] present in test data point [True]
61 Text feature [using] present in test data point [True]
62 Text feature [treatment] present in test data point [True]
63 Text feature [mutation] present in test data point [True]
64 Text feature [figure] present in test data point [True]
65 Text feature [absence] present in test data point [True]
66 Text feature [fig] present in test data point [True]
67 Text feature [sensitive] present in test data point [True]
68 Text feature [may] present in test data point [True]
69 Text feature [mechanism] present in test data point [True]
70 Text feature [obtained] present in test data point [True]
71 Text feature [inhibited] present in test data point [True]
72 Text feature [reported] present in test data point [True]
73 Text feature [serum] present in test data point [True]
74 Text feature [including] present in test data point [True]
75 Text feature [without] present in test data point [True]
76 Text feature [enhanced] present in test data point [True]
77 Text feature [inhibitors] present in test data point [True]
78 Text feature [inhibition] present in test data point [True]
79 Text feature [furthermore] present in test data point [True]
80 Text feature [activating] present in test data point [True]
81 Text feature [approximately] present in test data point [True]
82 Text feature [demonstrated] present in test data point [True]
83 Text feature [interestingly] present in test data point [True]
84 Text feature [pathways] present in test data point [True]

```
85 Text feature [3a] present in test data point [True]
86 Text feature [followed] present in test data point [True]
87 Text feature [observed] present in test data point [True]
88 Text feature [expression] present in test data point [True]
89 Text feature [confirmed] present in test data point [True]
90 Text feature [activate] present in test data point [True]
91 Text feature [total] present in test data point [True]
92 Text feature [various] present in test data point [True]
93 Text feature [increase] present in test data point [True]
94 Text feature [expressed] present in test data point [True]
95 Text feature [antibodies] present in test data point [True]
96 Text feature [performed] present in test data point [True]
97 Text feature [4a] present in test data point [True]
98 Text feature [proliferation] present in test data point [True]
99 Text feature [induced] present in test data point [True]
Out of the top 100 features 81 are present in query point
```

K Nearest Neighbour Classification

Hyper parameter tuning

In [60]:

```
# find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

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


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


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

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

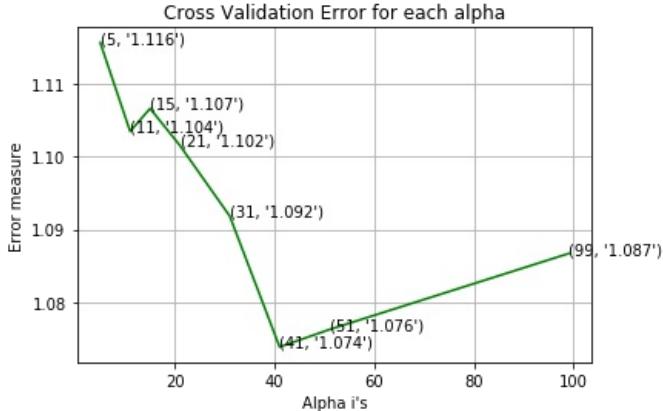
best_alpha = np.argmin(cv_log_error_array)
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

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

```

for alpha = 5
Log Loss : 1.1157424693118543
for alpha = 11
Log Loss : 1.1035060508890109
for alpha = 15
Log Loss : 1.106665692527968
for alpha = 21
Log Loss : 1.101558248643624
for alpha = 31
Log Loss : 1.091825344450353
for alpha = 41
Log Loss : 1.0737973379149046
for alpha = 51
Log Loss : 1.0761995009736451
for alpha = 99
Log Loss : 1.0866731567753485

```



```

For values of best alpha = 41 The train log loss is: 0.8749737005203456
For values of best alpha = 41 The cross validation log loss is: 1.0737973379149046
For values of best alpha = 41 The test log loss is: 1.1127189007296343

```

Testing the model with best hyper parameters

In [61]:

```

# find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
# -----
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
# metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

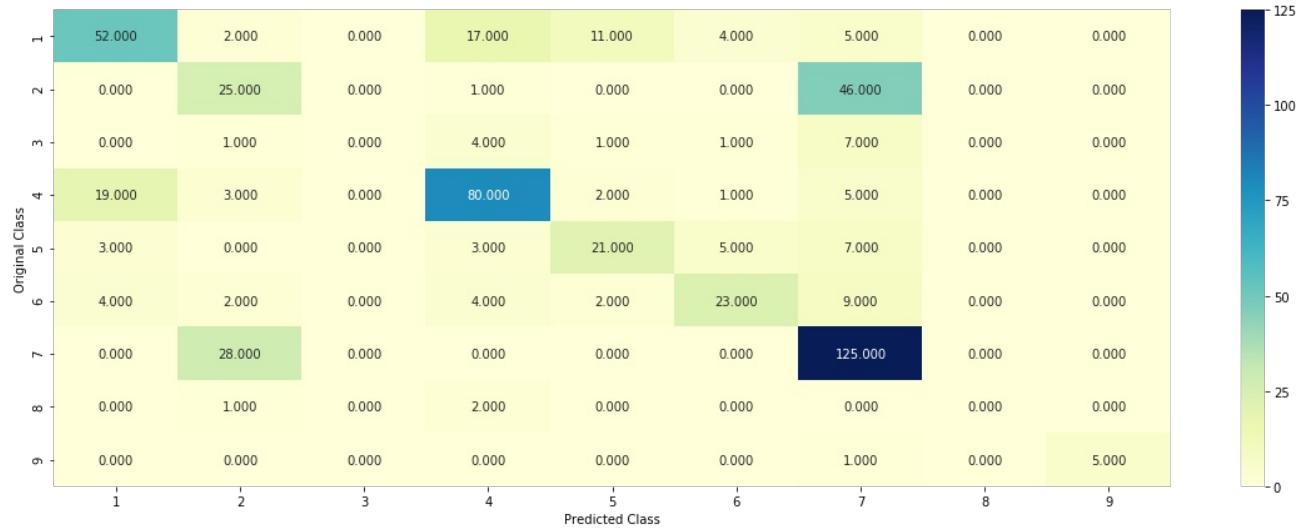
# methods of
# fit(X, y) : Fit the model using X as training data and y as target values
# predict(X):Predict the class labels for the provided data
# predict_proba(X):Return probability estimates for the test data X.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-nearest-neighbors-geometric-intuition-with-a-toy-e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

```

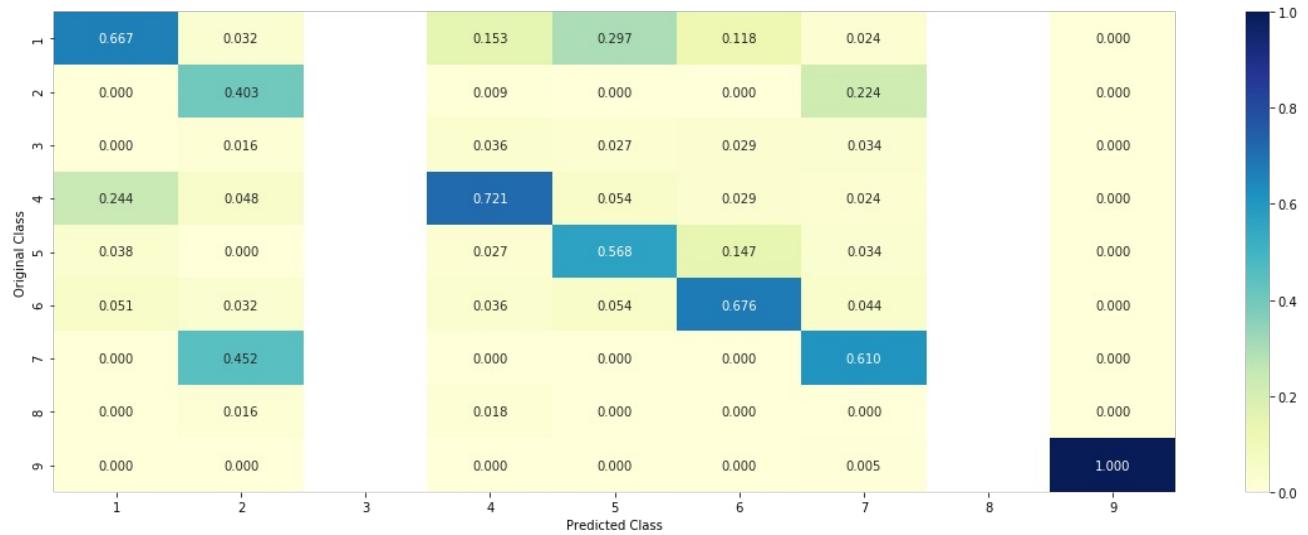
```

Log loss : 1.0737973379149046
Number of mis-classified points : 0.37781954887218044
----- Confusion matrix -----

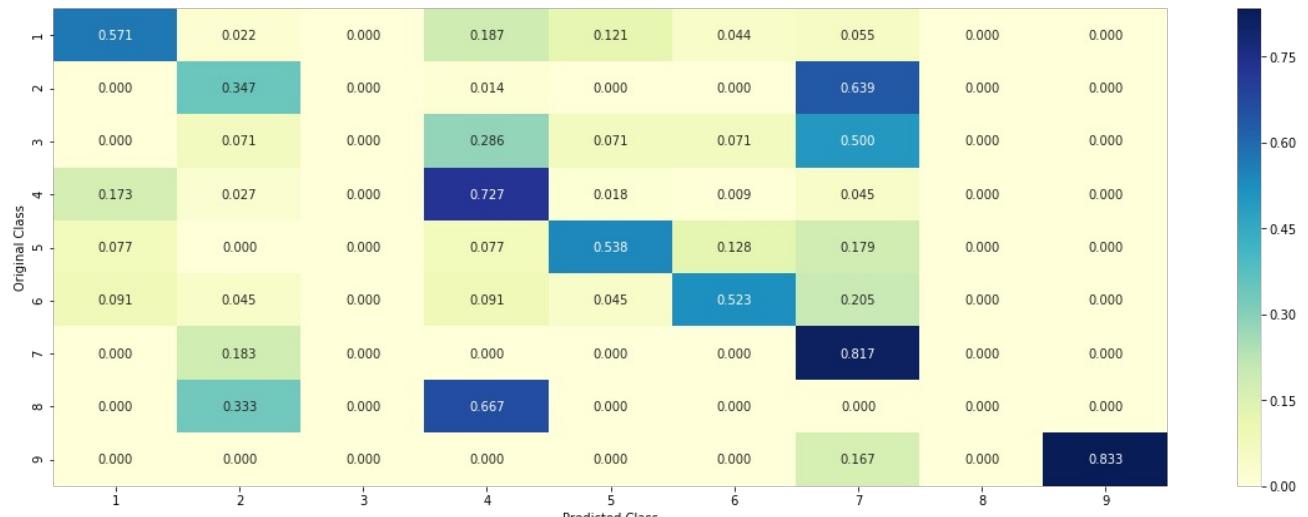
```



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Sample Query point -1

In [62]:

```

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

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

```

Predicted Class : 6

Actual Class : 3

The 41 nearest neighbours of the test points belongs to classes [3 3 7 5 5 3 7 3 3 5 3 3 3 7 5 7 7
7 7 7 7 7 7 2 7 7 2 7 7 2 7 7 7 7 7 5 7 7
7 5 5 7]

Frequency of nearest points : Counter({7: 23, 3: 8, 5: 7, 2: 3})

Sample Query Point-2

In [63]:

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

test_point_index = 100

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

Predicted Class

Frequency of nearest points : Counter([7, 36, 2, 3, 3, 1, 5, 1])

Logistic Regression

with class balancing

Hyper parameter tuning

In [64]:

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

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

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


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


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

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

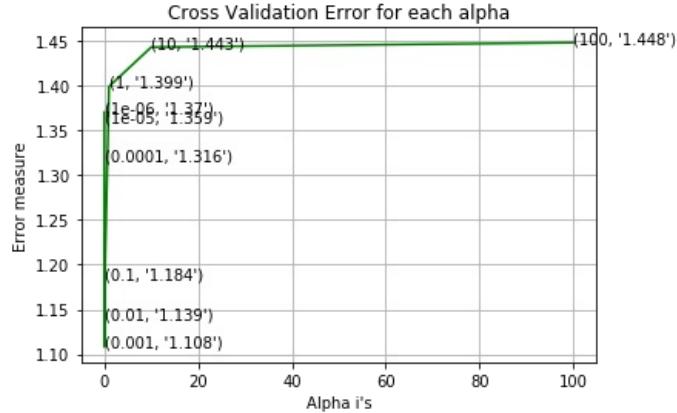
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

```

for alpha = 1e-06
Log Loss : 1.3702490574106325
for alpha = 1e-05
Log Loss : 1.3587623945883234
for alpha = 0.0001
Log Loss : 1.3164784476545648
for alpha = 0.001
Log Loss : 1.108053700492278
for alpha = 0.01
Log Loss : 1.1393391456839073
for alpha = 0.1
Log Loss : 1.1843694841445518
for alpha = 1
Log Loss : 1.3989409861498876
for alpha = 10
Log Loss : 1.442877787991565
for alpha = 100
Log Loss : 1.4479221330657857

```



```

For values of best alpha = 0.001 The train log loss is: 0.5720081045050858
For values of best alpha = 0.001 The cross validation log loss is: 1.108053700492278
For values of best alpha = 0.001 The test log loss is: 1.03351042907678

```

Testing the model with best hyper parameters

In [65]:

```

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

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

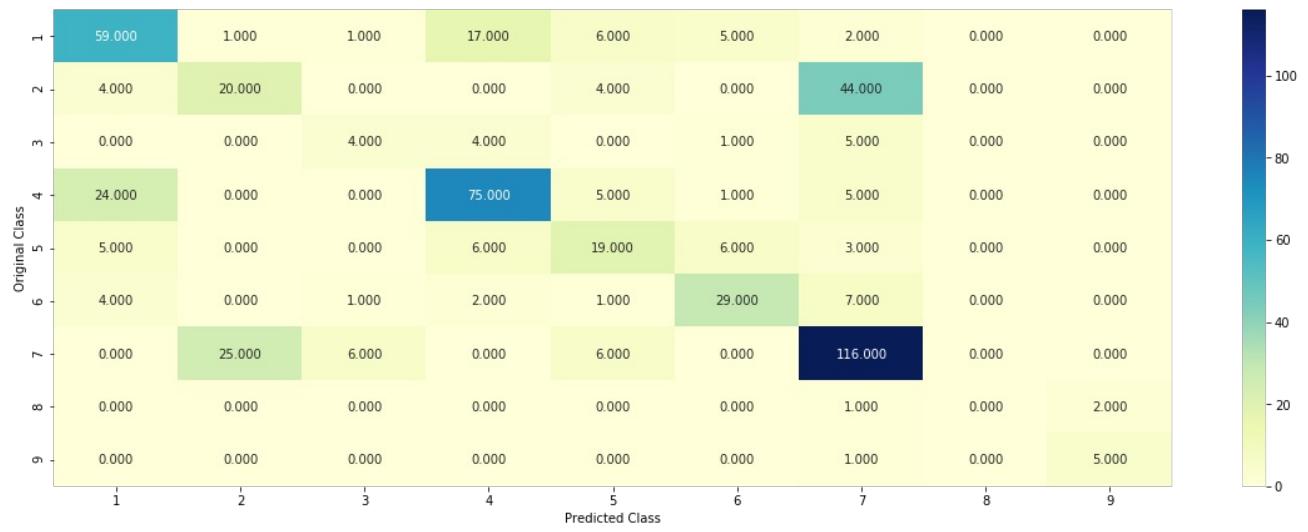
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,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

```

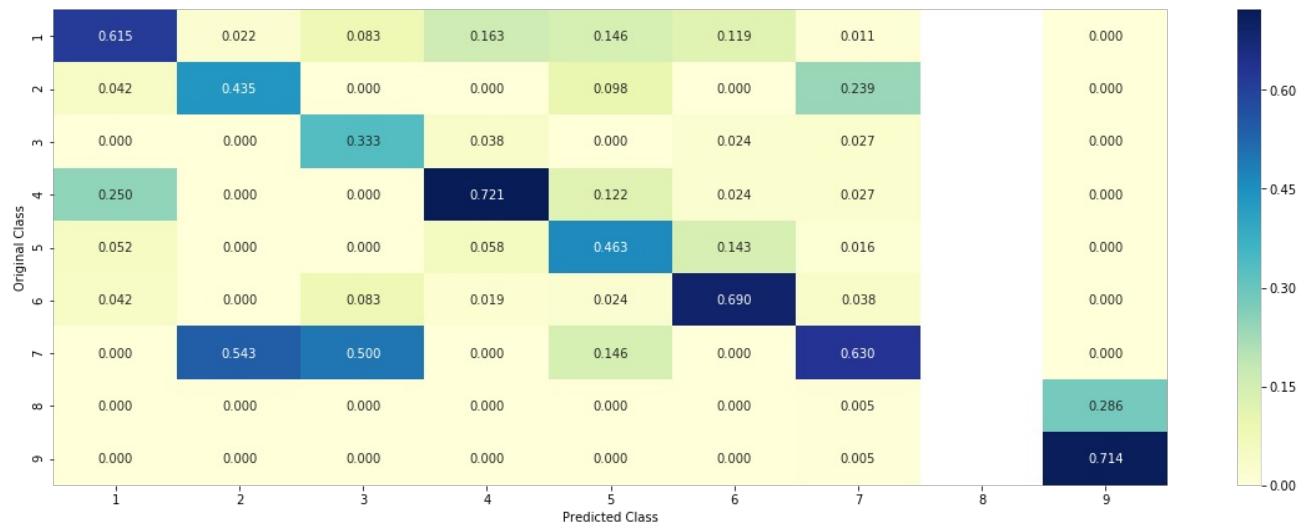
```

Log loss : 1.108053700492278
Number of mis-classified points : 0.38533834586466165
----- Confusion matrix -----

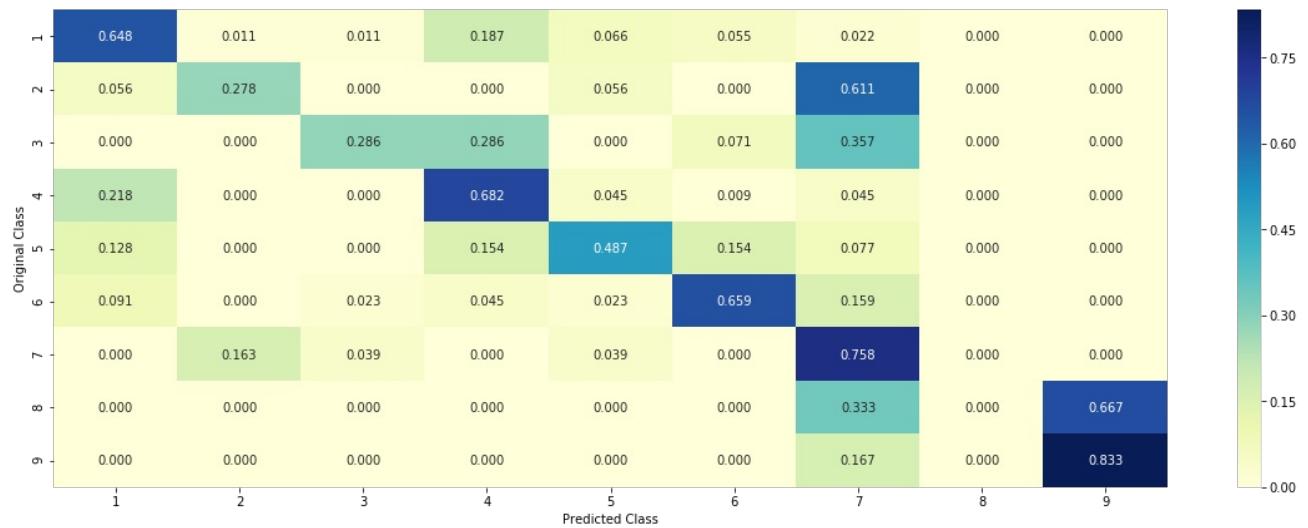
```



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Feature Importance

In [66]:

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

Correctly Classified point

In [67]:

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

```
Predicted Class : 3
Predicted Class Probabilities: [[0.0029 0.008 0.5325 0.0052 0.3709 0.0035 0.0702 0.0051 0.0018]]
Actual Class : 3
-----
292 Text feature [efficiencies] present in test data point [True]
329 Text feature [nherr] present in test data point [True]
382 Text feature [bossi] present in test data point [True]
395 Text feature [palmer] present in test data point [True]
427 Text feature [iwahara] present in test data point [True]
436 Text feature [sch] present in test data point [True]
453 Text feature [2011a] present in test data point [True]
Out of the top 500 features 7 are present in query point
```

Incorrectly Classified point

In [68]:

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

Predicted Class : 7
Predicted Class Probabilities: [[0.0035 0.0139 0.004 0.0039 0.0071 0.003 0.9582 0.0048 0.0016]]
Actual Class : 7

20 Text feature [constitutive] present in test data point [True]
23 Text feature [constitutively] present in test data point [True]
29 Text feature [activated] present in test data point [True]
53 Text feature [mitogen] present in test data point [True]
69 Text feature [ligand] present in test data point [True]
76 Text feature [3t3] present in test data point [True]
93 Text feature [murine] present in test data point [True]
134 Text feature [oncogenes] present in test data point [True]
150 Text feature [oncogene] present in test data point [True]
157 Text feature [activation] present in test data point [True]
175 Text feature [ba] present in test data point [True]
179 Text feature [nude] present in test data point [True]
192 Text feature [manageable] present in test data point [True]
194 Text feature [f3] present in test data point [True]
197 Text feature [phosphorylation] present in test data point [True]
201 Text feature [phospho] present in test data point [True]
209 Text feature [downstream] present in test data point [True]
215 Text feature [allylamino] present in test data point [True]
218 Text feature [transforming] present in test data point [True]
255 Text feature [serum] present in test data point [True]
258 Text feature [transformed] present in test data point [True]
259 Text feature [responsiveness] present in test data point [True]
301 Text feature [treating] present in test data point [True]
325 Text feature [noninvasive] present in test data point [True]
326 Text feature [expressing] present in test data point [True]
330 Text feature [withdrawn] present in test data point [True]
361 Text feature [epidermal] present in test data point [True]
384 Text feature [akt] present in test data point [True]
395 Text feature [reversed] present in test data point [True]
406 Text feature [activating] present in test data point [True]
411 Text feature [il] present in test data point [True]
416 Text feature [kim] present in test data point [True]
434 Text feature [experienced] present in test data point [True]
439 Text feature [brodeur] present in test data point [True]
440 Text feature [sci] present in test data point [True]
441 Text feature [interleukin] present in test data point [True]
452 Text feature [mgh] present in test data point [True]
455 Text feature [adenocarcinoma] present in test data point [True]
491 Text feature [rarely] present in test data point [True]
492 Text feature [cascades] present in test data point [True]
495 Text feature [technology] present in test data point [True]
499 Text feature [activate] present in test data point [True]
Out of the top 500 features 42 are present in query point

Without Class balancing

Hyper parameter tuning

In [69]:

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

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

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


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


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

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

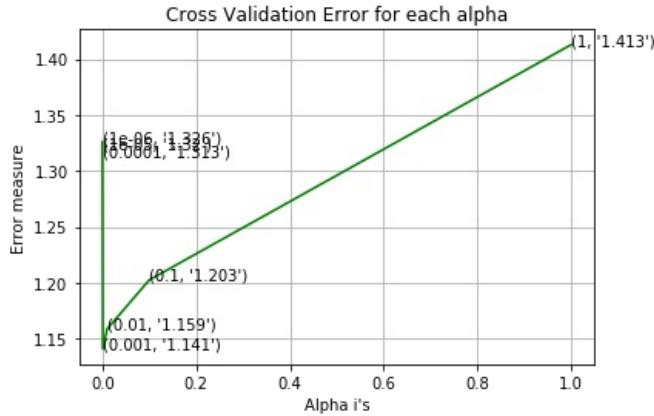
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_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 : 1.326017119977709
for alpha = 1e-05
Log Loss : 1.3199794194673449
for alpha = 0.0001
Log Loss : 1.313258188297746
for alpha = 0.001
Log Loss : 1.1410297893267967
for alpha = 0.01
Log Loss : 1.1587217223411794
for alpha = 0.1
Log Loss : 1.2026275181377326
for alpha = 1
Log Loss : 1.412724061012384

```



For values of best alpha = 0.001 The train log loss is: 0.5662707444066366
 For values of best alpha = 0.001 The cross validation log loss is: 1.1410297893267967
 For values of best alpha = 0.001 The test log loss is: 1.0581644208076553

Testing model with best hyper parameters

In [70]:

```

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

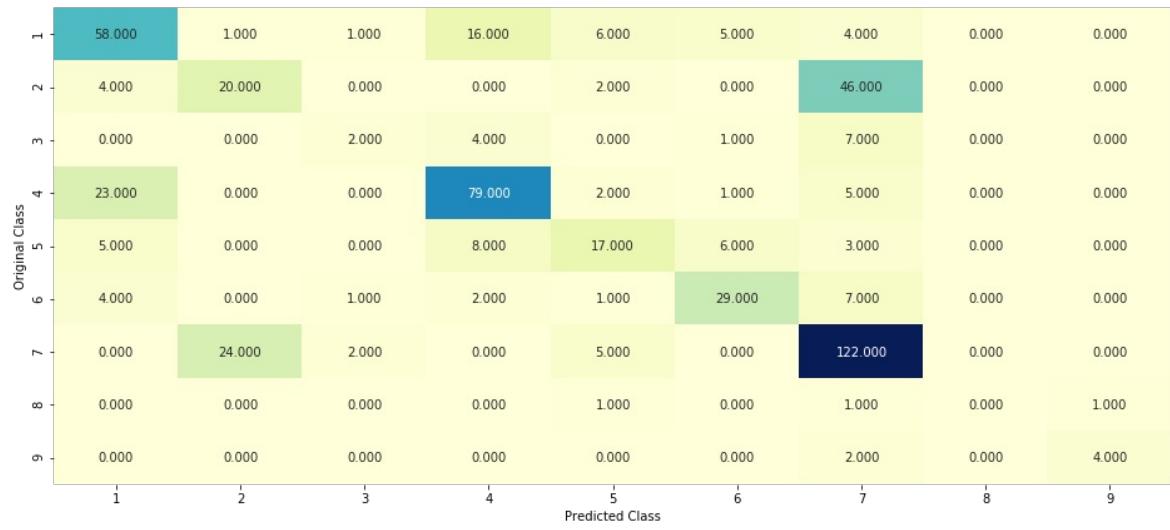
# some of methods
# fit(X, y[, coef_init, intercept_init, ...])      Fit linear model with Stochastic Gradient Descent.
# predict(X)    Predict class labels for samples in X.

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

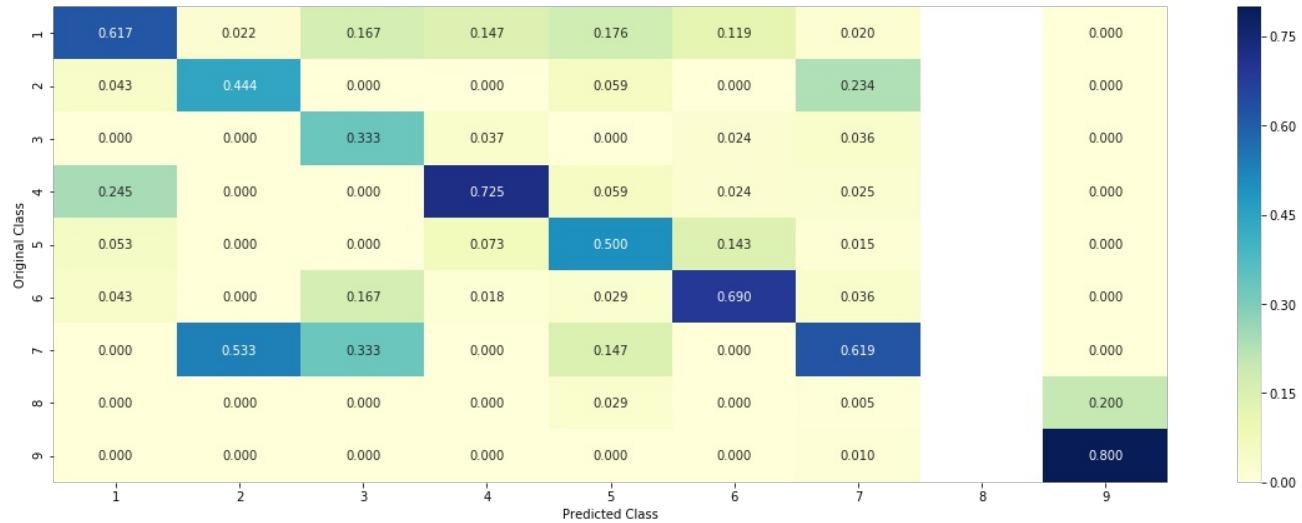
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

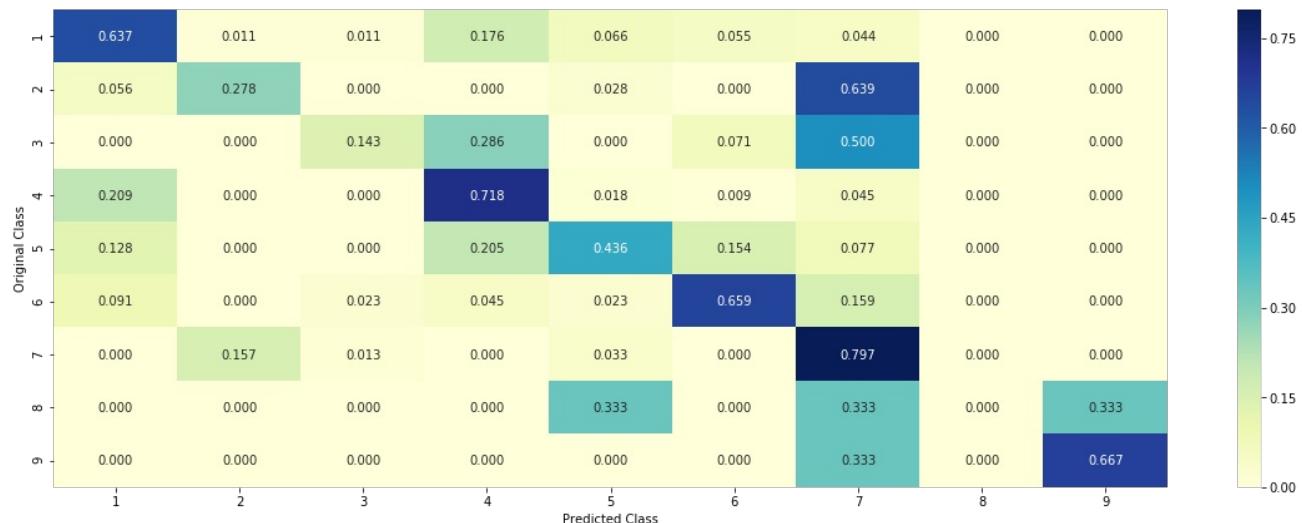
Log loss : 1.1410297893267967
 Number of mis-classified points : 0.37781954887218044
 ----- Confusion matrix -----



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Feature Importance, Correctly Classified point

In [71]:

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

```
Predicted Class : 5
Predicted Class Probabilities: [[0.0066 0.0107 0.395  0.0089 0.4048 0.0035 0.1649 0.0057 0.          ]]
Actual Class : 3
-----
Out of the top 500 features 0 are present in query point
```

Feature Importance, Inorrectly 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[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[5.300e-03 1.650e-02 2.400e-03 5.600e-03 5.700e-03 2.500e-03 9.581e-01
3.900e-03 1.000e-04]]
Actual Class : 7
-----
```

```
59 Text feature [constitutive] present in test data point [True]
68 Text feature [constitutively] present in test data point [True]
91 Text feature [activated] present in test data point [True]
122 Text feature [3t3] present in test data point [True]
136 Text feature [mitogen] present in test data point [True]
154 Text feature [ligand] present in test data point [True]
179 Text feature [phosphorylation] present in test data point [True]
194 Text feature [activation] present in test data point [True]
200 Text feature [manageable] present in test data point [True]
220 Text feature [ba] present in test data point [True]
244 Text feature [phospho] present in test data point [True]
245 Text feature [f3] present in test data point [True]
246 Text feature [expressing] present in test data point [True]
249 Text feature [downstream] present in test data point [True]
256 Text feature [transformed] present in test data point [True]
264 Text feature [murine] present in test data point [True]
275 Text feature [oncogene] present in test data point [True]
292 Text feature [transforming] present in test data point [True]
296 Text feature [serum] present in test data point [True]
300 Text feature [activating] present in test data point [True]
307 Text feature [oncogenes] present in test data point [True]
344 Text feature [brodeur] present in test data point [True]
359 Text feature [reversed] present in test data point [True]
397 Text feature [responsiveness] present in test data point [True]
399 Text feature [epidermal] present in test data point [True]
407 Text feature [activate] present in test data point [True]
409 Text feature [nude] present in test data point [True]
427 Text feature [noninvasive] present in test data point [True]
432 Text feature [interleukin] present in test data point [True]
436 Text feature [treating] present in test data point [True]
440 Text feature [experienced] present in test data point [True]
442 Text feature [kinase] present in test data point [True]
446 Text feature [cells6] present in test data point [True]
472 Text feature [adenocarcinoma] present in test data point [True]
Out of the top 500 features 34 are present in query point
```

Linear Support Vector Machines

Hyper parameter tuning

In [73]:

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

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

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


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


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

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

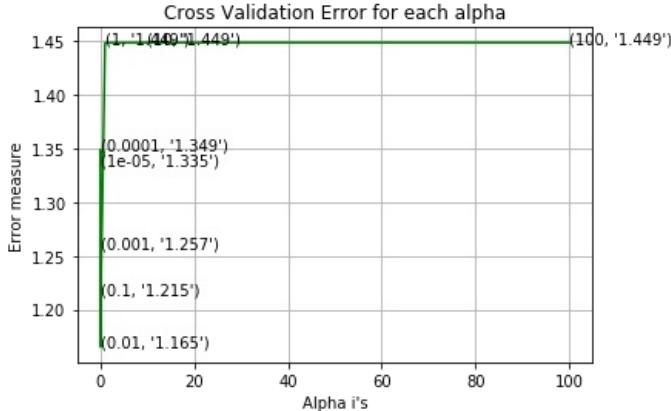
best_alpha = np.argmin(cv_log_error_array)
# clf = SVC(C=i,kernel='linear',probability=True, class_weight='balanced')
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf,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

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

```

for C = 1e-05
Log Loss : 1.3350995350418027
for C = 0.0001
Log Loss : 1.3489754598929704
for C = 0.001
Log Loss : 1.2567680505744536
for C = 0.01
Log Loss : 1.1649730293038694
for C = 0.1
Log Loss : 1.21461654146705
for C = 1
Log Loss : 1.4487091274008266
for C = 10
Log Loss : 1.4489261157597928
for C = 100
Log Loss : 1.448926102328059

```



```

For values of best alpha = 0.01 The train log loss is: 0.6839179989591164
For values of best alpha = 0.01 The cross validation log loss is: 1.1649730293038694
For values of best alpha = 0.01 The test log loss is: 1.0889523314465388

```

Testing model with best hyper parameters

In [74]:

```

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

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

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

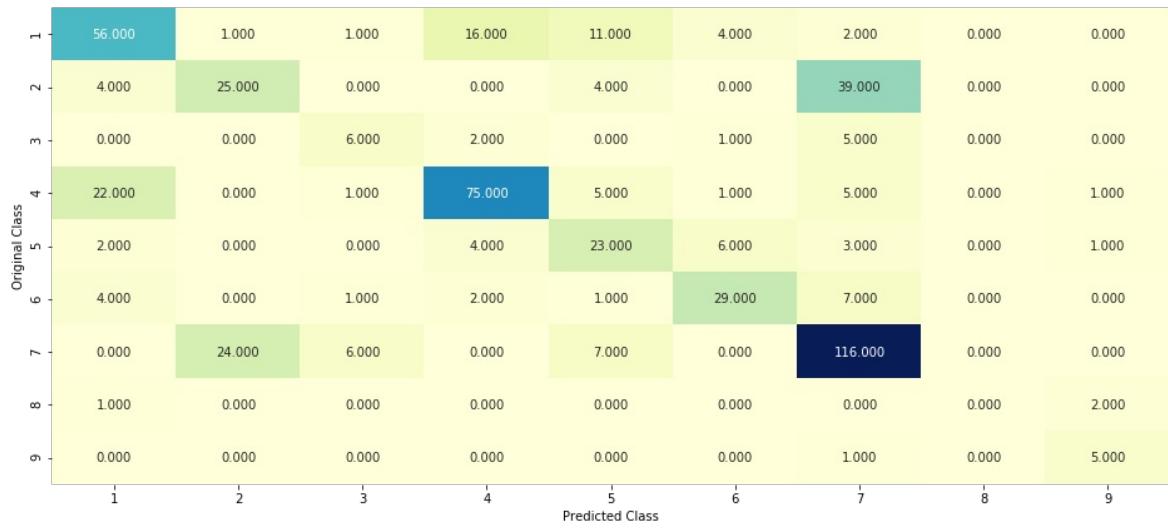

# clf = SVC(C=alpha[best_alpha],kernel='linear',probability=True, class_weight='balanced')
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42, class_weight='balanced')
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

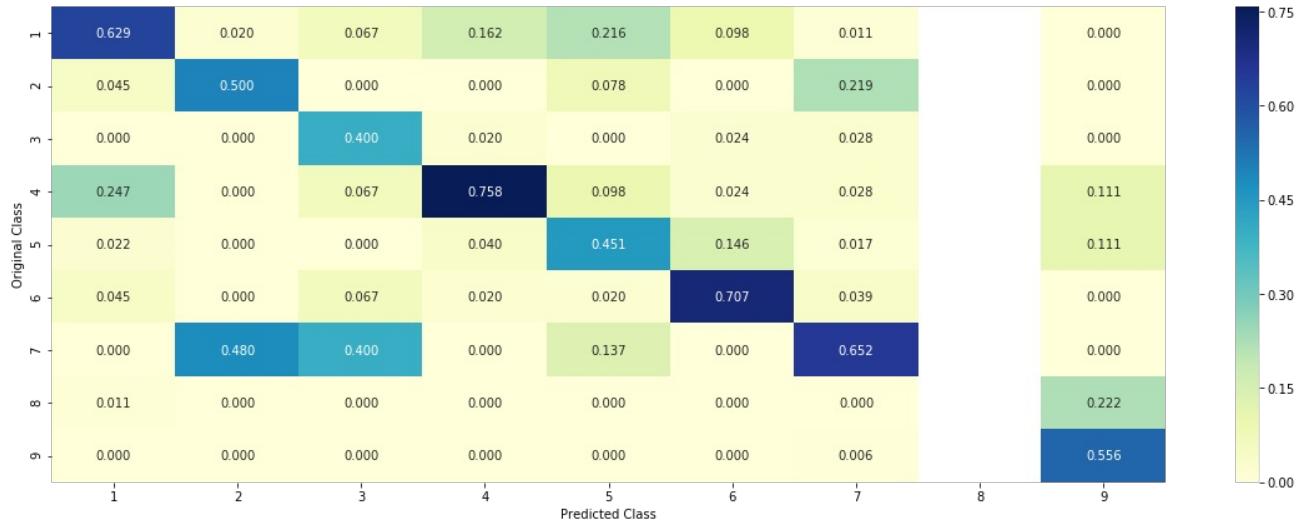
```

Log loss : 1.1649730293038694
Number of mis-classified points : 0.37030075187969924
----- Confusion matrix -----

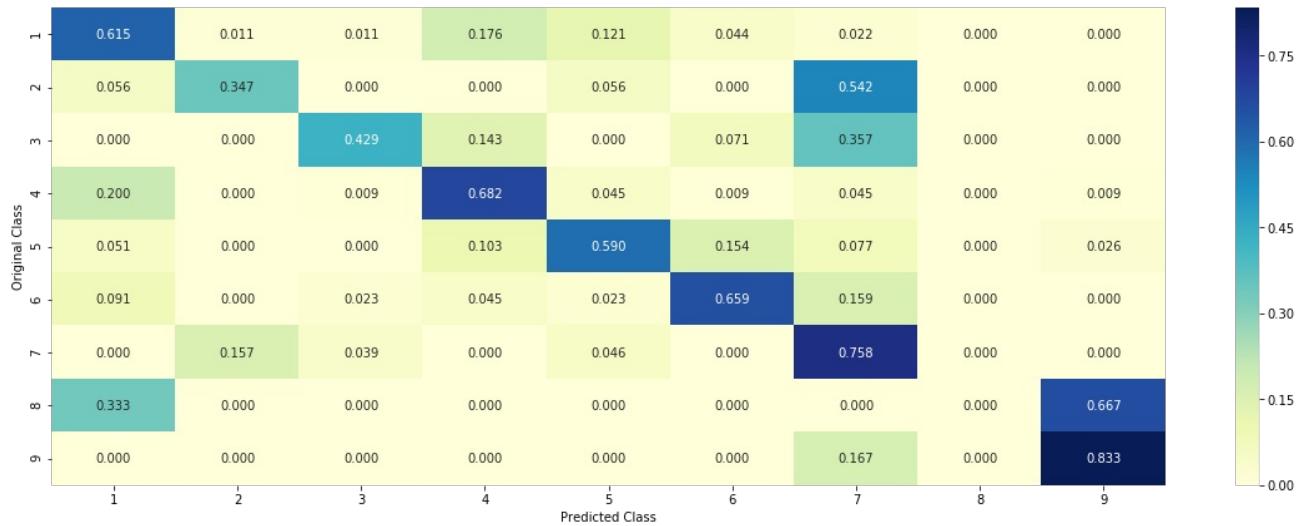
```



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



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



Feature Importance

For Correctly classified point

In [75]:

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

Predicted Class : 3

Predicted Class Probabilities: [[0.0322 0.0182 0.6944 0.0225 0.1773 0.0129 0.0349 0.0049 0.0027]]

Actual Class : 3

```
-----  
114 Text feature [nherr] present in test data point [True]  
128 Text feature [sch] present in test data point [True]  
136 Text feature [2011b] present in test data point [True]  
137 Text feature [bossi] present in test data point [True]  
146 Text feature [2011a] present in test data point [True]  
155 Text feature [palmer] present in test data point [True]  
158 Text feature [iwahara] present in test data point [True]  
312 Text feature [hallberg] present in test data point [True]  
335 Text feature [soda] present in test data point [True]  
451 Text feature [m1166r] present in test data point [True]  
486 Text feature [morris] present in test data point [True]  
Out of the top 500 features 11 are present in query point
```

For 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[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_imfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation'].iloc[test_point_index], no_feature)
```

```
Predicted Class : 7  
Predicted Class Probabilities: [[0.0173 0.0107 0.0058 0.0172 0.0147 0.0117 0.9184 0.0027 0.0014]]
```

```
Actual Class : 7
```

```
# -----  
# default parameters  
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,  
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,  
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,  
# class_weight=None)  
  
# Some of methods of RandomForestClassifier()  
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.  
# predict(X) Perform classification on samples in X.  
# predict_proba (X) Perform classification on samples in X.  
  
# some of attributes of RandomForestClassifier()  
# feature_importances_ : array of shape = [n_features]  
# The feature importances (the higher, the more important the feature).  
# -----  
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-con
```

Random Forest Classifier

Hyper parameter tuning (With One hot Encoding)

```
In [77]:
```

```
# -----  
# default parameters  
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,  
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,  
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,  
# class_weight=None)  
  
# Some of methods of RandomForestClassifier()  
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.  
# predict(X) Perform classification on samples in X.  
# predict_proba (X) Perform classification on samples in X.  
  
# some of attributes of RandomForestClassifier()  
# feature_importances_ : array of shape = [n_features]  
# The feature importances (the higher, the more important the feature).  
# -----  
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-con
```

```

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

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

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

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

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

```

```

for n_estimators = 100 and max depth =  5
Log Loss : 1.1900928555323633
for n_estimators = 100 and max depth =  10
Log Loss : 1.1377002836286825
for n_estimators = 200 and max depth =  5
Log Loss : 1.1722181474983866
for n_estimators = 200 and max depth =  10
Log Loss : 1.127928388626611
for n_estimators = 500 and max depth =  5
Log Loss : 1.1503997040786935
for n_estimators = 500 and max depth =  10
Log Loss : 1.1143460705956068
for n_estimators = 1000 and max depth =  5
Log Loss : 1.147680635588759
for n_estimators = 1000 and max depth =  10
Log Loss : 1.115646111403648
for n_estimators = 2000 and max depth =  5
Log Loss : 1.1460519728187284
for n_estimators = 2000 and max depth =  10
Log Loss : 1.1135619865417052
For values of best estimator =  2000 The train log loss is: 0.6433901042342053
For values of best estimator =  2000 The cross validation log loss is: 1.1135619865417052
For values of best estimator =  2000 The test log loss is: 1.1381637420770192

```

Testing model with best hyper parameters (One Hot Encoding)

In [78]:

```

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

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

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

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

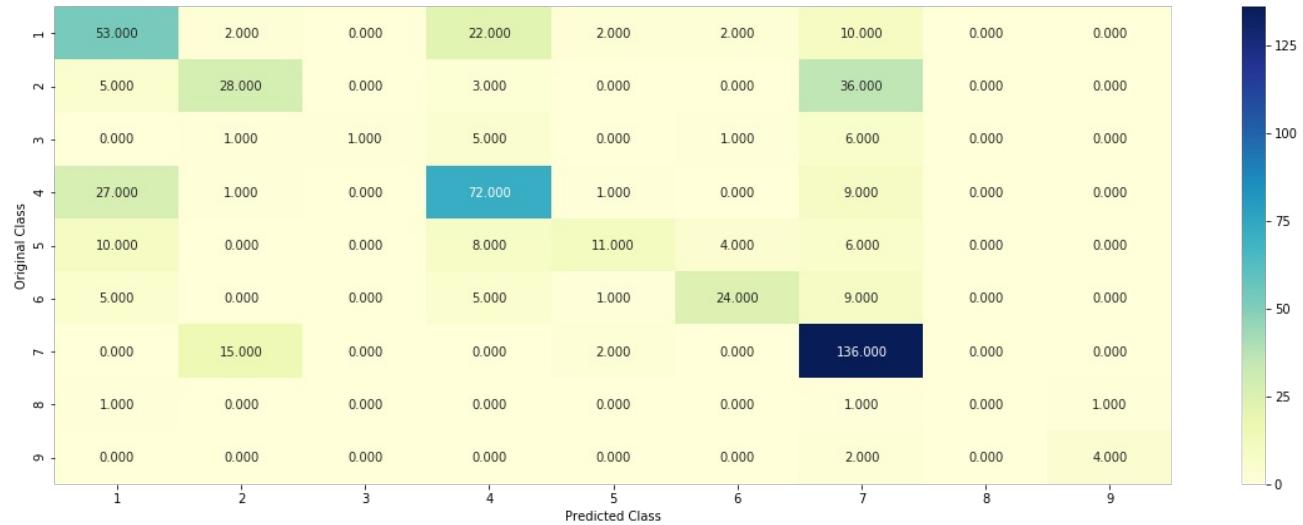

clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha%2)], random_state=42, n_jobs=-1)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

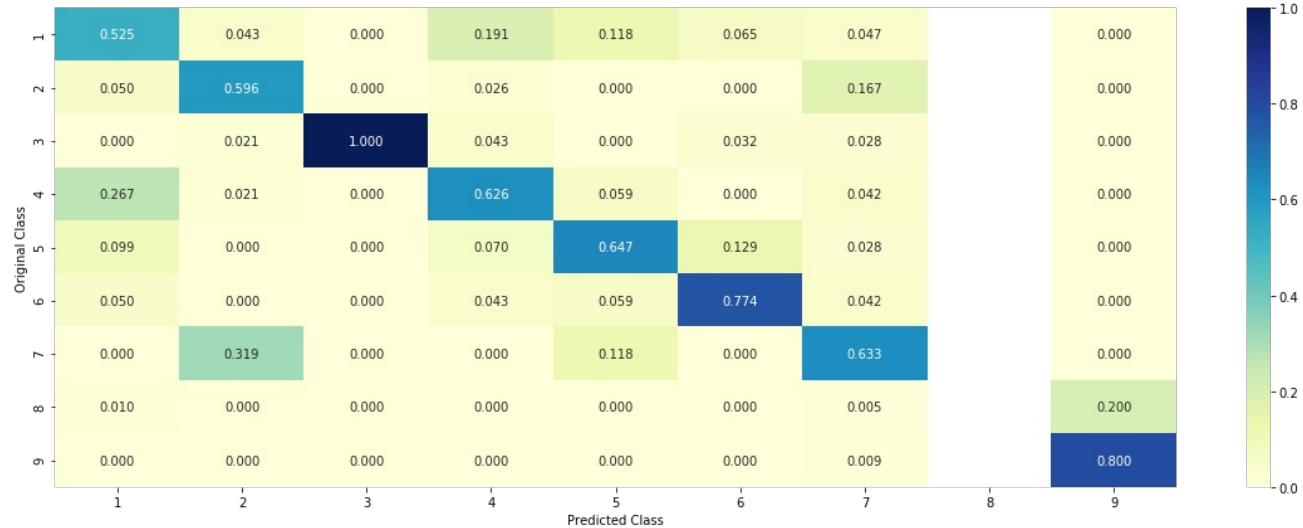
Log loss : 1.1135619865417052

Number of mis-classified points : 0.3815789473684211

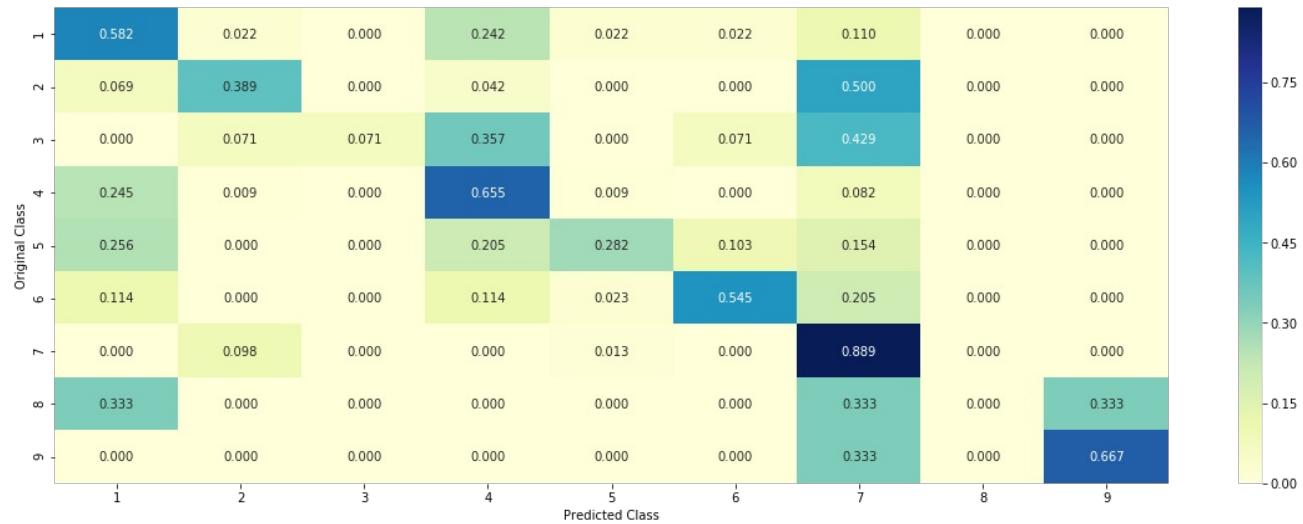
----- Confusion matrix -----



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



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



Feature Importance

Correctly Classified point

In [79]:

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

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

Predicted Class : 7
Predicted Class Probabilities: [[0.0424 0.1581 0.155 0.0376 0.1275 0.0484 0.4209 0.0058 0.0042]]

Actual Class : 3

0 Text feature [tyrosine] present in test data point [True]
1 Text feature [kinase] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [activated] present in test data point [True]
4 Text feature [activation] present in test data point [True]
5 Text feature [phosphorylation] present in test data point [True]
6 Text feature [constitutive] present in test data point [True]
7 Text feature [inhibitors] present in test data point [True]
8 Text feature [suppressor] present in test data point [True]
9 Text feature [inhibitor] present in test data point [True]
10 Text feature [missense] present in test data point [True]
11 Text feature [treatment] present in test data point [True]
12 Text feature [signaling] present in test data point [True]
13 Text feature [function] present in test data point [True]
14 Text feature [drug] present in test data point [True]
17 Text feature [akt] present in test data point [True]
18 Text feature [oncogenic] present in test data point [True]
20 Text feature [loss] present in test data point [True]
21 Text feature [therapy] present in test data point [True]
22 Text feature [receptor] present in test data point [True]
23 Text feature [kinases] present in test data point [True]
24 Text feature [inhibition] present in test data point [True]
25 Text feature [f3] present in test data point [True]
26 Text feature [activate] present in test data point [True]
27 Text feature [transforming] present in test data point [True]
28 Text feature [downstream] present in test data point [True]
30 Text feature [ba] present in test data point [True]
33 Text feature [extracellular] present in test data point [True]
34 Text feature [stability] present in test data point [True]
35 Text feature [deleterious] present in test data point [True]
37 Text feature [growth] present in test data point [True]
38 Text feature [treated] present in test data point [True]
39 Text feature [cells] present in test data point [True]
40 Text feature [ic50] present in test data point [True]
41 Text feature [constitutively] present in test data point [True]
43 Text feature [clinical] present in test data point [True]
46 Text feature [trials] present in test data point [True]
47 Text feature [patients] present in test data point [True]
49 Text feature [protein] present in test data point [True]
50 Text feature [resistance] present in test data point [True]
53 Text feature [3t3] present in test data point [True]
55 Text feature [efficacy] present in test data point [True]
56 Text feature [dose] present in test data point [True]
57 Text feature [transformation] present in test data point [True]
58 Text feature [functional] present in test data point [True]
59 Text feature [autophosphorylation] present in test data point [True]
60 Text feature [harboring] present in test data point [True]
61 Text feature [respond] present in test data point [True]
62 Text feature [predicted] present in test data point [True]
63 Text feature [therapeutic] present in test data point [True]
64 Text feature [sensitive] present in test data point [True]
65 Text feature [nsclc] present in test data point [True]
68 Text feature [cell] present in test data point [True]
69 Text feature [imatinib] present in test data point [True]
70 Text feature [phosphorylated] present in test data point [True]
71 Text feature [ligand] present in test data point [True]
74 Text feature [expression] present in test data point [True]
76 Text feature [phosphatase] present in test data point [True]
77 Text feature [transform] present in test data point [True]
78 Text feature [oncogene] present in test data point [True]
79 Text feature [dna] present in test data point [True]
84 Text feature [variants] present in test data point [True]
85 Text feature [patient] present in test data point [True]
86 Text feature [response] present in test data point [True]
87 Text feature [advanced] present in test data point [True]
89 Text feature [resistant] present in test data point [True]
90 Text feature [amplification] present in test data point [True]
94 Text feature [proteins] present in test data point [True]
96 Text feature [trial] present in test data point [True]
98 Text feature [independence] present in test data point [True]

Out of the top 100 features 70 are present in query point

Inorrectly Classified point

In [80]:

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

Predicted Class : 7

Predicted Class Probabilities: [[0.0188 0.1422 0.014 0.0164 0.0318 0.0264 0.7445 0.0036 0.0023]]

Actual Class : 7

0 Text feature [tyrosine] present in test data point [True]
1 Text feature [kinase] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [activated] present in test data point [True]
4 Text feature [activation] present in test data point [True]
5 Text feature [phosphorylation] present in test data point [True]
6 Text feature [constitutive] present in test data point [True]
7 Text feature [inhibitors] present in test data point [True]
8 Text feature [suppressor] present in test data point [True]
9 Text feature [inhibitor] present in test data point [True]
10 Text feature [missense] present in test data point [True]
11 Text feature [treatment] present in test data point [True]
12 Text feature [signaling] present in test data point [True]
13 Text feature [function] present in test data point [True]
14 Text feature [drug] present in test data point [True]
15 Text feature [phospho] present in test data point [True]
17 Text feature [akt] present in test data point [True]
18 Text feature [oncogenic] present in test data point [True]
19 Text feature [erk] present in test data point [True]
20 Text feature [loss] present in test data point [True]
21 Text feature [therapy] present in test data point [True]
22 Text feature [receptor] present in test data point [True]
23 Text feature [kinases] present in test data point [True]
24 Text feature [inhibition] present in test data point [True]
25 Text feature [f3] present in test data point [True]
26 Text feature [activate] present in test data point [True]
27 Text feature [transforming] present in test data point [True]
28 Text feature [downstream] present in test data point [True]
30 Text feature [ba] present in test data point [True]
32 Text feature [expressing] present in test data point [True]
33 Text feature [extracellular] present in test data point [True]
36 Text feature [proliferation] present in test data point [True]
37 Text feature [growth] present in test data point [True]
38 Text feature [treated] present in test data point [True]
39 Text feature [cells] present in test data point [True]
40 Text feature [ic50] present in test data point [True]
41 Text feature [constitutively] present in test data point [True]
43 Text feature [clinical] present in test data point [True]
46 Text feature [trials] present in test data point [True]
47 Text feature [patients] present in test data point [True]
48 Text feature [inhibited] present in test data point [True]
49 Text feature [protein] present in test data point [True]
50 Text feature [resistance] present in test data point [True]
51 Text feature [serum] present in test data point [True]
53 Text feature [3t3] present in test data point [True]
55 Text feature [efficacy] present in test data point [True]
56 Text feature [dose] present in test data point [True]
58 Text feature [functional] present in test data point [True]
59 Text feature [autophosphorylation] present in test data point [True]
60 Text feature [harboring] present in test data point [True]
62 Text feature [predicted] present in test data point [True]
63 Text feature [therapeutic] present in test data point [True]
64 Text feature [sensitive] present in test data point [True]
65 Text feature [nsclc] present in test data point [True]
66 Text feature [mitogen] present in test data point [True]
67 Text feature [mek] present in test data point [True]
68 Text feature [cell] present in test data point [True]
69 Text feature [imatinib] present in test data point [True]
70 Text feature [phosphorylated] present in test data point [True]
71 Text feature [ligand] present in test data point [True]
72 Text feature [months] present in test data point [True]
74 Text feature [expression] present in test data point [True]
78 Text feature [oncogene] present in test data point [True]
79 Text feature [dna] present in test data point [True]
81 Text feature [factor] present in test data point [True]
83 Text feature [il] present in test data point [True]
85 Text feature [patient] present in test data point [True]
86 Text feature [response] present in test data point [True]
87 Text feature [advanced] present in test data point [True]
89 Text feature [resistant] present in test data point [True]
90 Text feature [amplification] present in test data point [True]
91 Text feature [kit] present in test data point [True]
94 Text feature [proteins] present in test data point [True]
96 Text feature [trial] present in test data point [True]
97 Text feature [retained] present in test data point [True]
Out of the top 100 features 75 are present in query point

Hyper parameter tuning (With Response Coding)

In [81]:

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

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

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

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


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


alpha = [10,50,100,200,500,1000]
max_depth = [2,3,5,10]
cv_log_error_array = []
for i in alpha:
    for j in max_depth:
        print("for n_estimators =", i,"and max depth = ", j)
        clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
        clf.fit(train_x_responseCoding, train_y)
        sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
        sig_clf.fit(train_x_responseCoding, train_y)
        sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
        cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
        print("Log Loss :",log_loss(cv_y, sig_clf_probs))
    ...
fig, ax = plt.subplots()
features = np.dot(np.array(alpha)[:,None],np.array(max_depth)[None]).ravel()
ax.plot(features, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[int(i/4)],max_depth[int(i%4)],str(txt)), (features[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()
```

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

predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
```

```

predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The test log loss is:",log_loss(y_test, predict_y
, labels=clf.classes_, eps=1e-15))

for n_estimators = 10 and max depth = 2
Log Loss : 2.1388879082362946
for n_estimators = 10 and max depth = 3
Log Loss : 1.6949502774277017
for n_estimators = 10 and max depth = 5
Log Loss : 1.5836733522621822
for n_estimators = 10 and max depth = 10
Log Loss : 2.3220125606556192
for n_estimators = 50 and max depth = 2
Log Loss : 1.5852386298579557
for n_estimators = 50 and max depth = 3
Log Loss : 1.4429884988576251
for n_estimators = 50 and max depth = 5
Log Loss : 1.4297692126443675
for n_estimators = 50 and max depth = 10
Log Loss : 1.9170427756135155
for n_estimators = 100 and max depth = 2
Log Loss : 1.523605303655311
for n_estimators = 100 and max depth = 3
Log Loss : 1.4222160000336952
for n_estimators = 100 and max depth = 5
Log Loss : 1.3436531414393724
for n_estimators = 100 and max depth = 10
Log Loss : 1.811447573224995
for n_estimators = 200 and max depth = 2
Log Loss : 1.565430228704022
for n_estimators = 200 and max depth = 3
Log Loss : 1.5067927383037305
for n_estimators = 200 and max depth = 5
Log Loss : 1.379134323242554
for n_estimators = 200 and max depth = 10
Log Loss : 1.7587740208806495
for n_estimators = 500 and max depth = 2
Log Loss : 1.6321601560582752
for n_estimators = 500 and max depth = 3
Log Loss : 1.5439564825025756
for n_estimators = 500 and max depth = 5
Log Loss : 1.3958029311712263
for n_estimators = 500 and max depth = 10
Log Loss : 1.8293437905857932
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6051519974680264
for n_estimators = 1000 and max depth = 3
Log Loss : 1.552077556263543
for n_estimators = 1000 and max depth = 5
Log Loss : 1.405224347590736
for n_estimators = 1000 and max depth = 10
Log Loss : 1.8091736424197091
For values of best alpha = 100 The train log loss is: 0.0566596918920669
For values of best alpha = 100 The cross validation log loss is: 1.3436531414393722
For values of best alpha = 100 The test log loss is: 1.4014686784217025

```

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

In [82]:

```

default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrea
se=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=Fa
lse,
class_weight=None)

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

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

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

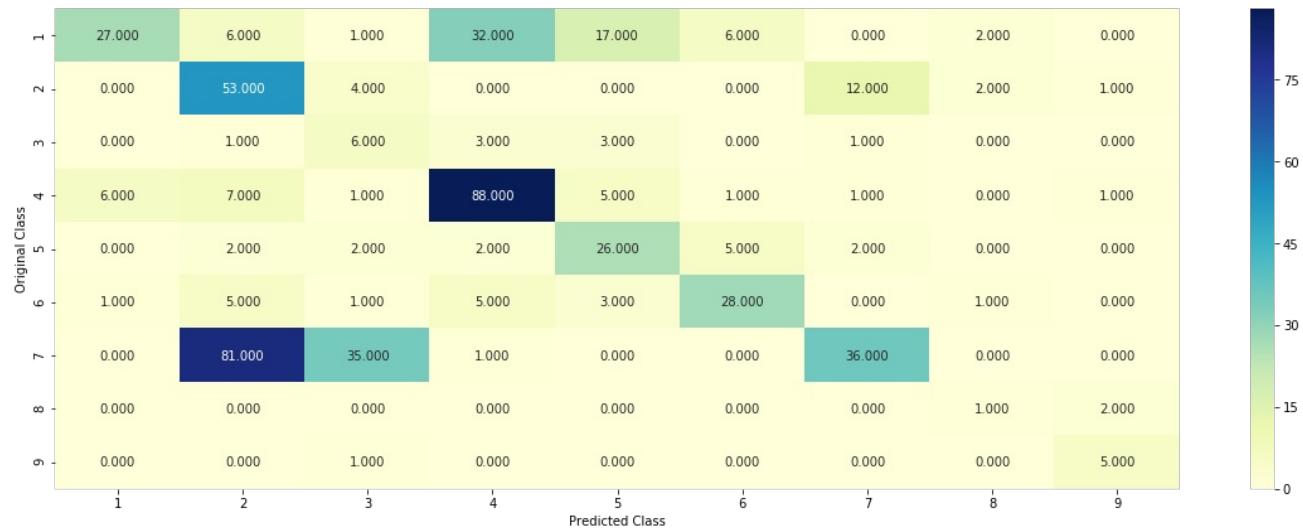
```

clf = RandomForestClassifier(max\_depth=max\_depth[int(best\_alpha%4)], n\_estimators=alpha[int(best\_alpha/4)], crite
rion='gini', max\_features='auto',random\_state=42)
predict\_and\_plot\_confusion\_matrix(train\_x\_responseCoding, train\_y, cv\_x\_responseCoding, cv\_y, clf)

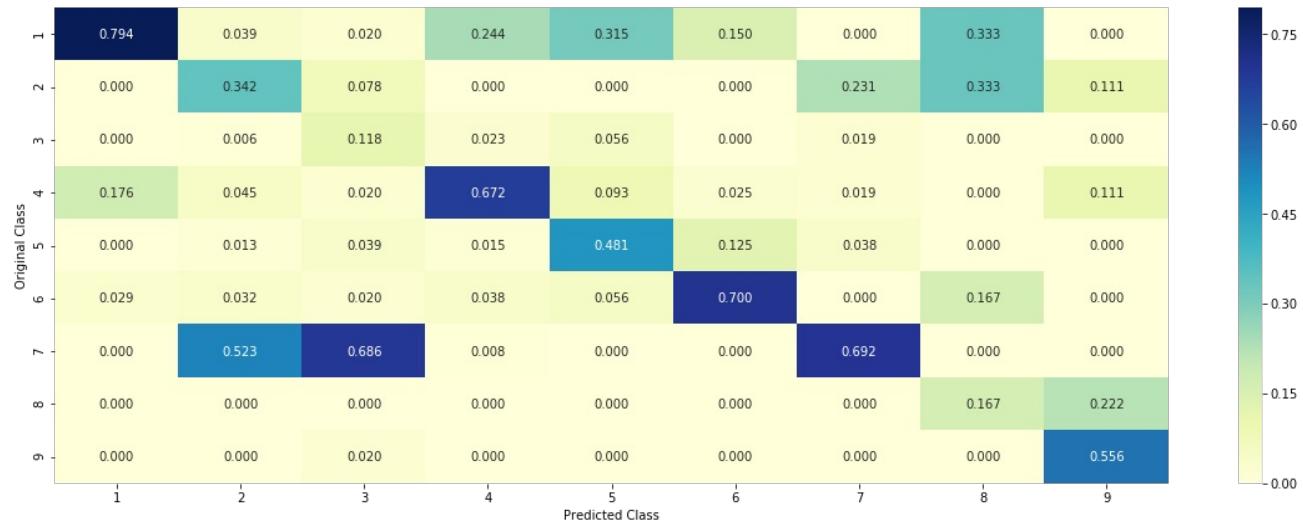
Log loss : 1.343653141439372

Number of mis-classified points : 0.4924812030075188

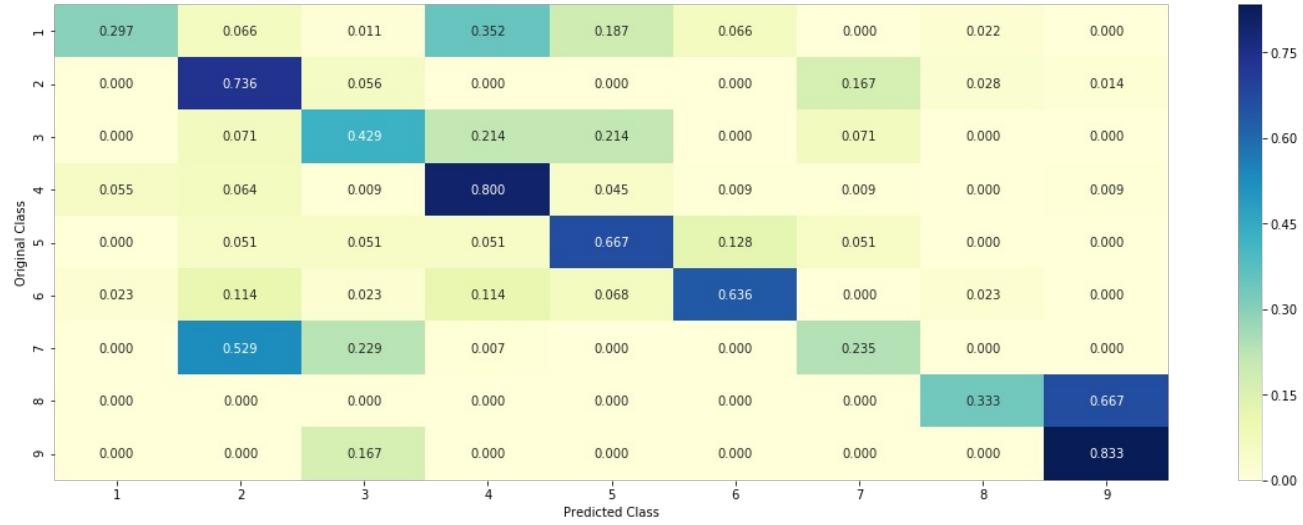
----- Confusion matrix -----



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



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



## Feature Importance

Correctly Classified point

In [83]:

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

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

Predicted Class : 3  
Predicted Class Probabilities: [[0.0056 0.0387 0.8214 0.0071 0.0287 0.0132 0.0726 0.0089 0.0038]]  
Actual Class : 3

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

**Incorrectly Classified point**

In [84]:

```
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].re
shape(1,-1)),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
for i in indices:
 if i<9:
 print("Gene is important feature")
 elif i<18:
 print("Variation is important feature")
 else:
 print("Text is important feature")
```

```
Predicted Class : 3
Predicted Class Probabilities: [[0.0056 0.0387 0.8214 0.0071 0.0287 0.0132 0.0726 0.0089 0.0038]]
Actual Class : 3
```

```

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

## Stack the models

testing with hyper parameter tuning

In [85]:

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

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

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

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

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

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

```

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

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

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

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

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

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

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

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

clf2 = SGDClassifier(alpha=1, penalty='l2', loss='hinge', class_weight='balanced', random_state=0)
clf2.fit(train_x_onehotCoding, train_y)
sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")

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

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

Logistic Regression : Log Loss: 1.06
Support vector machines : Log Loss: 1.45
Naive Bayes : Log Loss: 1.27

Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.178
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.035
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.502
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.112
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.205
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.422

```

**testing the model with the best hyper parameters**

In [86]:

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

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

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

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

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

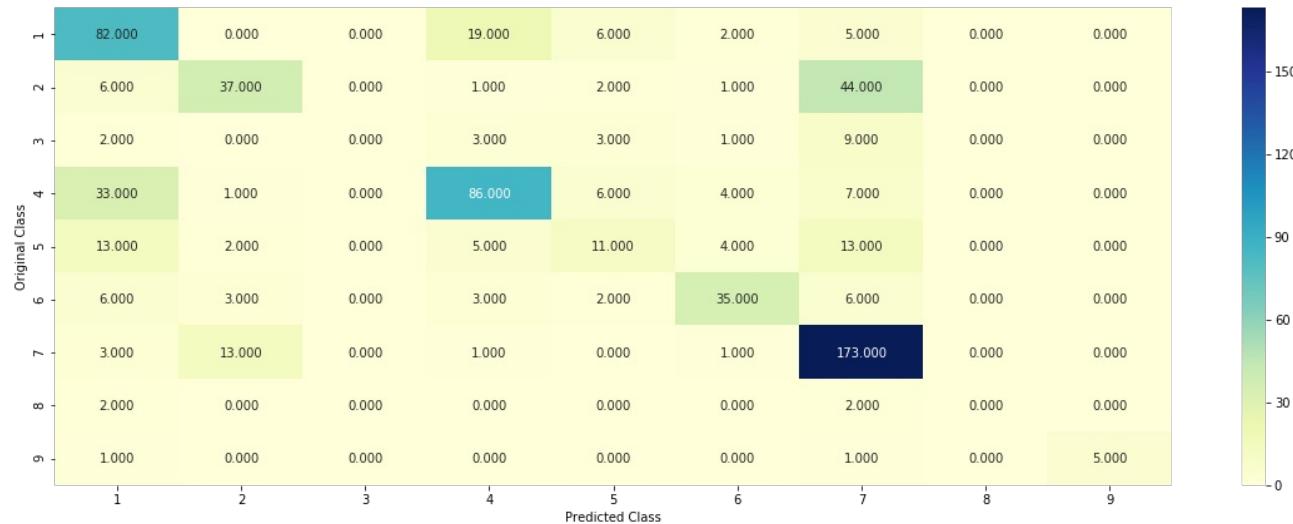
Log loss (train) on the stacking classifier : 0.6345194655828197

Log loss (CV) on the stacking classifier : 1.111637205618287

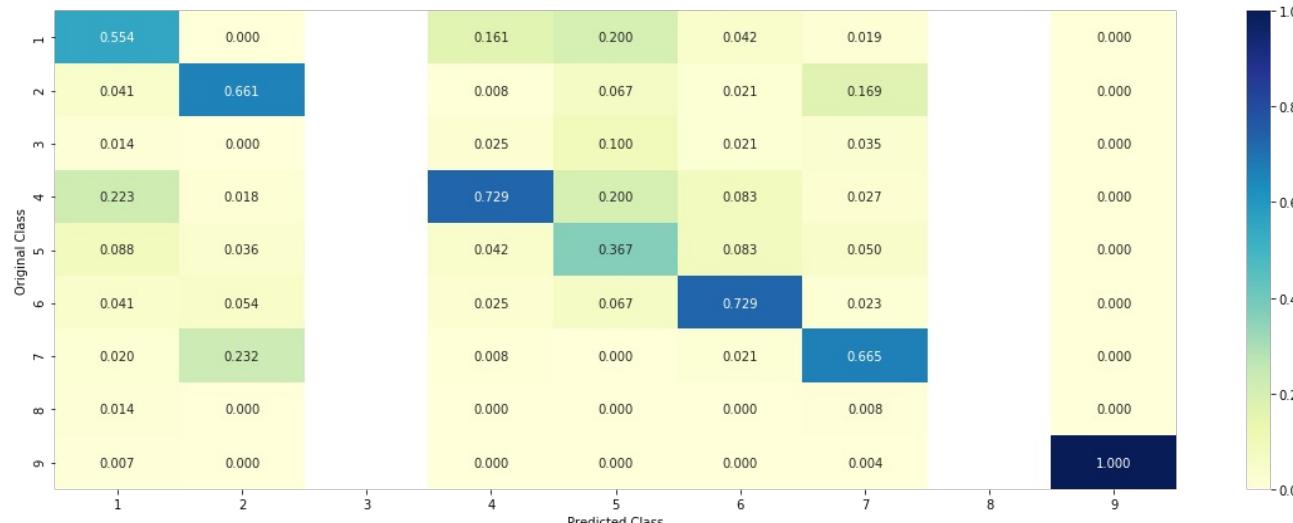
Log loss (test) on the stacking classifier : 1.105867898328202

Number of missclassified point : 0.3548872180451128

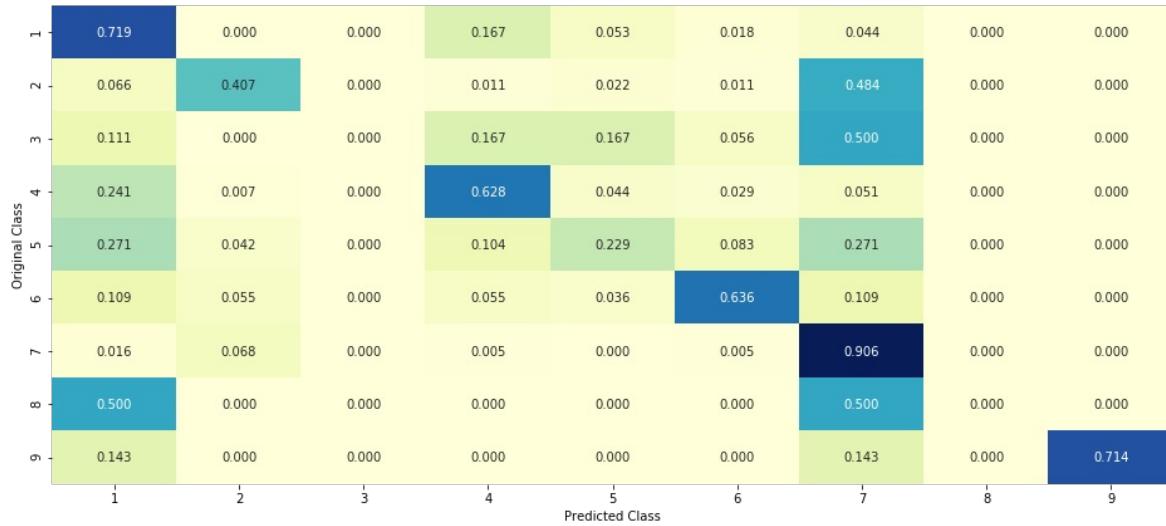
----- Confusion matrix -----



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



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



## Maximum Voting classifier

In [87]:

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

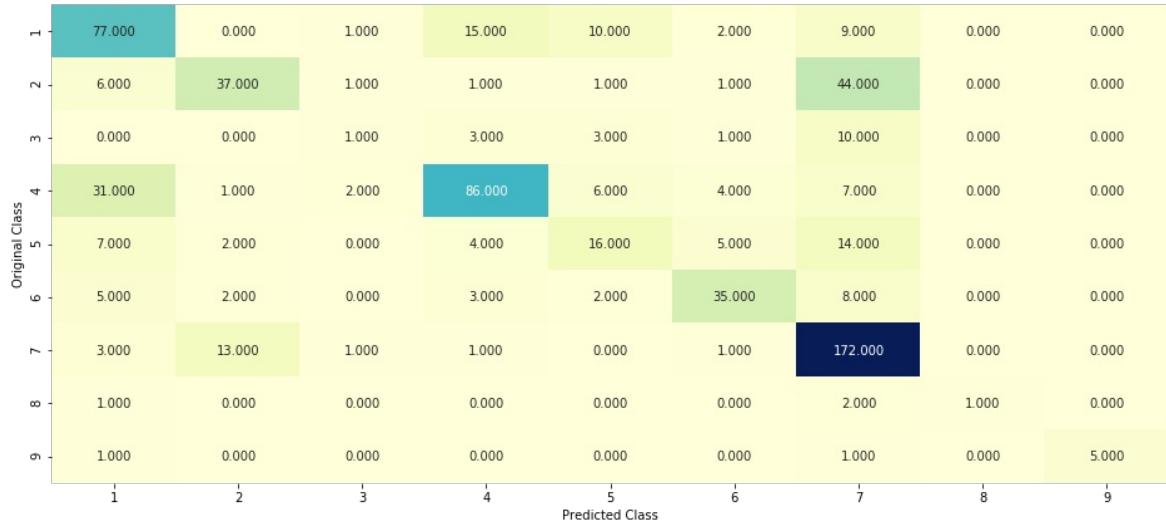
Log loss (train) on the VotingClassifier : 0.862199808726066

Log loss (CV) on the VotingClassifier : 1.1355536397409507

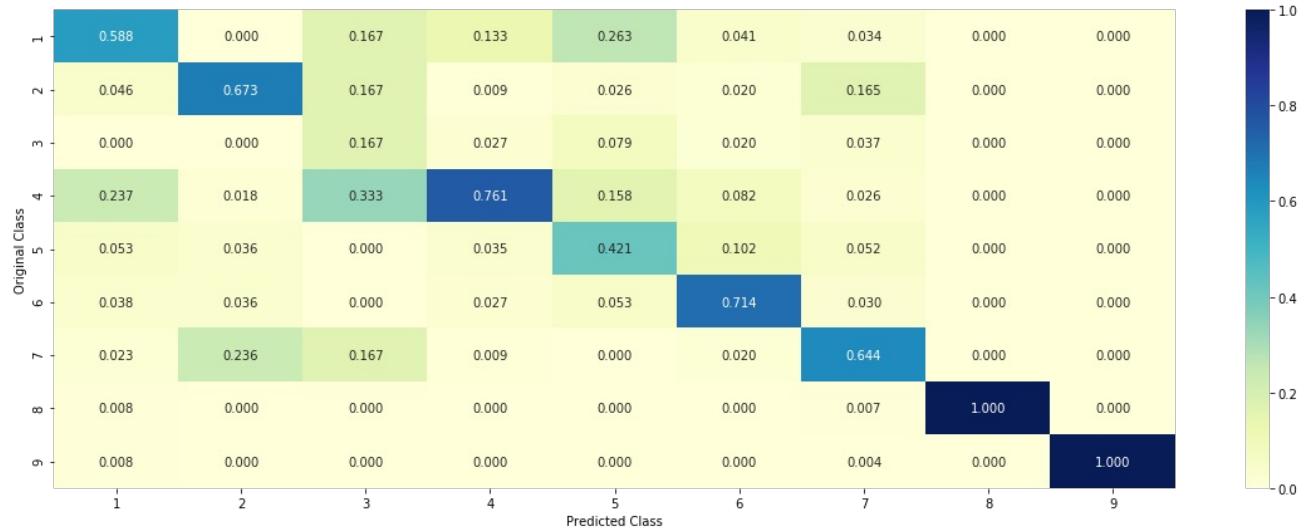
Log loss (test) on the VotingClassifier : 1.1371114955802752

Number of missclassified point : 0.3533834586466165

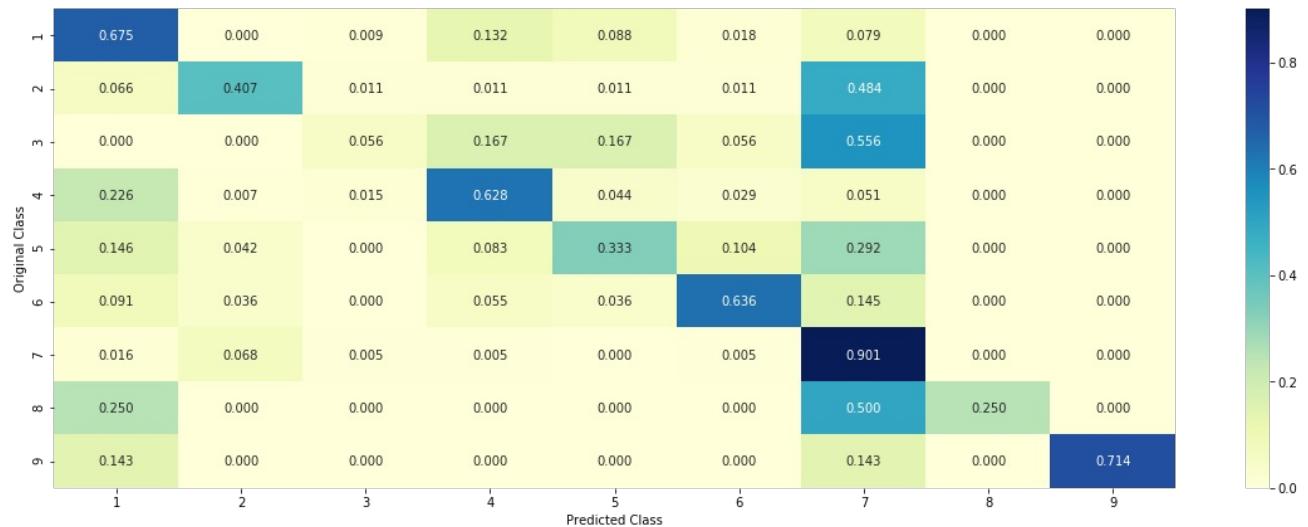
----- Confusion matrix -----



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



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



## 2. Applying All the models with tfidf and taking top 1000 words of tfidf

### Univariate analysis on gene feature

#### One hot encoding on gene feature

In [22]:

```
one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer(max_features=1000)
train_gene_feature_onehotCoding_1000 = gene_vectorizer.fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding_1000 = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_onehotCoding_1000 = gene_vectorizer.transform(cv_df['Gene'])
```

In [23]:

```
print("train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature:", train_gene_feature_onehotCoding_1000.shape)
```

train\_gene\_feature\_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 236)

### LR on gene feature

In [17]:

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

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

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

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

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

cv_log_error_array=[]
for i in alpha:
 clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
 clf.fit(train_gene_feature_onehotCoding_1000, y_train)
 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
 sig_clf.fit(train_gene_feature_onehotCoding_1000, y_train)
 predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding_1000)
 cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
 print('For values of alpha = ', i, "The log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-1
5))

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()

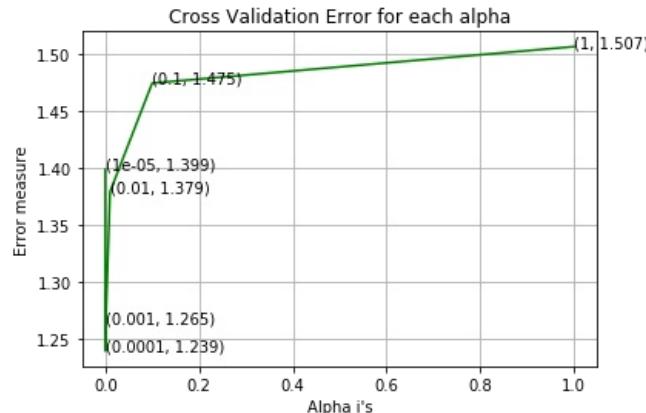
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_gene_feature_onehotCoding_1000, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_gene_feature_onehotCoding_1000, y_train)

predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding_1000)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, lab
els=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding_1000)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predic
t_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding_1000)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, label
s=clf.classes_, eps=1e-15))
```

```

For values of alpha = 1e-05 The log loss is: 1.399013546841295
For values of alpha = 0.0001 The log loss is: 1.2394500545181077
For values of alpha = 0.001 The log loss is: 1.2647607847834386
For values of alpha = 0.01 The log loss is: 1.3792664414674847
For values of alpha = 0.1 The log loss is: 1.4747531588885514
For values of alpha = 1 The log loss is: 1.5068348038842083

```



```

For values of best alpha = 0.0001 The train log loss is: 1.0436041419683073
For values of best alpha = 0.0001 The cross validation log loss is: 1.2394500545181077
For values of best alpha = 0.0001 The test log loss is: 1.1677261214385788

```

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

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.shape[0], " genes in t
rain dataset?")

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

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

```

Q6. How many data points in Test and CV datasets are covered by the 230 genes in train dataset?

Ans

1. In test data 647 out of 665 : 97.29323308270676
2. In cross validation data 512 out of 532 : 96.2406015037594

## Univariate analysis on variation feature

### One hot encoding on variation feature

In [24]:

```

one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer(max_features=1000)
train_variation_feature_onehotCoding_1000 = variation_vectorizer.fit_transform(train_df['Variation'])
test_variation_feature_onehotCoding_1000 = variation_vectorizer.transform(test_df['Variation'])
cv_variation_feature_onehotCoding_1000 = variation_vectorizer.transform(cv_df['Variation'])

```

In [25]:

```

print("train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape o
f Variation feature:", train_variation_feature_onehotCoding_1000.shape)

```

```

train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The s
hape of Variation feature: (2124, 1000)

```

### LR on variation feature

In [31]:

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

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

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

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

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

cv_log_error_array=[]
for i in alpha:
 clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
 clf.fit(train_variation_feature_onehotCoding_1000, y_train)

 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
 sig_clf.fit(train_variation_feature_onehotCoding_1000, y_train)
 predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding_1000)

 cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
 print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
5)

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()

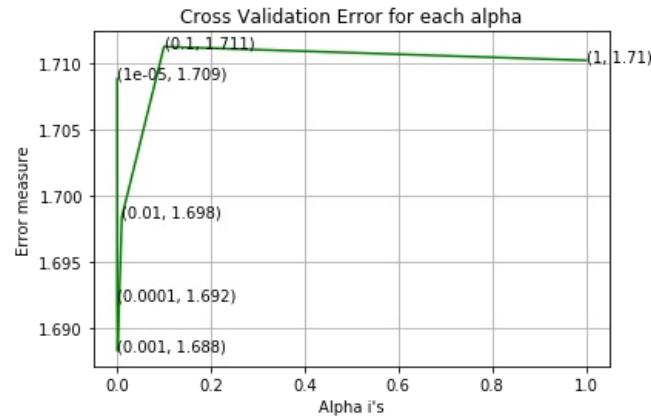
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_variation_feature_onehotCoding_1000, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_variation_feature_onehotCoding_1000, y_train)

predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding_1000)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding_1000)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding_1000)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```

For values of alpha = 1e-05 The log loss is: 1.7088400289253183
For values of alpha = 0.0001 The log loss is: 1.6920937183566835
For values of alpha = 0.001 The log loss is: 1.6882206162503572
For values of alpha = 0.01 The log loss is: 1.6983811232066328
For values of alpha = 0.1 The log loss is: 1.7112794144393952
For values of alpha = 1 The log loss is: 1.7102388757687779

```



```

For values of best alpha = 0.001 The train log loss is: 1.317740597514664
For values of best alpha = 0.001 The cross validation log loss is: 1.6882206162503572
For values of best alpha = 0.001 The test log loss is: 1.7324710060373456

```

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

No it is not stable

In [32]:

```

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

```

Q12. How many data points are covered by total 1938 genes in test and cross validation data sets?

Ans

1. In test data 67 out of 665 : 10.075187969924812
2. In cross validation data 64 out of 532 : 12.030075187969924

## Univariate Analysis on Text Feature

### One hot coding on text feature

In [26]:

```

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 [27]:

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

In [28]:

```
building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer(max_features=1000,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 features) vector
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 : 1000

In [29]:

```
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 [30]:

```
#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 [31]:

```
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 [32]:

```
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 [33]:

```
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()))
```

LR on text feature

In [42]:

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

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

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

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

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

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

 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
 sig_clf.fit(train_text_feature_onehotCoding, y_train)
 predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
 cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
 print('For values of alpha = ', i, "The log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-1
5))

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

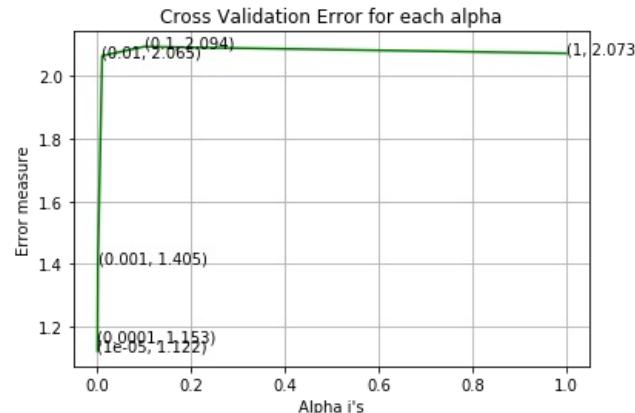
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_text_feature_onehotCoding, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_text_feature_onehotCoding, y_train)

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

```

For values of alpha = 1e-05 The log loss is: 1.1216849419901498
For values of alpha = 0.0001 The log loss is: 1.1533105868561473
For values of alpha = 0.001 The log loss is: 1.404907655614507
For values of alpha = 0.01 The log loss is: 2.0650426794148693
For values of alpha = 0.1 The log loss is: 2.0942532724317346
For values of alpha = 1 The log loss is: 2.0729077112453034

```



```

For values of best alpha = 1e-05 The train log loss is: 0.7696054002219983
For values of best alpha = 1e-05 The cross validation log loss is: 1.1216849419901498
For values of best alpha = 1e-05 The test log loss is: 1.1490287046356498

```

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

Yes, it seems like!

In [43]:

```

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

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

```

In [44]:

```

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

```

3.452 % of word of test data appeared in train data  
3.879 % of word of Cross Validation appeared in train data

## Machine Learning Models

In [45]:

```

#Data preparation for ML models.

#Misc. functionns for ML models

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

 # for calculating log_loss we will provide the array of probabilities belongs to each class
 print("Log loss :",log_loss(test_y, sig_clf.predict_proba(test_x)))
 # calculating the number of data points that are misclassified
 print("Number of mis-classified points :", np.count_nonzero((pred_y- test_y))/test_y.shape[0])
 plot_confusion_matrix(test_y, pred_y)

```

In [46]:

```
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 [47]:

```
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_imptfeature_names(indices, text, gene, var, no_features):
 gene_count_vec = CountVectorizer()
 var_count_vec = CountVectorizer()
 text_count_vec = CountVectorizer(min_df=3)

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

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

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

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

## Stacking the three types of features

In [35]:

```
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_1000,train_variation_feature_onehotCoding_1000))
test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding_1000,test_variation_feature_onehotCoding_1000))
cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding_1000, cv_variation_feature_onehotCoding_1000))

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

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

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

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

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

In [36]:

```
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, 2236)
(number of data points * number of features) in test data = (665, 2236)
(number of data points * number of features) in cross validation data = (532, 2236)
```

In [37]:

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

Response encoding features :

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

## Base Line Model

### Naive Bayes

### Hyper parameter tuning

In [51]:

```
find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html

default paramters
sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)

some of methods of MultinomialNB()
fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
predict(X) Perform classification on an array of test vectors X.
predict_log_proba(X) Return log-probability estimates for the test vector X.

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification

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

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

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

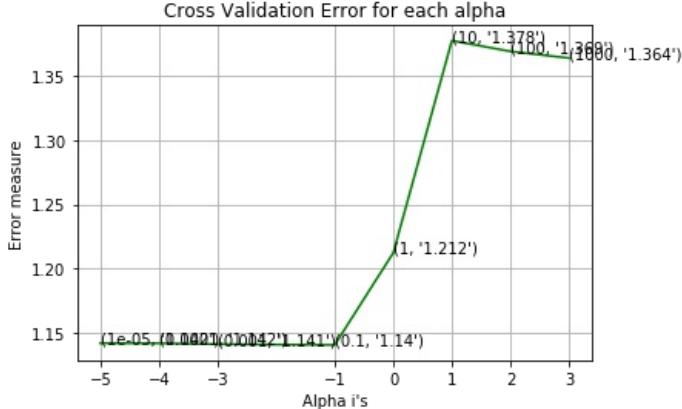
best_alpha = np.argmin(cv_log_error_array)
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_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-05
Log Loss : 1.141642690993172
for alpha = 0.0001
Log Loss : 1.1415619035369442
for alpha = 0.001
Log Loss : 1.1409024463526787
for alpha = 0.1
Log Loss : 1.1402553990302782
for alpha = 1
Log Loss : 1.2121645991707204
for alpha = 10
Log Loss : 1.3778072419402283
for alpha = 100
Log Loss : 1.3694236406377283
for alpha = 1000
Log Loss : 1.3642102656385848

```



```

For values of best alpha = 0.1 The train log loss is: 0.7826513048305865
For values of best alpha = 0.1 The cross validation log loss is: 1.1402553990302782
For values of best alpha = 0.1 The test log loss is: 1.1567710403660776

```

## Testing the model with best hyper parameters

In [52]:

```

find more about Multinomial Naive base function here http://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html

default paramters
sklearn.naive_bayes.MultinomialNB(alpha=1.0, fit_prior=True, class_prior=None)

some of methods of MultinomialNB()
fit(X, y[, sample_weight]) Fit Naive Bayes classifier according to X, y
predict(X) Perform classification on an array of test vectors X.
predict_log_proba(X) Return log-probability estimates for the test vector X.

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification

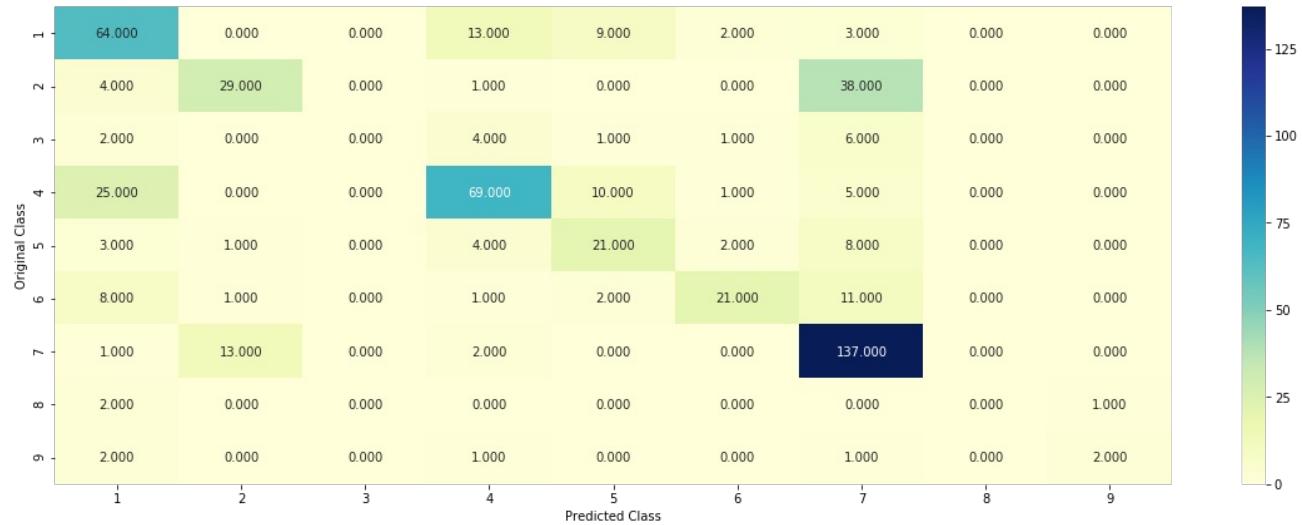
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
to avoid rounding error while multiplying probabilites we use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

```

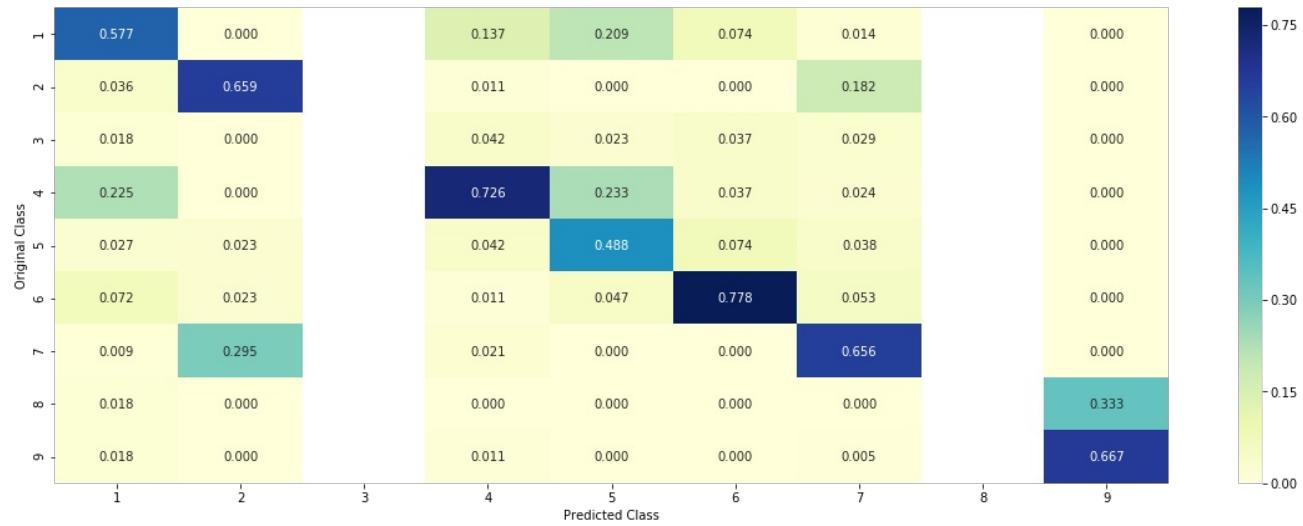
Log Loss : 1.1402553990302782

Number of missclassified point : 0.35526315789473684

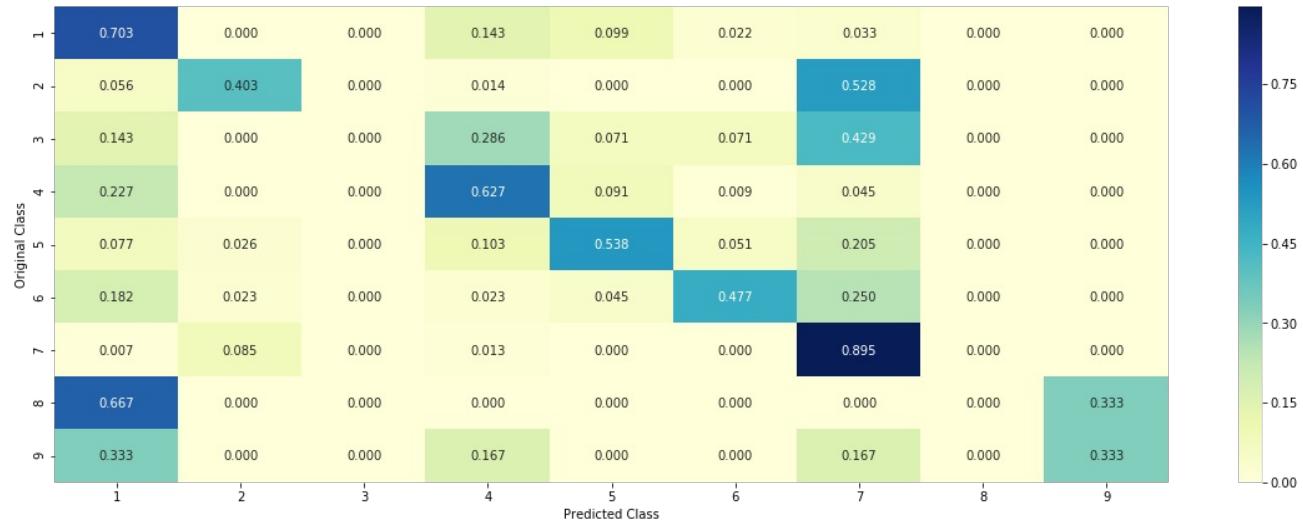
----- Confusion matrix -----



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



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



**Feature Importance, Correctly classified point**

In [53]:

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

```
Predicted Class : 5
Predicted Class Probabilities: [[0.1759 0.066 0.0232 0.1034 0.4483 0.1013 0.0704 0.0067 0.0049]]
Actual Class : 4

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

### Feature Importance, Incorrectly classified point

In [54]:

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.083 0.0532 0.0172 0.6256 0.0404 0.0371 0.1349 0.0049 0.0036]]
Actual Class : 4

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

## K Nearest Neighbour Classification

### Hyper parameter tuning

In [55]:

```
find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html

default parameter
KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

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

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

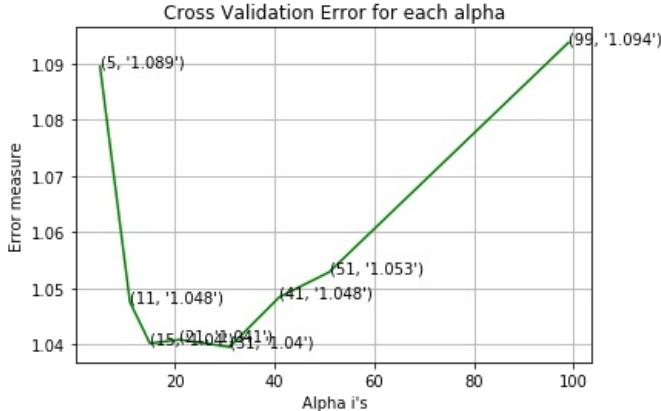
best_alpha = np.argmin(cv_log_error_array)
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

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

```

for alpha = 5
Log Loss : 1.0894545183775952
for alpha = 11
Log Loss : 1.0475476465019273
for alpha = 15
Log Loss : 1.0402111633324478
for alpha = 21
Log Loss : 1.0408659710357733
for alpha = 31
Log Loss : 1.0395637556224726
for alpha = 41
Log Loss : 1.0484483283997548
for alpha = 51
Log Loss : 1.0528924322582458
for alpha = 99
Log Loss : 1.0937073496966288

```



```

For values of best alpha = 31 The train log loss is: 0.8291299381972101
For values of best alpha = 31 The cross validation log loss is: 1.0395637556224726
For values of best alpha = 31 The test log loss is: 1.0422285191558496

```

### Testing the model with best hyper parameters

In [56]:

```

find more about KNeighborsClassifier() here http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html

default parameter
KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2,
metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

methods of
fit(X, y) : Fit the model using X as training data and y as target values
predict(X):Predict the class labels for the provided data
predict_proba(X):Return probability estimates for the test data X.

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

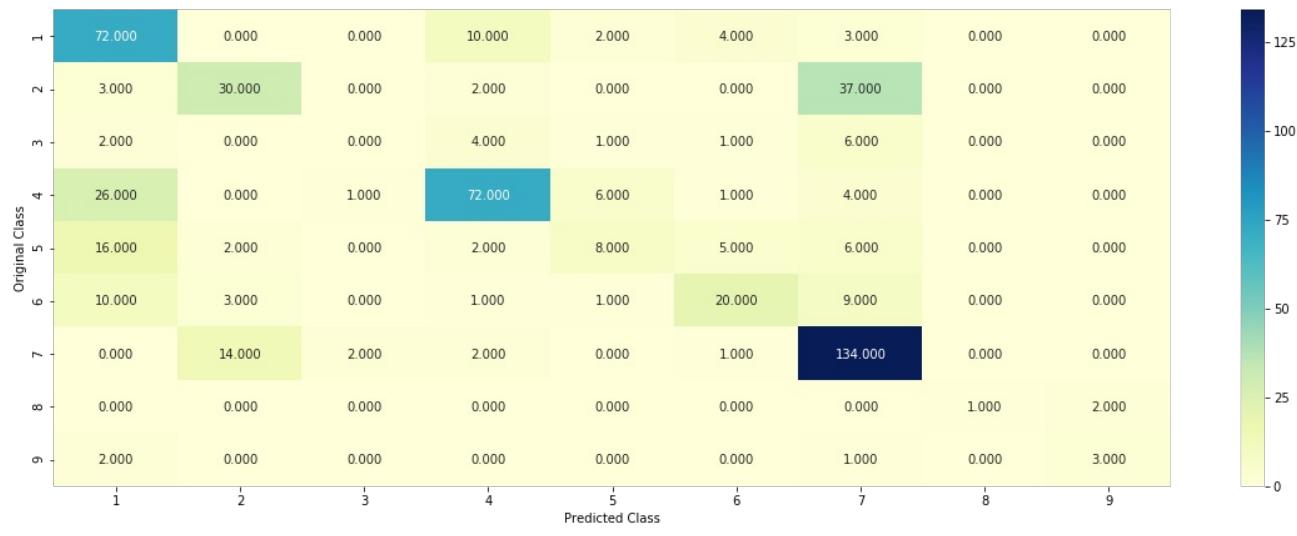
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseCoding, cv_y, clf)

```

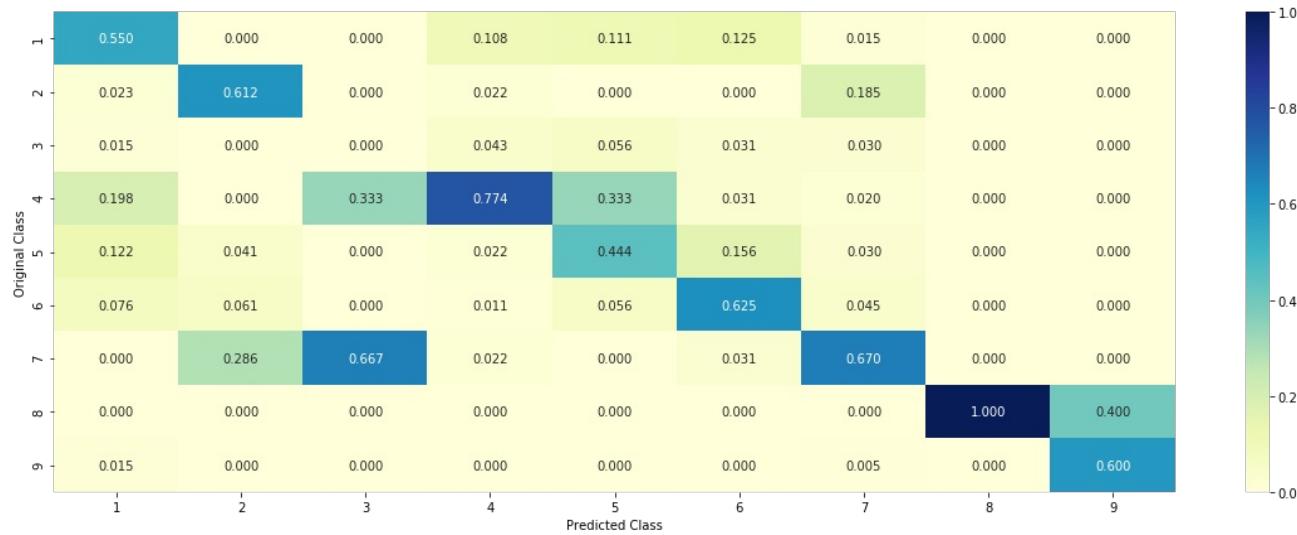
```

Log loss : 1.0395637556224726
Number of mis-classified points : 0.3609022556390977
----- Confusion matrix -----

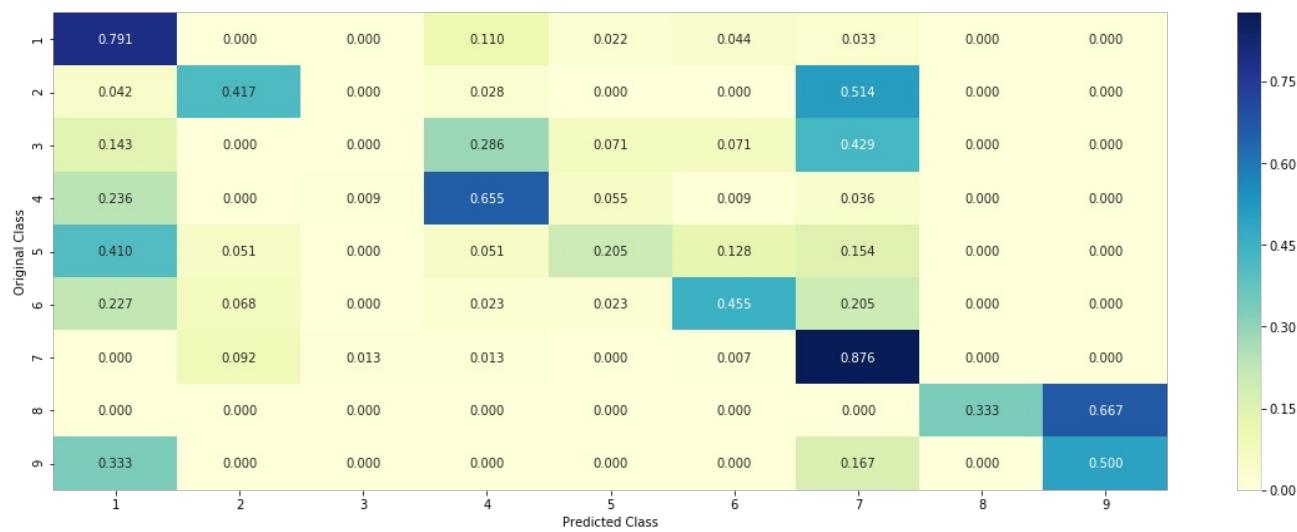
```



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



Sample Query point -1

In [57]:

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

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

Predicted Class : 4
Actual Class : 4
The 31 nearest neighbours of the test points belongs to classes [1 1 6 1 5 4 5 4 1 5 5 1 5 6 6 6 1
6 5 1 4 1 1 4 6 6 6 6 4 5]
Frequency of nearest points : Counter({6: 10, 1: 9, 5: 7, 4: 5})
```

## Sample Query Point-2

In [58]:

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

test_point_index = 100
predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
print("Predicted Class :", predicted_cls[0])
print("Actual Class :", test_y[test_point_index])
neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha])
print("the k value for knn is",alpha[best_alpha],"and the nearest neighbours of the test points belongs to classe
s",train_y[neighbors[1][0]])
print("Frequency of nearest points :",Counter(train_y[neighbors[1][0]]))

Predicted Class : 4
Actual Class : 4
the k value for knn is 31 and the nearest neighbours of the test points belongs to classes [4 7 8 4
9 4 4 7 4 4 4 4 4 4 7 7 4 4 1 7 4 4 4 7 1 1 7 2 8 1]
Frequency of nearest points : Counter({4: 16, 7: 7, 1: 4, 8: 2, 9: 1, 2: 1})
```

## Logistic Regression

### With Class balancing

#### Hyper paramter tuning

In [59]:

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

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

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

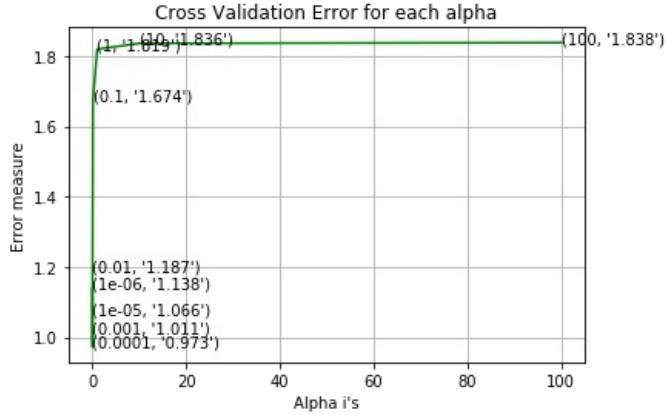
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

```

for alpha = 1e-06
Log Loss : 1.1378973204776976
for alpha = 1e-05
Log Loss : 1.0655131027682159
for alpha = 0.0001
Log Loss : 0.9726017451766378
for alpha = 0.001
Log Loss : 1.0111245413285104
for alpha = 0.01
Log Loss : 1.187408557806137
for alpha = 0.1
Log Loss : 1.6736229259902087
for alpha = 1
Log Loss : 1.8191100104797913
for alpha = 10
Log Loss : 1.8357369900966047
for alpha = 100
Log Loss : 1.8377569585929927

```



For values of best alpha = 0.0001 The train log loss is: 0.5948252788294383  
 For values of best alpha = 0.0001 The cross validation log loss is: 0.9726017451766378  
 For values of best alpha = 0.0001 The test log loss is: 0.9795367200026402

### Testing the model with best hyper parameters

In [60]:

```

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

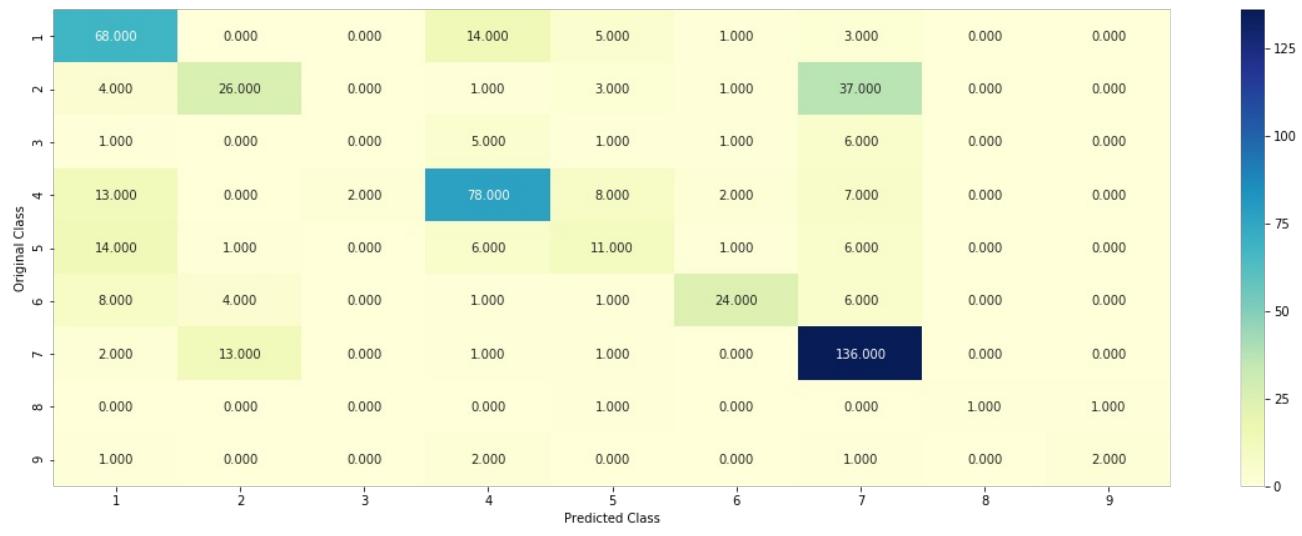
default parameters
SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
class_weight=None, warm_start=False, average=False,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

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

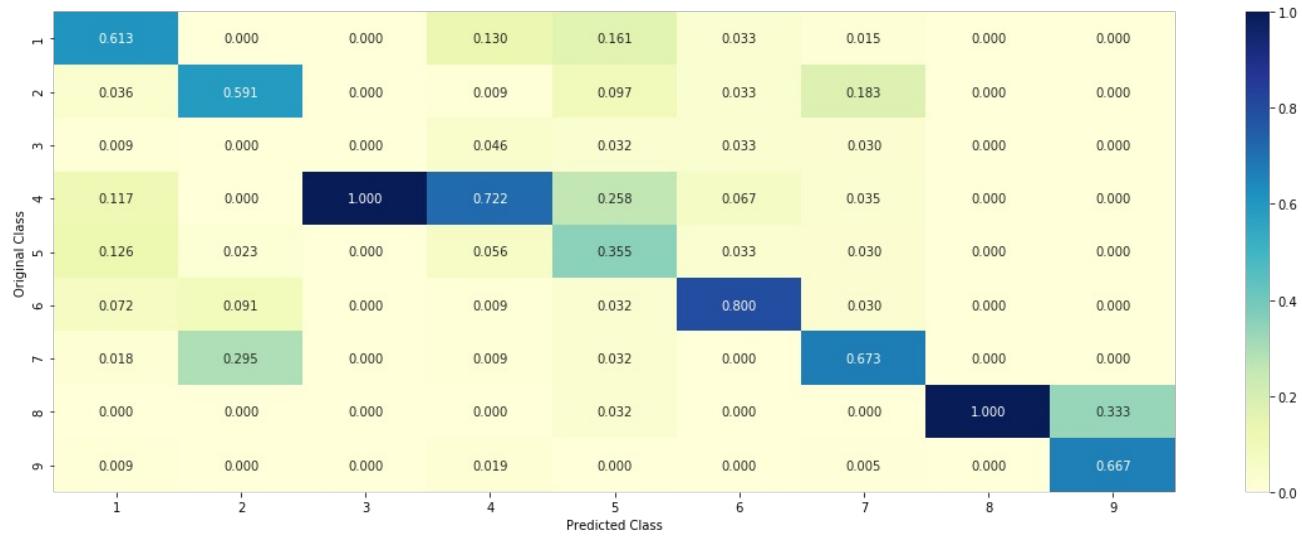
#-----
video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding,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

```

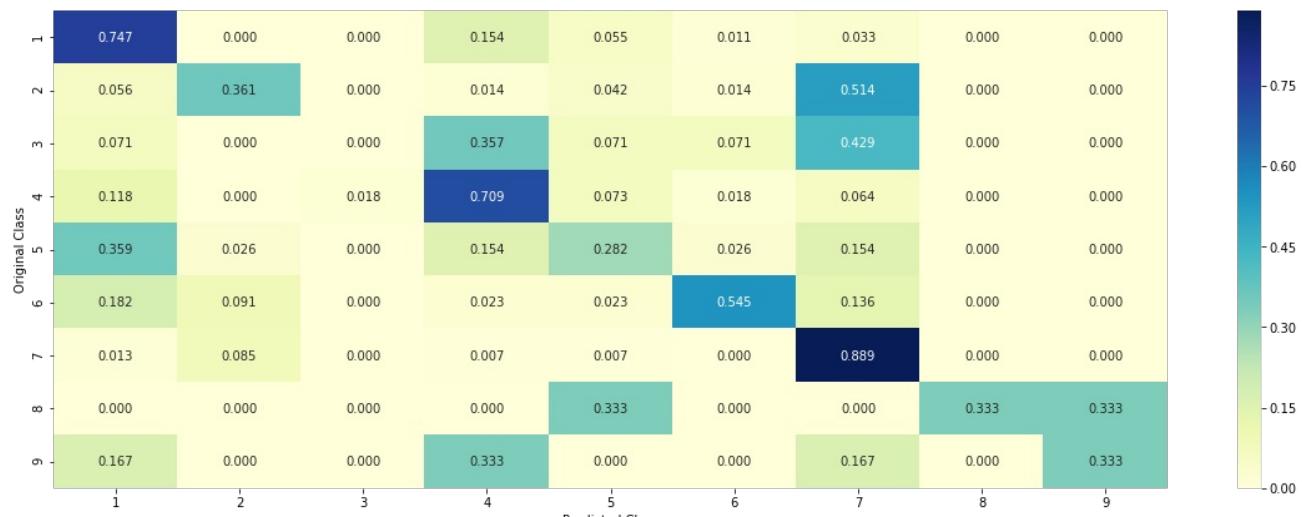
Log loss : 0.9726017451766378  
 Number of mis-classified points : 0.34962406015037595  
 ----- Confusion matrix -----



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



## Feature Importance

In [61]:

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

## Correctly Classified point

In [62]:

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

Predicted Class : 5
Predicted Class Probabilities: [[0.2959 0.0187 0.0264 0.0613 0.4014 0.1886 0.0026 0.0036 0.0016]]
Actual Class : 4

399 Text feature [000249] present in test data point [True]
Out of the top 500 features 1 are present in query point
```

## Incorrectly Classified point

In [63]:

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

Predicted Class : 4
Predicted Class Probabilities: [[0.0544 0.0246 0.0071 0.4849 0.0205 0.0043 0.3967 0.0046 0.0028]]
Actual Class : 4

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

## Without Class balancing

### Hyper parameter tuning

In [64]:

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification

video link:

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

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

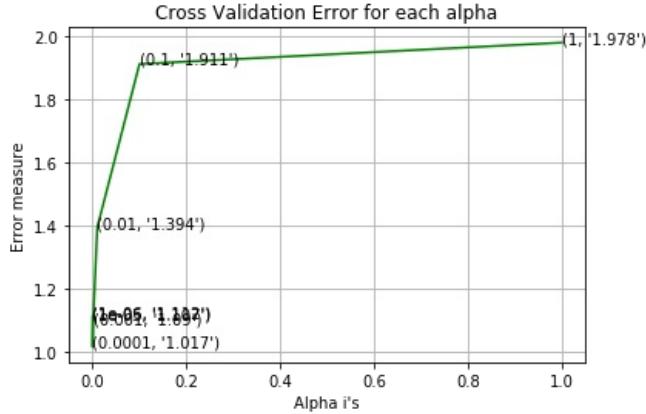
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_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 : 1.1119108157677706
for alpha = 1e-05
Log Loss : 1.1074048232362774
for alpha = 0.0001
Log Loss : 1.0168095145556055
for alpha = 0.001
Log Loss : 1.090005458579951
for alpha = 0.01
Log Loss : 1.3944084842291729
for alpha = 0.1
Log Loss : 1.9106820707698657
for alpha = 1
Log Loss : 1.9781664834971244

```



```

For values of best alpha = 0.0001 The train log loss is: 0.5808504775440413
For values of best alpha = 0.0001 The cross validation log loss is: 1.0168095145556055
For values of best alpha = 0.0001 The test log loss is: 1.0097076505457607

```

### Testing model with best hyper parameters

In [65]:

```

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

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

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

video link:

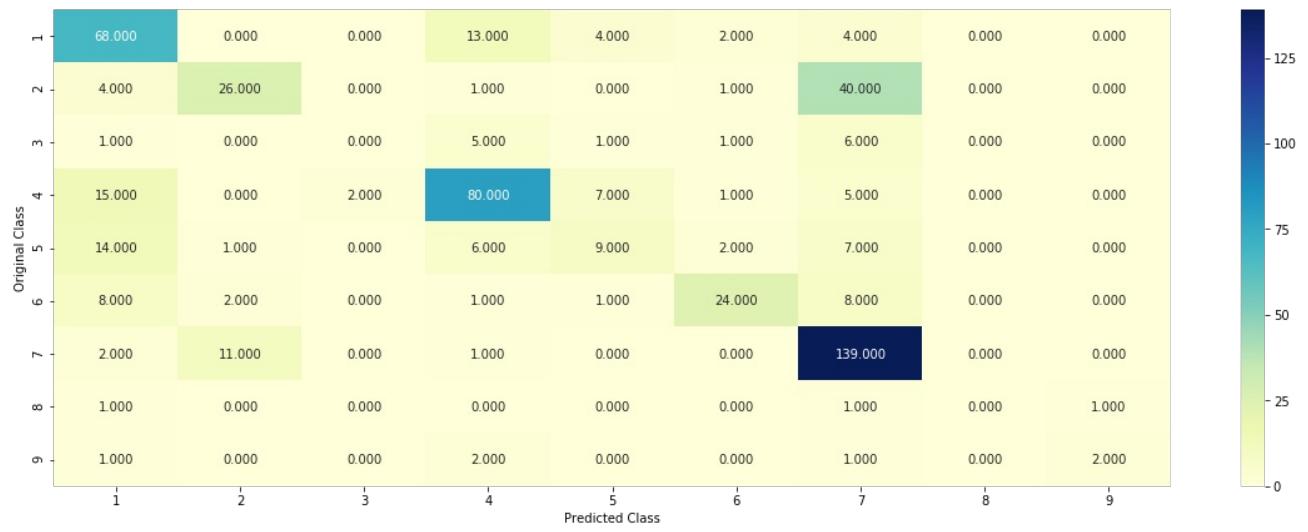
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

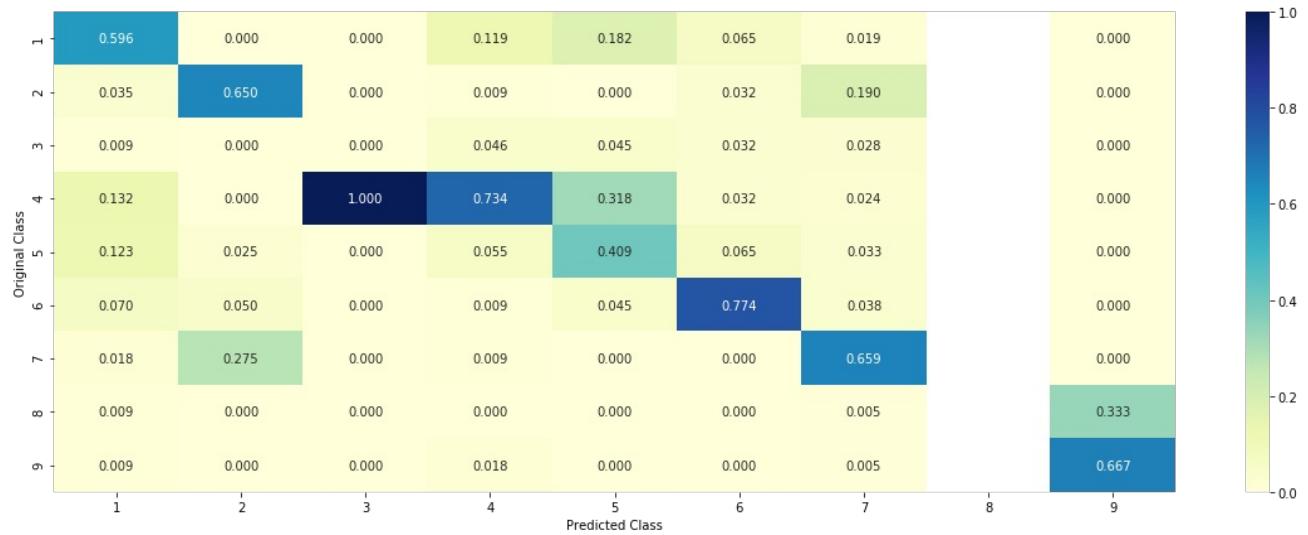
```

Log loss : 1.0168095145556055
Number of mis-classified points : 0.3458646616541353
----- Confusion matrix -----

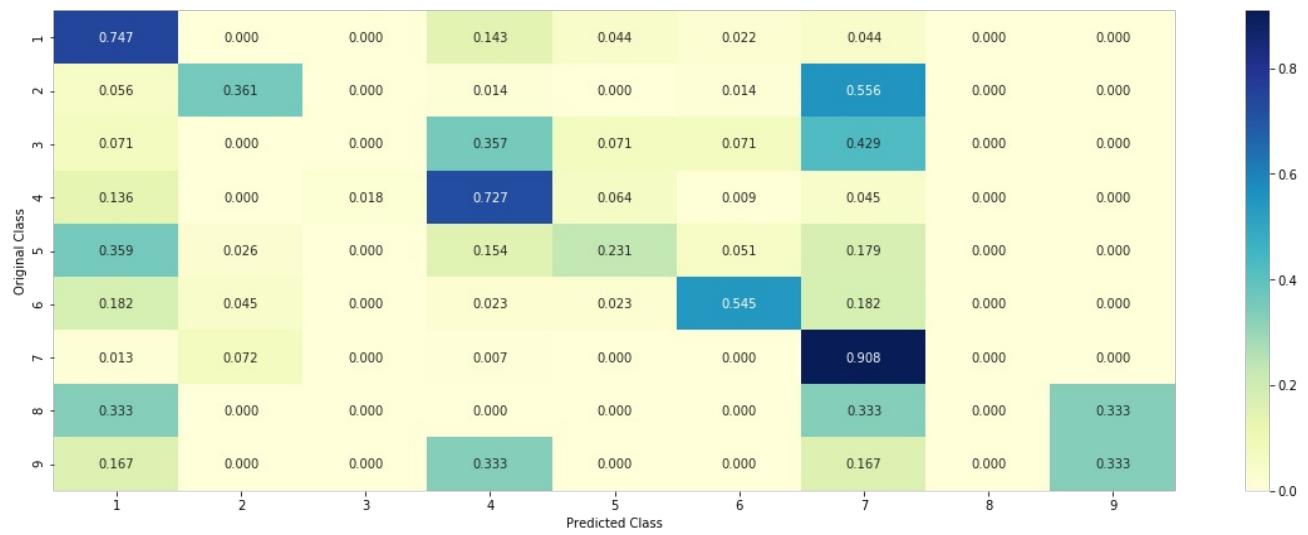
```



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



### Feature Importance, Correctly Classified point

In [66]:

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

```
Predicted Class : 5
Predicted Class Probabilities: [[3.408e-01 2.280e-02 4.100e-03 7.290e-02 3.833e-01 1.695e-01 3.700e-03
2.500e-03 3.000e-04]]
Actual Class : 4

380 Text feature [000249] present in test data point [True]
Out of the top 500 features 1 are present in query point
```

## Feature Importance, Inorrectly Classified point

In [67]:

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

Predicted Class : 4
Predicted Class Probabilities: [[0.0602 0.025 0.004 0.4868 0.0202 0.0047 0.3939 0.0039 0.0013]]
Actual Class : 4

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

## Linear Support Vector Machines

### Hyper paramter tuning

In [68]:

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibra
tion.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

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

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

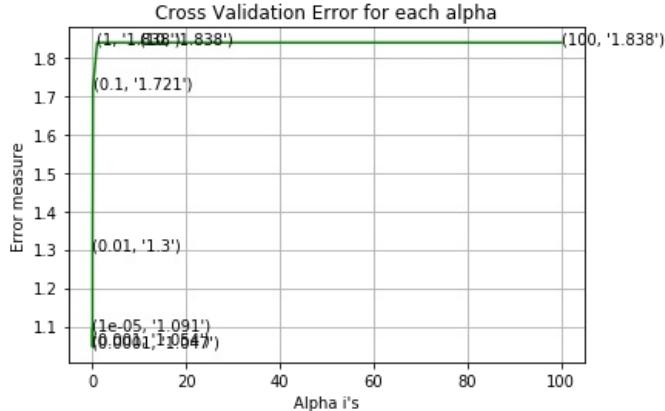
best_alpha = np.argmin(cv_log_error_array)
clf = SVC(C=i,kernel='linear',probability=True, class_weight='balanced')
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf,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

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

```

for C = 1e-05
Log Loss : 1.0909108656504967
for C = 0.0001
Log Loss : 1.0465506122146153
for C = 0.001
Log Loss : 1.053979076289716
for C = 0.01
Log Loss : 1.3001474962423754
for C = 0.1
Log Loss : 1.7214636161146284
for C = 1
Log Loss : 1.8383450325310866
for C = 10
Log Loss : 1.8383455017561428
for C = 100
Log Loss : 1.8383451976112968

```



```

For values of best alpha = 0.0001 The train log loss is: 0.6767518845684507
For values of best alpha = 0.0001 The cross validation log loss is: 1.0465506122146153
For values of best alpha = 0.0001 The test log loss is: 1.0785021787123679

```

## Testing model with best hyper parameters

In [69]:

```

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

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

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

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

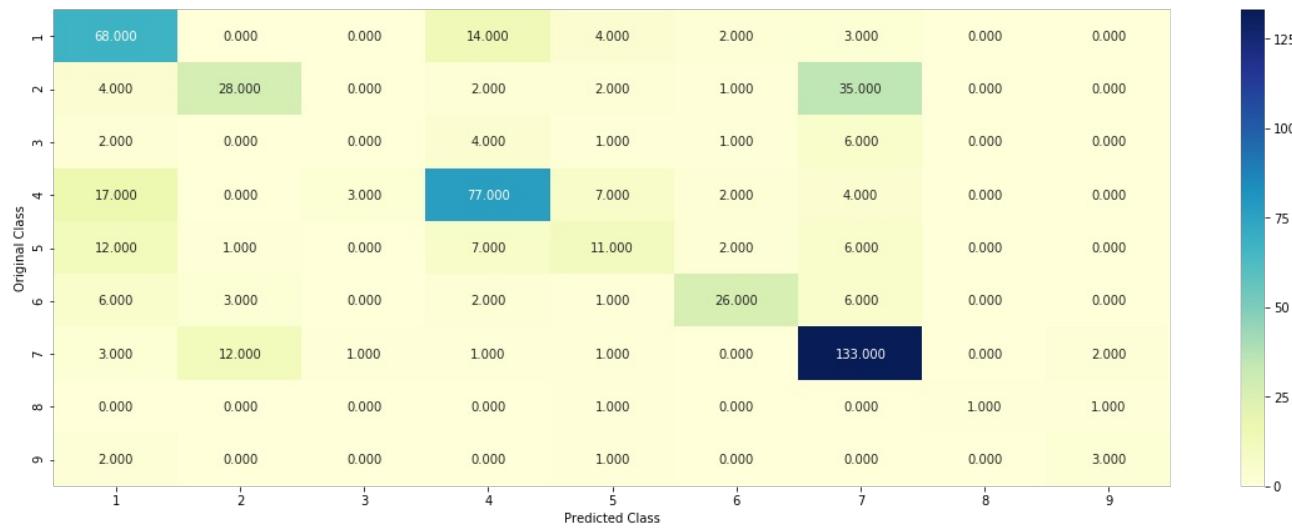
clf = SVC(C=alpha[best_alpha],kernel='linear',probability=True, class_weight='balanced')
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_state=42, class_weight='balanced')
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

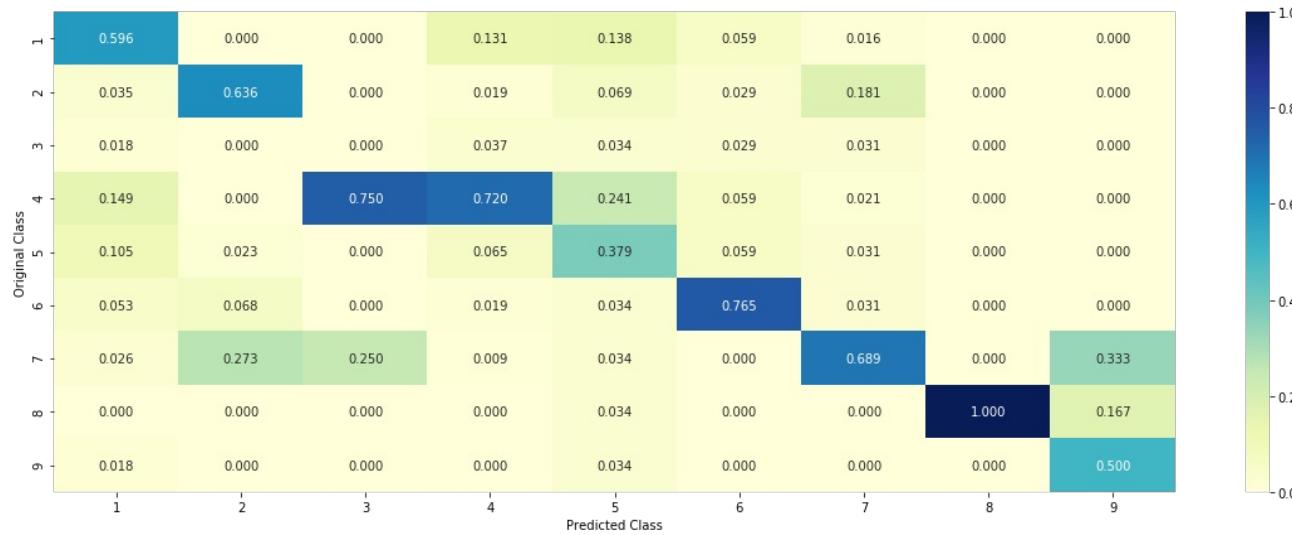
```

Log loss : 1.0465506122146153
Number of mis-classified points : 0.34774436090225563
----- Confusion matrix -----

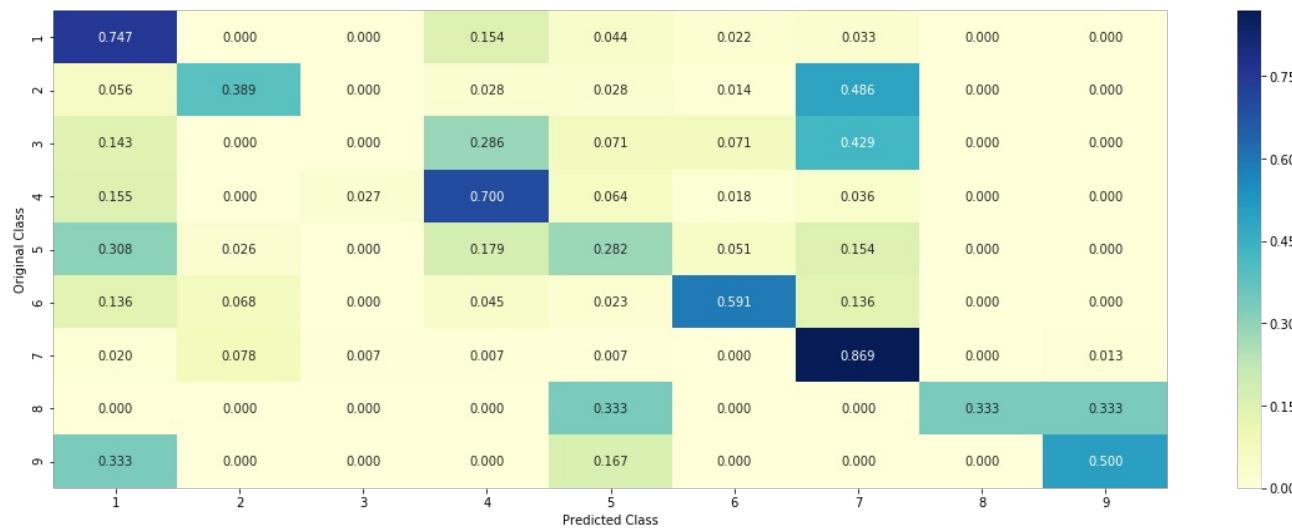
```



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



## Feature Importance

**For Correctly classified point**

In [70]:

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

```
Predicted Class : 6
Predicted Class Probabilities: [[0.1381 0.0496 0.0527 0.1055 0.2828 0.3235 0.0423 0.0035 0.0018]]
Actual Class : 4

73 Text feature [000] present in test data point [True]
Out of the top 500 features 1 are present in query point
```

#### For Incorrectly classified point

In [71]:

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

```
Predicted Class : 4
Predicted Class Probabilities: [[0.0604 0.0282 0.0108 0.482 0.0313 0.0046 0.3771 0.0035 0.0021]]
Actual Class : 4

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

## 4.5 Random Forest Classifier

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

In [72]:

```

default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrea
se=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=Fa
lse,
class_weight=None)

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibra
tion.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

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

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

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

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

```

for n_estimators = 100 and max depth = 5
Log Loss : 1.239948403365552
for n_estimators = 100 and max depth = 10
Log Loss : 1.2407343670249138
for n_estimators = 200 and max depth = 5
Log Loss : 1.224090314573389
for n_estimators = 200 and max depth = 10
Log Loss : 1.2306299279008606
for n_estimators = 500 and max depth = 5
Log Loss : 1.217698068576323
for n_estimators = 500 and max depth = 10
Log Loss : 1.2243302024967533
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2117571932451994
for n_estimators = 1000 and max depth = 10
Log Loss : 1.2221676303855031
for n_estimators = 2000 and max depth = 5
Log Loss : 1.2082978093221546
for n_estimators = 2000 and max depth = 10
Log Loss : 1.2204876833464997
For values of best estimator = 2000 The train log loss is: 0.8566719980344872
For values of best estimator = 2000 The cross validation log loss is: 1.2082978093221546
For values of best estimator = 2000 The test log loss is: 1.175293325156972

```

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

In [73]:

```

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

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

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

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

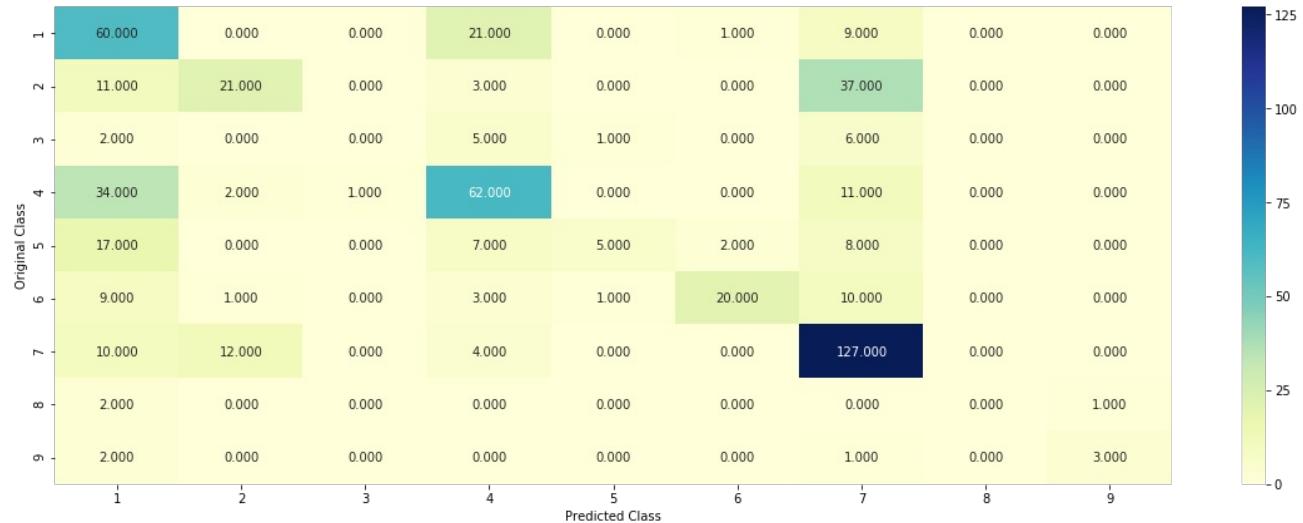
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha%2)], random_state=42, n_jobs=-1)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

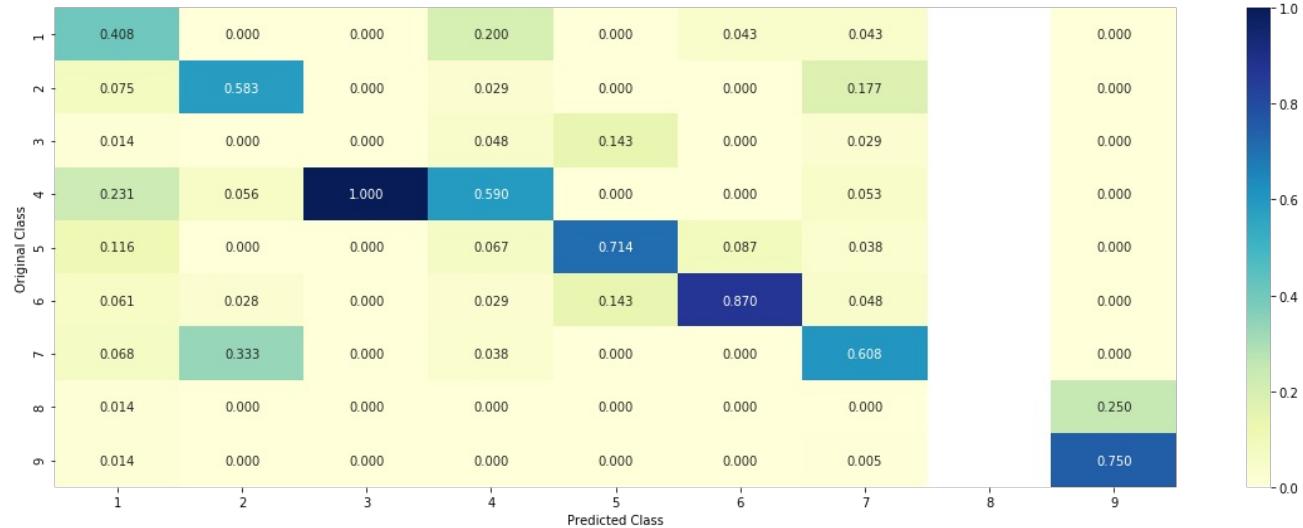
Log loss : 1.2082978093221546

Number of mis-classified points : 0.4398496240601504

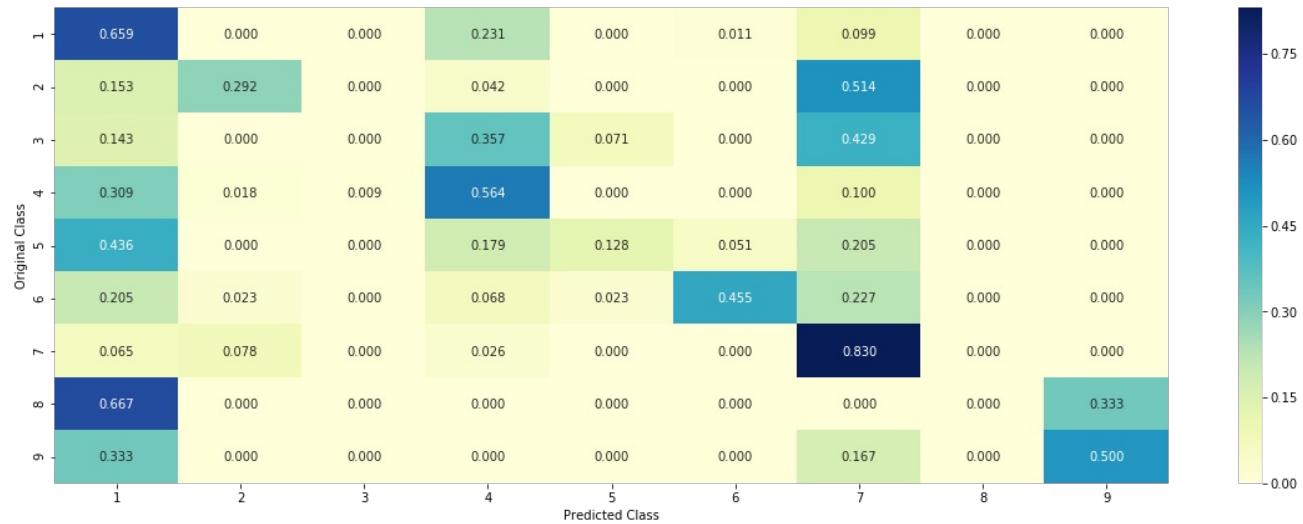
----- Confusion matrix -----



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



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



## Feature Importance

### Correctly Classified point

In [74]:

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

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

```
Predicted Class : 1
Predicted Class Probabilities: [[0.4613 0.0147 0.0136 0.2234 0.1786 0.0781 0.0196 0.0046 0.0062]]
Actual Class : 4
```

```

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

### Incorrectly Classified point

In [75]:

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

Predicted Class : 4
Predicted Class Probabilities: [[0.1959 0.0808 0.019 0.3506 0.0675 0.0556 0.2012 0.0171 0.0124]]
Actual Class : 4
```

```

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

### Hyper parameter tuning (With Response Coding)

In [76]:

```

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)

```

```

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

alpha = [10,50,100,200,500,1000]
max_depth = [2,3,5,10]
cv_log_error_array = []
for i in alpha:
 for j in max_depth:
 print("for n_estimators =", i,"and max depth = ", j)
 clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, random_state=42, n_jobs=-1)
 clf.fit(train_x_responseCoding, train_y)
 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
 sig_clf.fit(train_x_responseCoding, train_y)
 sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
 cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
 print("Log Loss :",log_loss(cv_y, sig_clf_probs))
...
fig, ax = plt.subplots()
features = np.dot(np.array(alpha)[:,None],np.array(max_depth)[None]).ravel()
ax.plot(features, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
 ax.annotate((alpha[int(i/4)],max_depth[int(i%4)],str(txt)), (features[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()
'''

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

predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

for n_estimators = 10 and max depth = 2
Log Loss : 1.9815573699007627
for n_estimators = 10 and max depth = 3
Log Loss : 1.7156198284848752
for n_estimators = 10 and max depth = 5
Log Loss : 1.2091557154072532
for n_estimators = 10 and max depth = 10
Log Loss : 1.6763972415537254
for n_estimators = 50 and max depth = 2
Log Loss : 1.5830686439626
for n_estimators = 50 and max depth = 3
Log Loss : 1.4110982985102154
for n_estimators = 50 and max depth = 5
Log Loss : 1.2863239652019216
for n_estimators = 50 and max depth = 10
Log Loss : 1.5736687336529367
for n_estimators = 100 and max depth = 2
Log Loss : 1.4807841266514132
for n_estimators = 100 and max depth = 3
Log Loss : 1.4382288860027586
for n_estimators = 100 and max depth = 5
Log Loss : 1.2741991406012791
for n_estimators = 100 and max depth = 10
Log Loss : 1.6209890699355554
for n_estimators = 200 and max depth = 2
Log Loss : 1.506989890649295
for n_estimators = 200 and max depth = 3
Log Loss : 1.4124207499482735
for n_estimators = 200 and max depth = 5
Log Loss : 1.3261881332897674
for n_estimators = 200 and max depth = 10
Log Loss : 1.5733895872483683
for n_estimators = 500 and max depth = 2
Log Loss : 1.542851292617801
for n_estimators = 500 and max depth = 3
Log Loss : 1.4674393613005776
for n_estimators = 500 and max depth = 5
Log Loss : 1.3957599093475692
for n_estimators = 500 and max depth = 10
Log Loss : 1.6367917954756621
for n_estimators = 1000 and max depth = 2
Log Loss : 1.5376343704480848
for n_estimators = 1000 and max depth = 3
Log Loss : 1.456198280947336
for n_estimators = 1000 and max depth = 5
Log Loss : 1.3870506065188752
for n_estimators = 1000 and max depth = 10
Log Loss : 1.6543880859012372
For values of best alpha = 10 The train log loss is: 0.07336466282971817
For values of best alpha = 10 The cross validation log loss is: 1.2091557154072534
For values of best alpha = 10 The test log loss is: 1.2563508184806496

```

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

In [77]:

```

default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrea
se=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=Fa
lse,
class_weight=None)

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

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

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

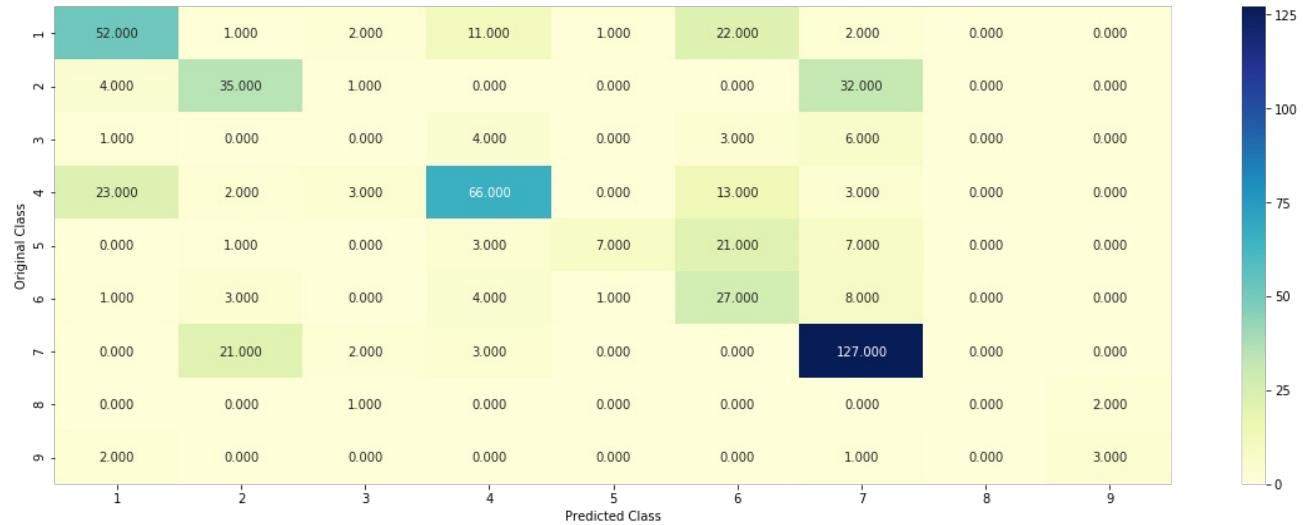
```

clf = RandomForestClassifier(max\_depth=max\_depth[int(best\_alpha%4)], n\_estimators=alpha[int(best\_alpha/4)], criter
ion='gini', max\_features='auto',random\_state=42)
predict\_and\_plot\_confusion\_matrix(train\_x\_responseCoding, train\_y, cv\_x\_responseCoding, cv\_y, clf)

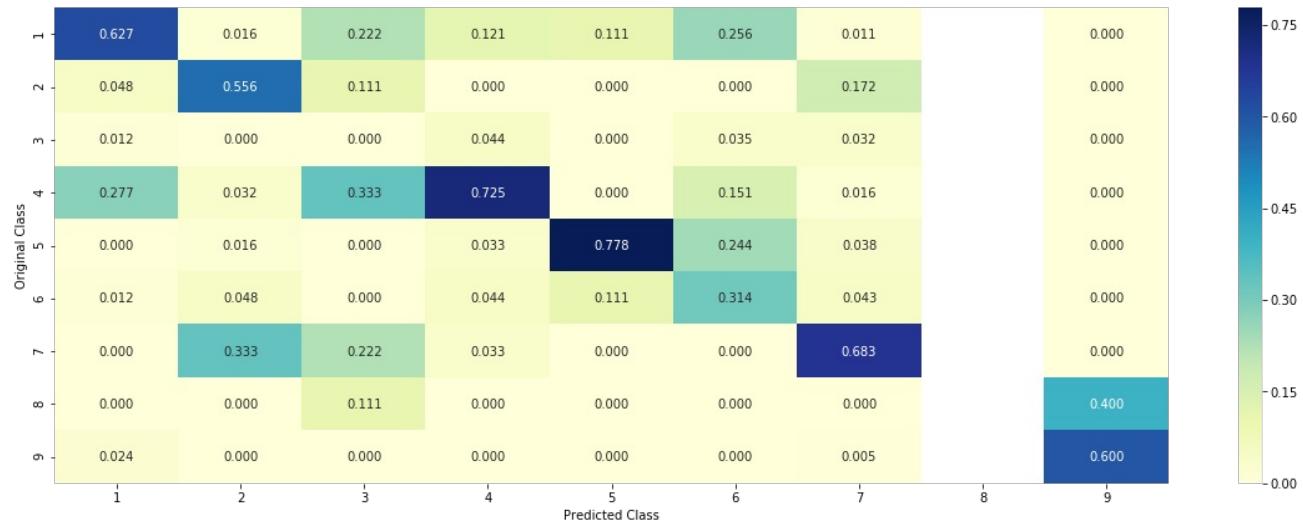
Log loss : 1.2091557154072532

Number of mis-classified points : 0.4041353383458647

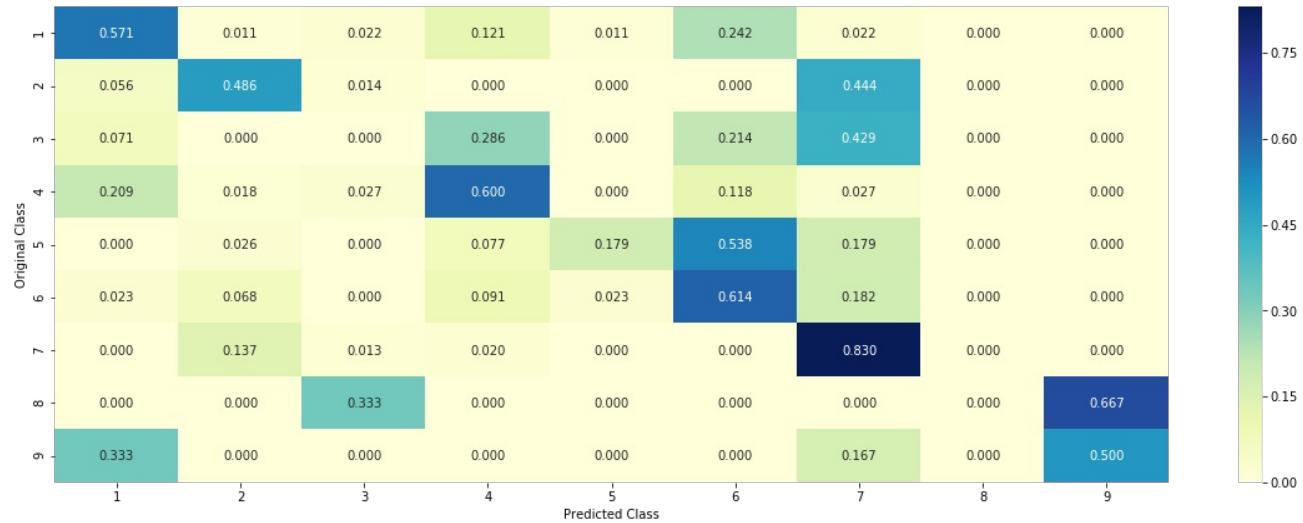
----- Confusion matrix -----



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



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



## Feature Importance

## Correctly Classified point

In [78]:

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

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

```
Predicted Class : 6
Predicted Class Probabilities: [[0.0822 0.0052 0.093 0.2627 0.1495 0.3923 0.0038 0.0054 0.0058]]
Actual Class : 4

```

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

**Incorrectly Classified point**

In [79]:

```
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].re
shape(1,-1)),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-" * 50)
for i in indices:
 if i < 9:
 print("Gene is important feature")
 elif i < 18:
 print("Variation is important feature")
 else:
 print("Text is important feature")
```

```
Predicted Class : 4
Predicted Class Probabilities: [[0.174 0.0281 0.1985 0.443 0.0187 0.0426 0.0265 0.0384 0.0301]]
Actual Class : 4
```

```

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

## Stack the models

### testing with hyper parameter tuning

In [38]:

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

default parameters
SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=N
one,
shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0
.5,
class_weight=None, warm_start=False, average=False, n_iter=None)

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

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

```

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

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

```

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

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

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

default parameters
sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrea
se=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=Fa
lse,
class_weight=None)

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

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

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

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

clf2 = SGDClassifier(alpha=1, penalty='l2', loss='hinge', class_weight='balanced', random_state=0)
clf2.fit(train_x_onehotCoding, train_y)
sig_clf2 = CalibratedClassifierCV(clf2, method="sigmoid")

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

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

```

Logistic Regression : Log Loss: 1.14  
Support vector machines : Log Loss: 1.83  
Naive Bayes : Log Loss: 1.25

Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.178  
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.040  
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.540  
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.205  
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.246  
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.374

**testing the model with the best hyper parameters**

In [39]:

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

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

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

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

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

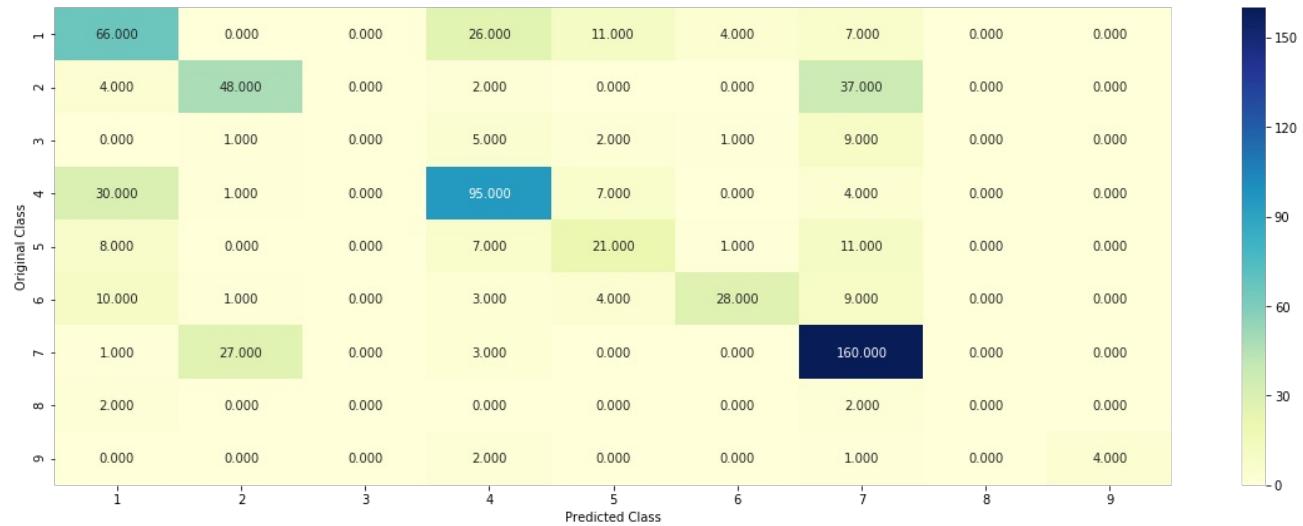
Log loss (train) on the stacking classifier : 0.8296709957190612

Log loss (CV) on the stacking classifier : 1.2048178545094657

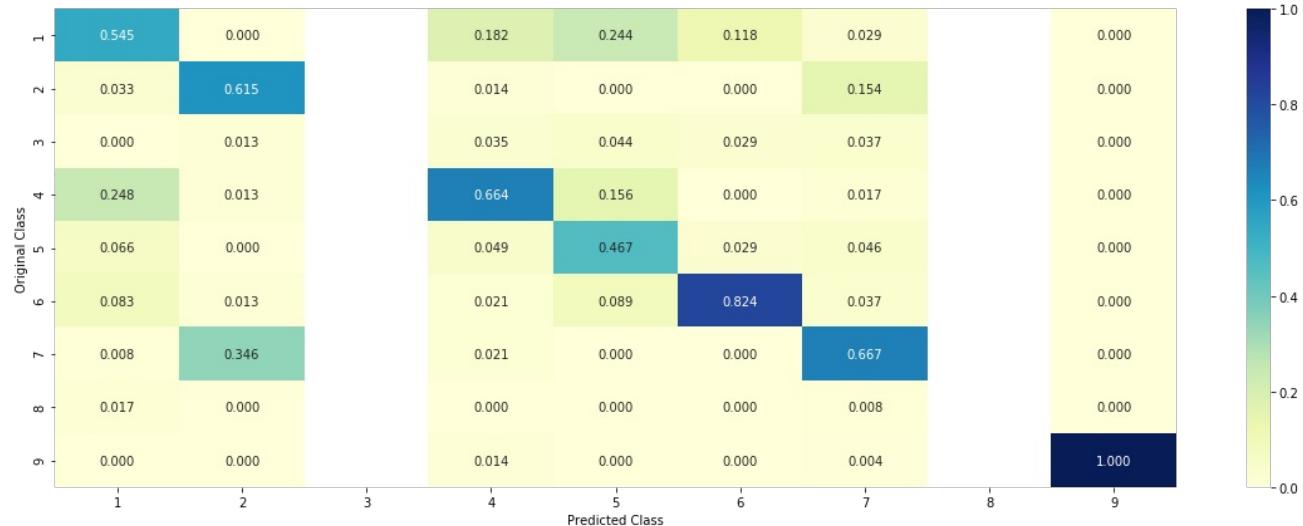
Log loss (test) on the stacking classifier : 1.1070106755024465

Number of missclassified point : 0.36541353383458647

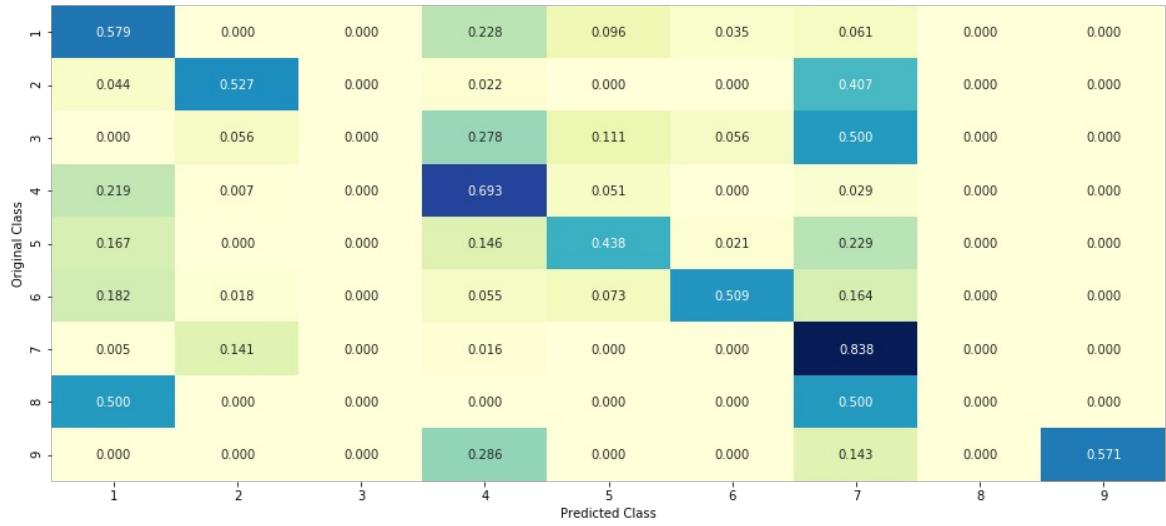
----- Confusion matrix -----



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



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



## Maximum Voting classifier

In [75]:

```
#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), ('rfi', sig_clf3)], voting='soft')
vclf.fit(train_x_onehotCoding, train_y)
print("Log loss (train) on the VotingClassifier : ", log_loss(train_y, vclf.predict_proba(train_x_onehotCoding)))
print("Log loss (CV) on the VotingClassifier : ", log_loss(cv_y, vclf.predict_proba(cv_x_onehotCoding)))
print("Log loss (test) on the VotingClassifier : ", log_loss(test_y, vclf.predict_proba(test_x_onehotCoding)))
print("Number of missclassified point : ", np.count_nonzero((vclf.predict(test_x_onehotCoding)- test_y)/test_y.shape[0]))
plot_confusion_matrix(test_y=test_y, predict_y=vclf.predict(test_x_onehotCoding))
```

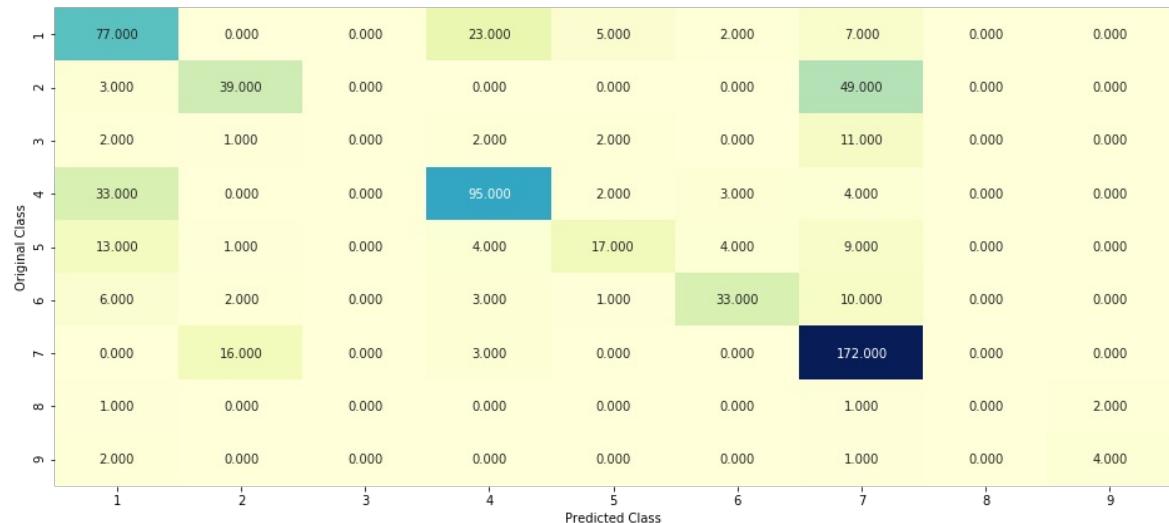
Log loss (train) on the VotingClassifier : 0.9695451507292676

Log loss (CV) on the VotingClassifier : 1.2034536765081423

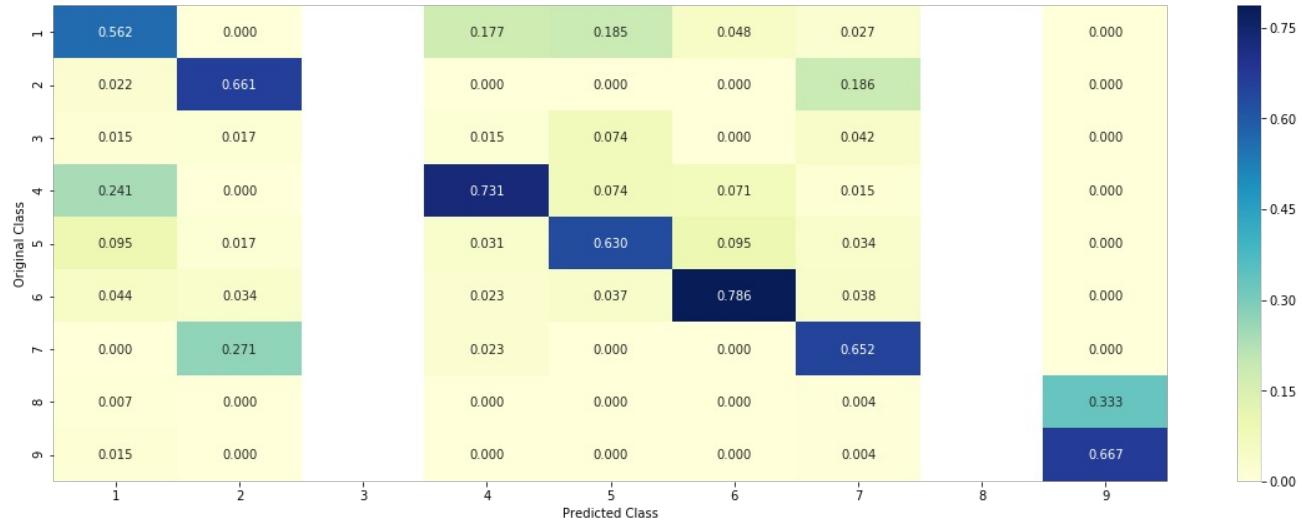
Log loss (test) on the VotingClassifier : 1.1656331333745162

Number of missclassified point : 0.34285714285714286

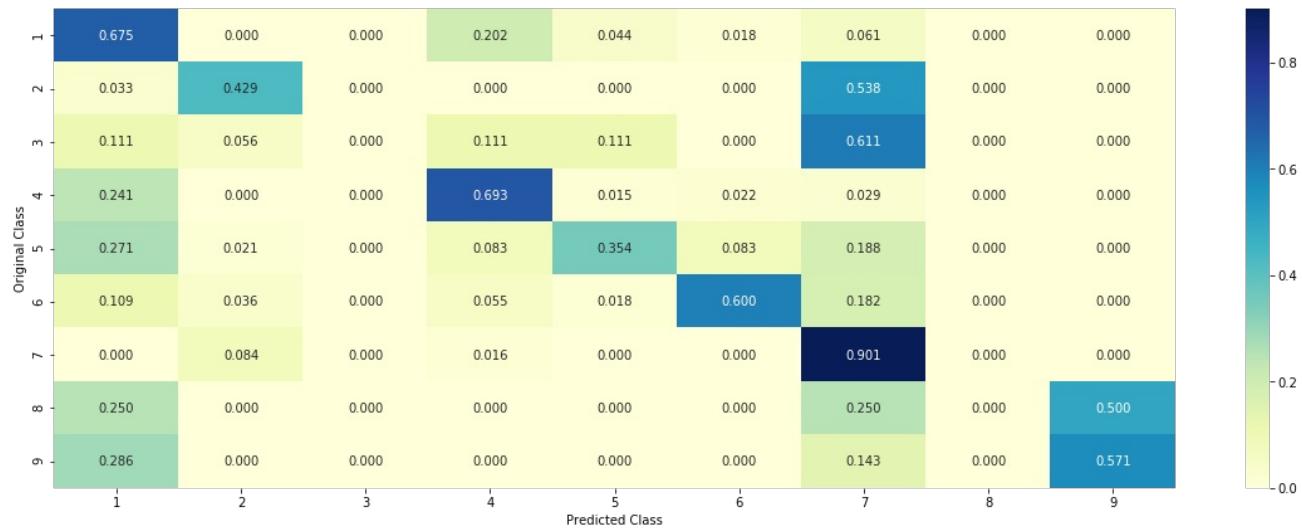
----- Confusion matrix -----



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



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



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

### gene feature

In [23]:

```
one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer(ngram_range=(1,2))
train_gene_feature_onehotCoding_count = gene_vectorizer.fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding_count = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_onehotCoding_count = gene_vectorizer.transform(cv_df['Gene'])
```

In [24]:

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

train\_gene\_feature\_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 229)

### variation feature

In [25]:

```
one-hot encoding of variation feature.
variation_vectorizer = CountVectorizer(ngram_range=(1,2))
train_variation_feature_onehotCoding_count = variation_vectorizer.fit_transform(train_df['Variation'])
test_variation_feature_onehotCoding_count = variation_vectorizer.transform(test_df['Variation'])
cv_variation_feature_onehotCoding_count = variation_vectorizer.transform(cv_df['Variation'])
```

In [26]:

```
print("train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The shape o
f Variation feature:", train_variation_feature_onehotCoding_count.shape)
```

train\_variation\_feature\_onehotEncoded is converted feature using the onne-hot encoding method. The s  
hape of Variation feature: (2124, 2070)

## Text feature

In [27]:

```
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 [28]:

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

In [29]:

```
building a CountVectorizer with all the words that occured minimum 3 times in train data
text_vectorizer = CountVectorizer(ngram_range=(1,2),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 features) vector
train_text_fea_counts = train_text_feature_onehotCoding.sum(axis=0).A1

zip(list(text_features),text_fea_counts) will zip a word with its number of times it occured
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 : 786314

In [30]:

```
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 [31]:

```
#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 [32]:

```
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 [33]:

```
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 [34]:

```
#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()))
```

## Machine Learning Model

In [36]:

```
#Data preparation for ML models.

#Misc. functionns for ML models

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

 # for calculating log_loss we will provide the array of probabilities belongs to each class
 print("Log loss :",log_loss(test_y, sig_clf.predict_proba(test_x)))
 # calculating the number of data points that are misclassified
 print("Number of mis-classified points :", np.count_nonzero((pred_y- test_y))/test_y.shape[0])
 plot_confusion_matrix(test_y, pred_y)
```

In [37]:

```
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 [38]:

```
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_imptfeature_names(indices, text, gene, var, no_features):
 gene_count_vec = CountVectorizer()
 var_count_vec = CountVectorizer()
 text_count_vec = CountVectorizer(min_df=3)

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

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

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

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

**Stacking the three types of features**

In [39]:

```
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_count,train_variation_feature_onehotCoding_count))
test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding_count,test_variation_feature_onehotCoding_count))
cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding_count,cv_variation_feature_onehotCoding_count))

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

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

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

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

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

## Logistic Regression

With Class balancing

*Hyper parameter tuning*

In [40]:

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification

video link:

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

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

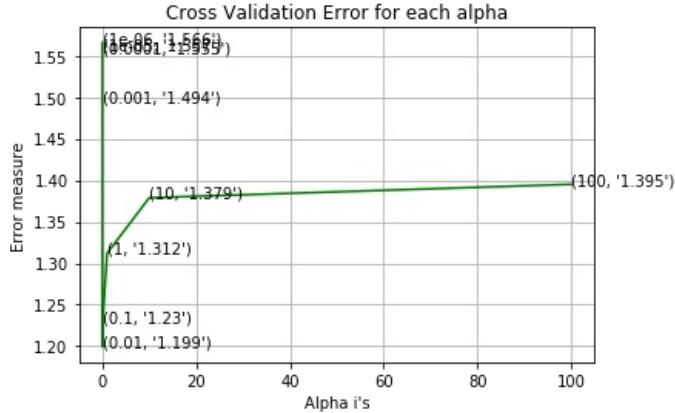
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

```

for alpha = 1e-06
Log Loss : 1.5657290200662106
for alpha = 1e-05
Log Loss : 1.556615895780145
for alpha = 0.0001
Log Loss : 1.5546079277453628
for alpha = 0.001
Log Loss : 1.4935793049267716
for alpha = 0.01
Log Loss : 1.1991340926644478
for alpha = 0.1
Log Loss : 1.2300696555953663
for alpha = 1
Log Loss : 1.3118199369721868
for alpha = 10
Log Loss : 1.3788215206669299
for alpha = 100
Log Loss : 1.39507506198723

```



For values of best alpha = 0.01 The train log loss is: 0.8441038279613815  
 For values of best alpha = 0.01 The cross validation log loss is: 1.1991340926644478  
 For values of best alpha = 0.01 The test log loss is: 1.3164304030144693

## Testing the model with best hyper parameters

In [41]:

```

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

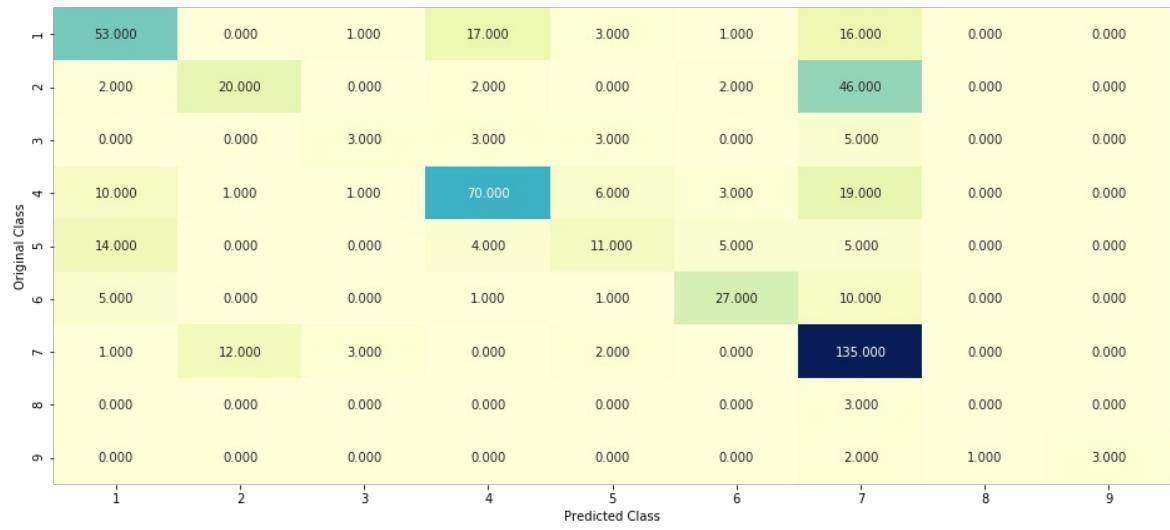
default parameters
SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
class_weight=None, warm_start=False, average=False,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

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

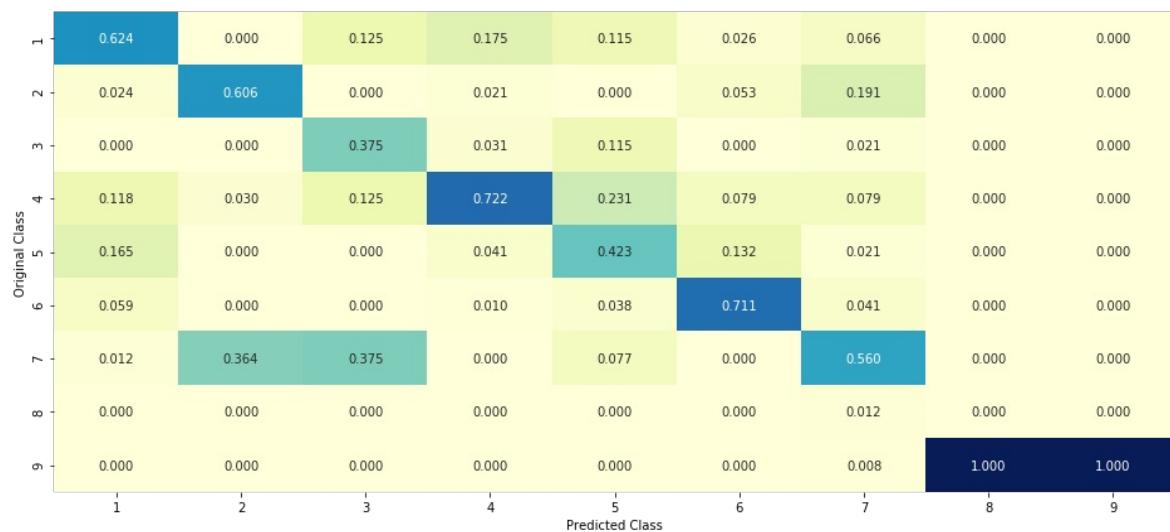
#-----
video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding,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

```

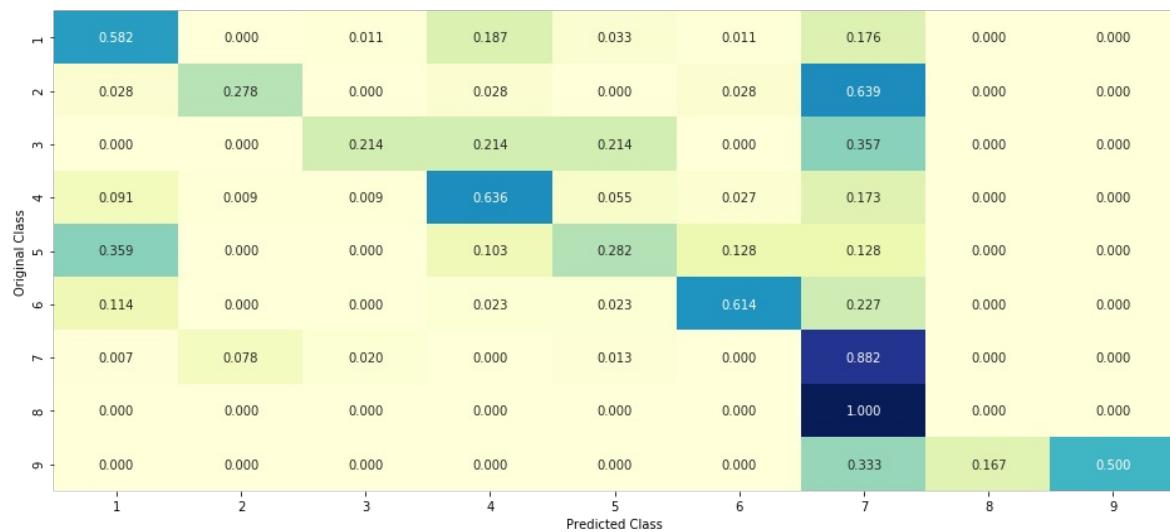
Log loss : 1.1991340926644478  
 Number of mis-classified points : 0.39473684210526316  
 ----- Confusion matrix -----



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



## Without Class balancing

### Hyper parameter tuning

In [47]:

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification
#-----
video link:
#-----

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

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

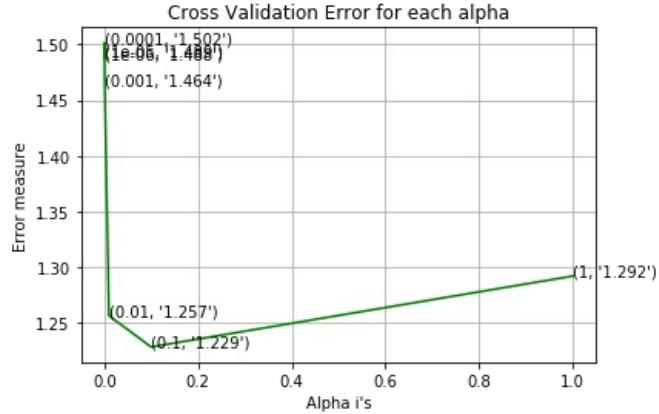
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_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 : 1.4879573852844477
for alpha = 1e-05
Log Loss : 1.4893591479930504
for alpha = 0.0001
Log Loss : 1.5016749825229934
for alpha = 0.001
Log Loss : 1.4638049962285935
for alpha = 0.01
Log Loss : 1.2566289553231964
for alpha = 0.1
Log Loss : 1.228575829419387
for alpha = 1
Log Loss : 1.29217813353908

```



For values of best alpha = 0.1 The train log loss is: 0.8267845017830673  
 For values of best alpha = 0.1 The cross validation log loss is: 1.228575829419387  
 For values of best alpha = 0.1 The test log loss is: 1.1291767706507838

### Testing model with best hyper parameters

In [47]:

```

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

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

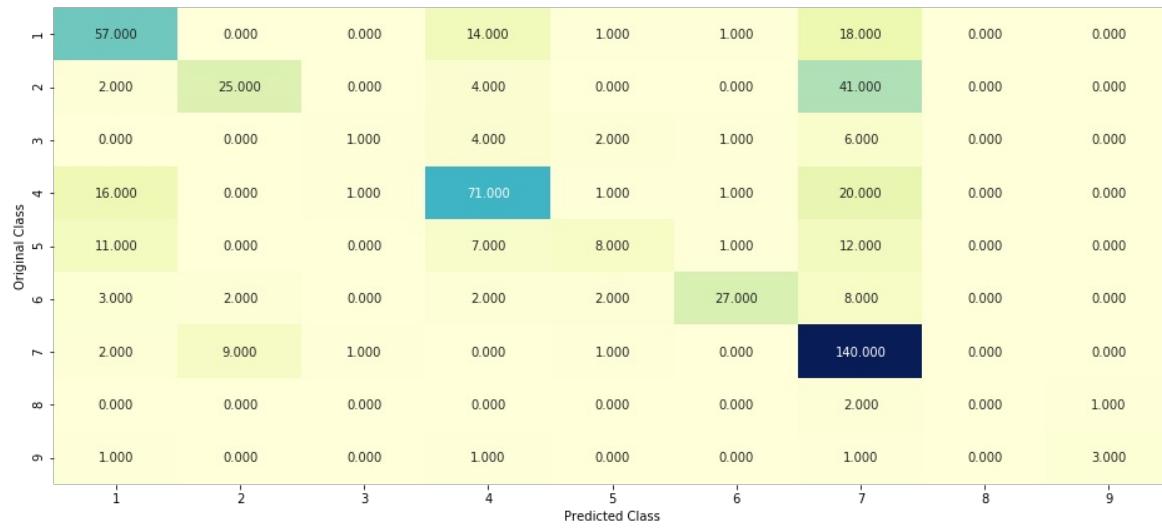
some of methods
fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
predict(X) Predict class labels for samples in X.

video link:

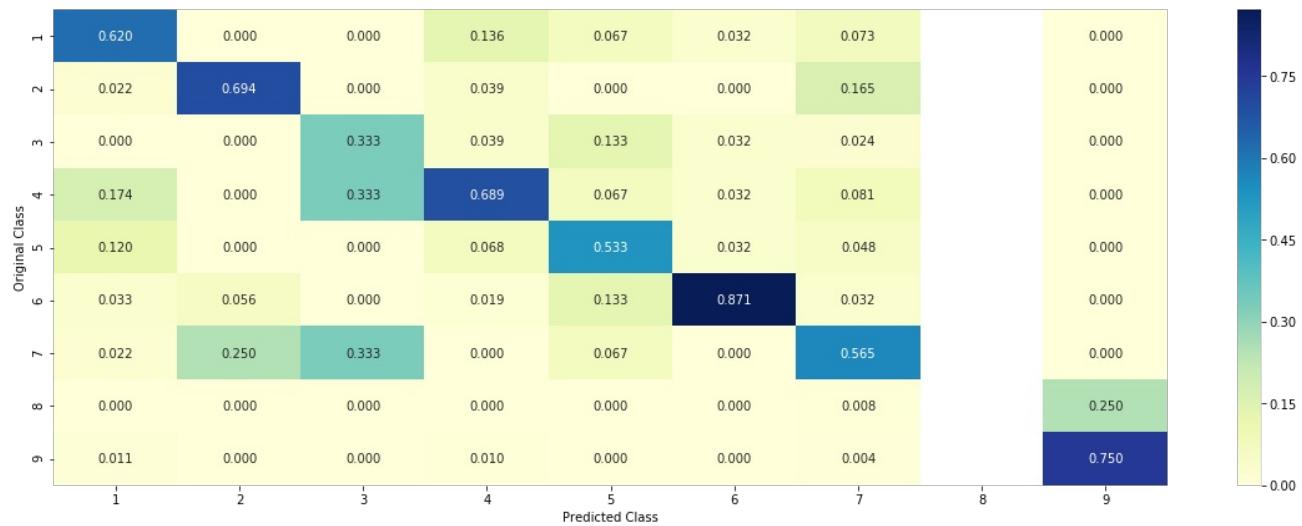
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)

```

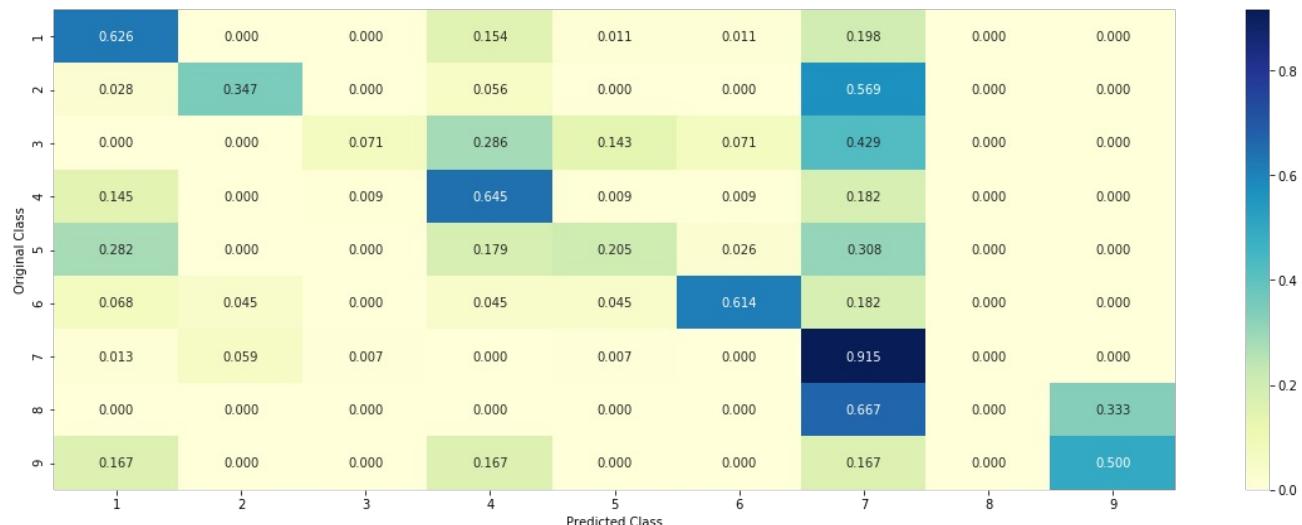
Log loss : 1.2272386750991784  
 Number of mis-classified points : 0.37593984962406013  
 ----- Confusion matrix -----



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



## Feature Engineering Technique

### Gene feature

In [14]:

```
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'])
```

In [15]:

```
gene_vectorizer.get_feature_names()
```

Out[15]:

```
['abl1',
 'acvrl1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ari',
 'araf',
 'arid1a',
 'arid1b',
 'arid2',
 'arid5b',
 'asxl2',
 'atm',
 'atrx',
 'aurka',
 'axin1',
 'b2m',
 'bap1',
 'bard1',
 'bcl10',
 'bcl2',
 'bcl2l11',
 'bcor',
 'braf',
 'brcal',
 'brca2',
 'brd4',
 'brip1',
 'btk',
 'card11',
 'carm1',
 'casp8',
 'cbl',
 'ccnd1',
 'ccnd2',
 'ccnd3',
 'ccne1',
 'cdh1',
 'cdk12',
 'cdk4',
 'cdk6',
 'cdk8',
 'cdkn1a',
 'cdkn1b',
 'cdkn2a',
 'cdkn2b',
 'cdkn2c',
 'cebpalpha',
 'chek2',
 'cic',
 'crebbp',
 'ctcf',
 'ctnnb1',
 'ddr2',
 'dicer1',
 'dnmt3a',
 'dnmt3b',
 'dusp4',
 'egfr',
 'eif1ax',
 'elf3',
 'ep300',
 'epas1',
 'epcam',
 'erbB2',
 'erbB3',
 'erbB4',
 'ercc2',
 'ercc3',
 'ercc4',
 'erg',
 'errfil',
 'esr1',
 'etv1',
```

'etv6',  
'ewsrl',  
'ezh2',  
'fam58a',  
'fanca',  
'fat1',  
'fbxw7',  
'fgf19',  
'fgf3',  
'fgf4',  
'fgfr1',  
'fgfr2',  
'fgfr3',  
'fgfr4',  
'flt3',  
'foxal1',  
'foxl2',  
'foxo1',  
'foxp1',  
'fubp1',  
'gata3',  
'glil1',  
'gnaq',  
'gnas',  
'h3f3a',  
'hla',  
'hnf1a',  
'hras',  
'idh1',  
'idh2',  
'igfir',  
'ikzf1',  
'inpp4b',  
'jak1',  
'jak2',  
'kdm5a',  
'kdm5c',  
'kdm6a',  
'kdr',  
'keap1',  
'kit',  
'klf4',  
'kmt2a',  
'kmt2b',  
'kmt2c',  
'kmt2d',  
'knstrn',  
'kras',  
'lats1',  
'lats2',  
'map2k1',  
'map2k2',  
'map2k4',  
'map3k1',  
'mapk1',  
'mdm2',  
'mdm4',  
'med12',  
'mef2b',  
'men1',  
'met',  
'mga',  
'mlh1',  
'mpl',  
'msh2',  
'msh6',  
'mtor',  
'myc',  
'mycn',  
'myd88',  
'nf1',  
'nf2',  
'nfe2l2',  
'nkbia',  
'nkx2',  
'notch1',  
'notch2',  
'npm1',  
'nras',  
'nsd1',  
'ntrk1',  
'ntrk3',  
'nup93',

```
'pakk',
'pbprm1',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pim1',
'pms1',
'pms2',
'pole',
'ppp2r1a',
'ppp6c',
'prdm1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rab35',
'rac1',
'rad21',
'rad50',
'rad51c',
'raf1',
'rara',
'rasal1',
'rb1',
'rbm10',
'ret',
'rheb',
'rhoa',
'rictor',
'rit1',
'ros1',
'rras2',
'runx1',
'rxra',
'sdhb',
'setd2',
'sf3b1',
'shq1',
'smad2',
'smad3',
'smad4',
'smarca4',
'smo',
'sos1',
'sox9',
'spop',
'src',
'stat3',
'stk11',
'tcf3',
'tert',
'tet1',
'tet2',
'tgfb1',
'tgfb2',
'tmprss2',
'tp53',
'tp53bp1',
'tsc1',
'tsc2',
'u2af1',
'vegfa',
'vh1',
'whsc1',
'xpol',
'xrcc2',
'yap1']
```

In [16]:

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

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 232)
```

## Variation Feature

In [17]:

```
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'])
```

In [18]:

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

```
train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The s
hape of Variation feature: (2124, 2069)
```

## Text feature

In [19]:

```
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 [20]:

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

In [21]:

```
building a CountVectorizer with all the words that occured minimum 3 times in train data
text_vectorizer = CountVectorizer(ngram_range=(1,4),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 features) vector
train_text_fea_counts = train_text_feature_onehotCoding.sum(axis=0).A1

zip(list(text_features),text_fea_counts) will zip a word with its number of times it occured
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 : 2881628

In [22]:

```
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 [23]:

```
#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 [24]:

```
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 [25]:

```
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 [26]:

```
#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 [27]:

```
#Data preparation for ML models.

#Misc. functionns for ML models

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

 # for calculating log_loss we will provide the array of probabilities belongs to each class
 print("Log loss :",log_loss(test_y, sig_clf.predict_proba(test_x)))
 # calculating the number of data points that are misclassified
 print("Number of mis-classified points :", np.count_nonzero((pred_y- test_y))/test_y.shape[0])
 plot_confusion_matrix(test_y, pred_y)
```

In [28]:

```
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 [29]:

```
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_imptfeature_names(indices, text, gene, var, no_features):
 gene_count_vec = CountVectorizer()
 var_count_vec = CountVectorizer()
 text_count_vec = CountVectorizer(min_df=3)

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

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

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

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

## Stacking the three types of features

In [38]:

```
merging gene, variance and text features

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

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

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

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

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

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

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

In [39]:

```
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, 2883929)
(number of data points * number of features) in test data = (665, 2883929)
(number of data points * number of features) in cross validation data = (532, 2883929)
```

In [40]:

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

Response encoding features :

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

## Logistic Regression

### With Class balancing

#### Hyper parameter tuning

In [42]:

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

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

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

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

find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html

default paramters
sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=3)
#
some of the methods of CalibratedClassifierCV()
fit(X, y[, sample_weight]) Fit the calibrated model
get_params([deep]) Get parameters for this estimator.
predict(X) Predict the target of new samples.
predict_proba(X) Posterior probabilities of classification

video link:

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

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

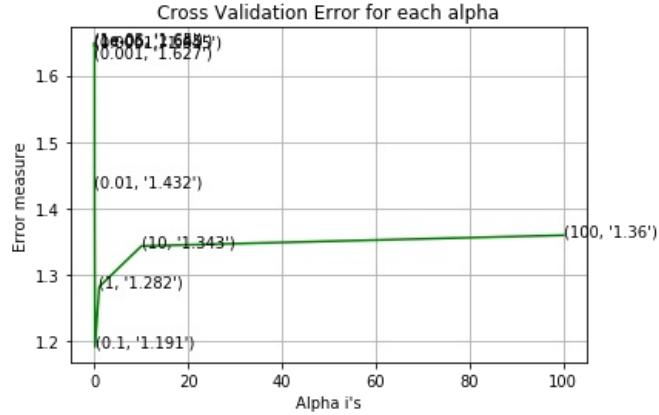
best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

```

for alpha = 1e-06
Log Loss : 1.6454082296645487
for alpha = 1e-05
Log Loss : 1.6496453000987308
for alpha = 0.0001
Log Loss : 1.6449148263389797
for alpha = 0.001
Log Loss : 1.6273863040929106
for alpha = 0.01
Log Loss : 1.4319798667951187
for alpha = 0.1
Log Loss : 1.1914005893415633
for alpha = 1
Log Loss : 1.282169036398474
for alpha = 10
Log Loss : 1.3433645728161845
for alpha = 100
Log Loss : 1.3595995178981684

```



For values of best alpha = 0.1 The train log loss is: 0.9188167909803133  
 For values of best alpha = 0.1 The cross validation log loss is: 1.1914005893415633  
 For values of best alpha = 0.1 The test log loss is: 1.2654942000313893

### Testing the model with best hyper parameters

In [46]:

```

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

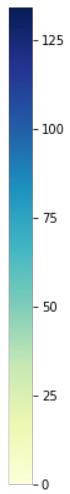
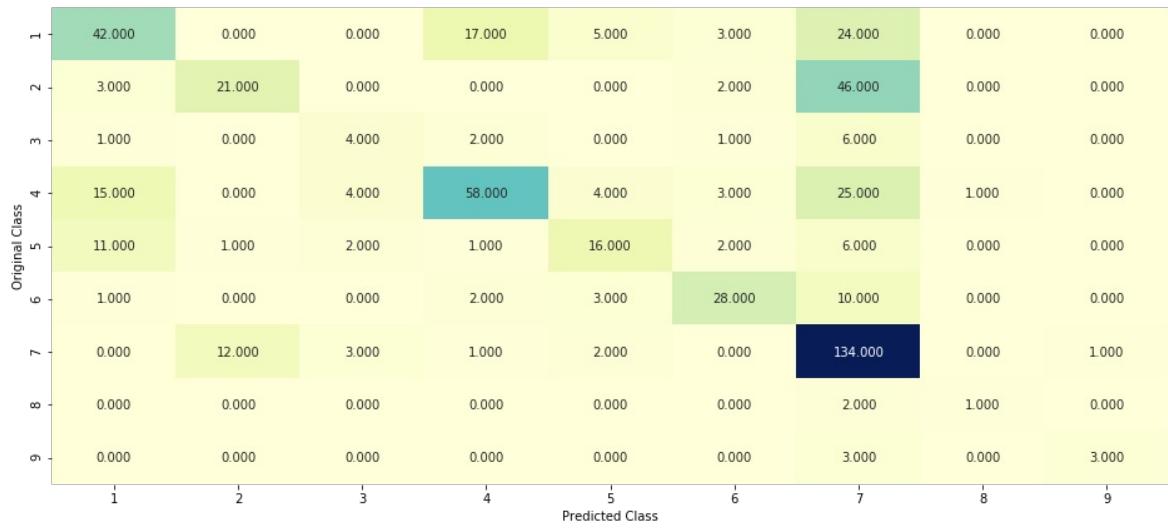
default parameters
SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
class_weight=None, warm_start=False, average=False,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

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

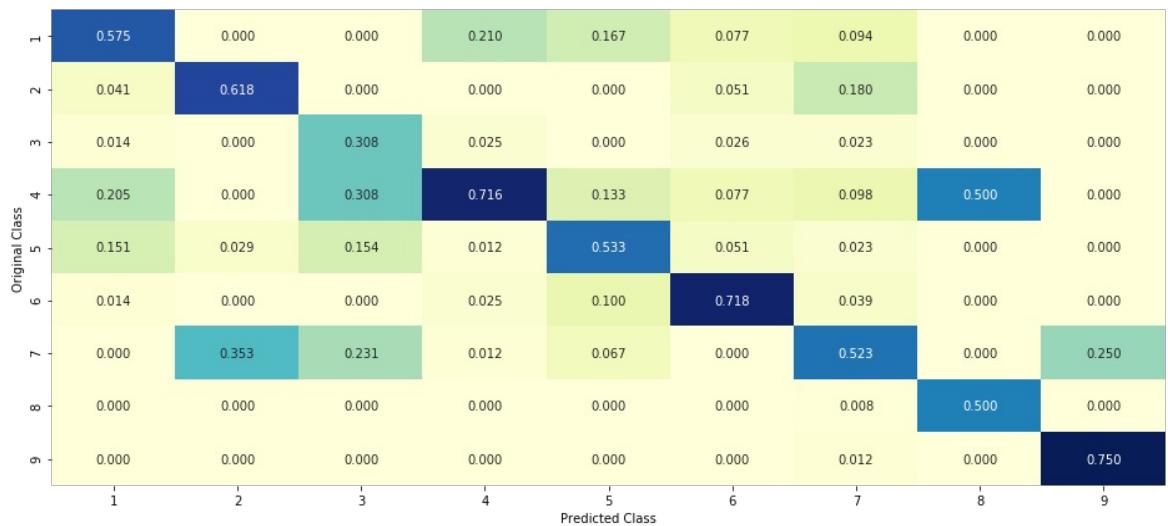
#-----
video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1/
#-----
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding,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

```

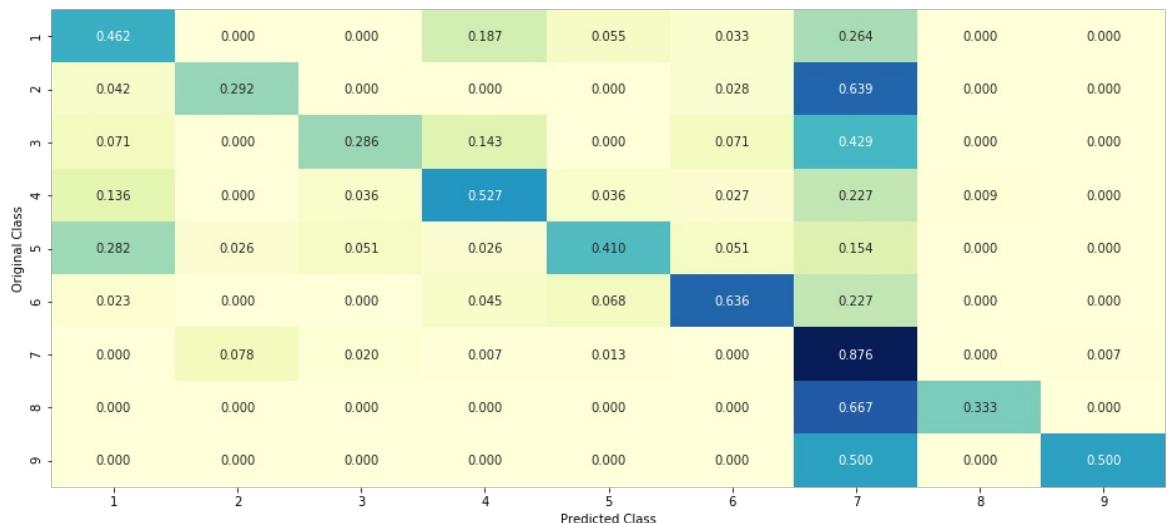
Log loss : 1.1914005893415633  
 Number of mis-classified points : 0.42293233082706766  
 ----- Confusion matrix -----



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



## Summary

In [1]:

```
#http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
X=PrettyTable()
X.field_names=["S.NO", "Model", "Train", "Cv", "Test", "% Misclass"]
X.add_row(["1", "Naive Bayes", 0.9206, 1.1816, 1.1818, 0.389])
X.add_row(["2", "KNN", 0.8749, 1.0737, 1.1127, 0.377])
X.add_row(["3", "LR(with class balancing)", 0.5720, 1.1080, 1.0335, 0.385])
X.add_row(["4", "LR(with class balancing)", 0.5662, 1.1410, 1.0581, 0.377])
X.add_row(["5", "Linear Support Vector Machines", 0.6839, 1.1649, 1.0889, 0.370])
X.add_row(["6", "RF(One hot encoding)", 0.6433, 1.1135, 1.1381, 0.381])
X.add_row(["7", "RF(Response coding)", 0.0566, 1.3436, 1.4014, 0.492])
X.add_row(["8", "stacking(lR+SVM+Naive bayes)", 0.6345, 1.1116, 1.1058, 0.354])
X.add_row(["9", "maximum Voting classifier", 0.8621, 1.1355, 1.1371, 0.353])
X.add_row(["10", "top 1000 tfidf words(Naive bayes)", 0.7826, 1.1402, 1.1567, 0.355])
X.add_row(["11", "top 1000 tfidf words(knn)", 0.8291, 1.0395, 1.0422, 0.360])
X.add_row(["12", "top 1000 tfidf words(LR with class balancing)", 0.5948, 0.9726, 0.9795, 0.349])
X.add_row(["13", "top 1000 tfidf words(LR with out class balancing)", 0.5808, 1.0168, 1.0097, 0.345])
X.add_row(["14", "top 1000 tfidf words(Linear svm)", 0.6767, 1.0465, 1.0785, 0.347])
X.add_row(["15", "top 1000 tfidf words(RF one hot encoding)", 0.8566, 1.2082, 1.1752, 0.439])
X.add_row(["16", "top 1000 tfidf words(RF response coding)", 0.0733, 1.2091, 1.2563, 0.404])
X.add_row(["17", "top 1000 tfidf words(stacking (LR+SVM+naive bayes))", 0.8296, 1.2048, 1.1070, 0.365])
X.add_row(["18", "top 1000 tfidf words(maximum voting classifier)", 0.9695, 1.2034, 1.1656, 0.342])
X.add_row(["19", "Countvectrizer including both unigrams bigrams ", 0.8441, 1.1919, 1.3164, 0.394])
X.add_row(["20", "countvecrizer including both unigrams bigrams ", 0.8267, 1.12285, 1.1291, 0.375])
X.add_row(["21", "feature engineering technique", 0.9188, 1.1914, 1.2654, 0.422])
print(X)
```

| S.NO | Model                                               | Train  | Cv      | Test   | % Misclas |
|------|-----------------------------------------------------|--------|---------|--------|-----------|
| 1    | Naive Bayes                                         | 0.9206 | 1.1816  | 1.1818 | 0.389     |
| 2    | KNN                                                 | 0.8749 | 1.0737  | 1.1127 | 0.377     |
| 3    | LR(with class balancing)                            | 0.572  | 1.108   | 1.0335 | 0.385     |
| 4    | LR(with class balancing)                            | 0.5662 | 1.141   | 1.0581 | 0.377     |
| 5    | Linear Support Vector Machines                      | 0.6839 | 1.1649  | 1.0889 | 0.37      |
| 6    | RF(One hot encoding)                                | 0.6433 | 1.1135  | 1.1381 | 0.381     |
| 7    | RF(Response coding)                                 | 0.0566 | 1.3436  | 1.4014 | 0.492     |
| 8    | stacking(lR+SVM+Naive bayes)                        | 0.6345 | 1.1116  | 1.1058 | 0.354     |
| 9    | maximum Voting classifier                           | 0.8621 | 1.1355  | 1.1371 | 0.353     |
| 10   | top 1000 tfidf words(Naive bayes)                   | 0.7826 | 1.1402  | 1.1567 | 0.355     |
| 11   | top 1000 tfidf words(knn)                           | 0.8291 | 1.0395  | 1.0422 | 0.36      |
| 12   | top 1000 tfidf words(LR with class balancing)       | 0.5948 | 0.9726  | 0.9795 | 0.349     |
| 13   | top 1000 tfidf words(LR with out class balancing)   | 0.5808 | 1.0168  | 1.0097 | 0.345     |
| 14   | top 1000 tfidf words(Linear svm)                    | 0.6767 | 1.0465  | 1.0785 | 0.347     |
| 15   | top 1000 tfidf words(RF one hot encoding)           | 0.8566 | 1.2082  | 1.1752 | 0.439     |
| 16   | top 1000 tfidf words(RF response coding)            | 0.0733 | 1.2091  | 1.2563 | 0.404     |
| 17   | top 1000 tfidf words(stacking (LR+SVM+naive bayes)) | 0.8296 | 1.2048  | 1.107  | 0.365     |
| 18   | top 1000 tfidf words(maximum voting classifier)     | 0.9695 | 1.2034  | 1.1656 | 0.342     |
| 19   | Countvectrizer including both unigrams bigrams      | 0.8441 | 1.1919  | 1.3164 | 0.394     |
| 20   | countvecrizer including both unigrams bigrams       | 0.8267 | 1.12285 | 1.1291 | 0.375     |
| 21   | feature engineering technique                       | 0.9188 | 1.1914  | 1.2654 | 0.422     |

