

Personalized cancer diagnosis

1. Business Problem

1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

Context:

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

Problem statement :

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

1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

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

1.3. Real-world/Business objectives and constraints.

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

2. Machine Learning Problem Formulation

2.1. Data

2.1.1. Data Overview

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

2.1.2. Example Data Point

training_variants

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

training_text

ID, Text

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

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

2.2.1. Type of Machine Learning Problem

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

2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

2.3. Train, CV and Test Datasets

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

3. Exploratory Data Analysis

In [138]:

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

from mlxtend.classifier import StackingClassifier

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

3.1. Reading Data

3.1.1. Reading Gene and Variation Data

In [139]:

```
data = pd.read_csv('training/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[139]:

ID	Gene	Variation	Class
----	------	-----------	-------

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

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

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

3.1.2. Reading Text Data

In [140]:

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

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

Out[140]:

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

3.1.3. Preprocessing of text

In [141]:

```
# 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:
```

```

    if not word in stop_words:
        string += word + " "

data_text[column][index] = string

```

In [142]:

```

#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 : 4418.609856500001 seconds

```

In [143]:

```

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

```

Out[143]:

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

In [144]:

```

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

```

Out[144]:

	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 [145]:

```

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

```

In [146]:

```

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

```

Out[146]:

	ID	Gene	Variation	Class	TEXT
--	----	------	-----------	-------	------

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

3.1.4. Test, Train and Cross Validation Split

3.1.4.1. Splitting data into train, test and cross validation (60:20:20)

In [252]:

```
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 [253]:

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

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

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

In [149]:

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

print(train_class_distribution)
print(test_class_distribution)
print(cv_class_distribution)
```

```
1    386
2    307
3     60
4    466
5    165
6    187
7    648
8     13
9     25
Name: Class, dtype: int64
1     85
2     68
3     14
4    103
5     36
6     41
7    143
8      3
-     -
```

```

9         6
Name: Class, dtype: int64
1         97
2         77
3         15
4        117
5         41
6         47
7        162
8          3
9          6
Name: Class, dtype: int64

```

In [150]:

```

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

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

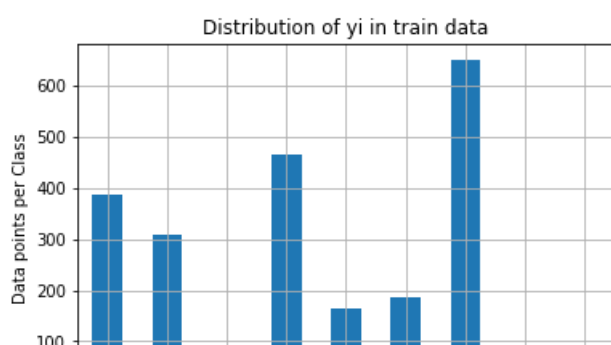
print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

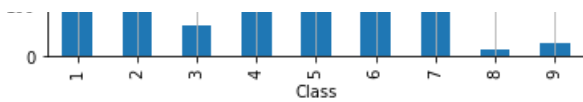
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

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

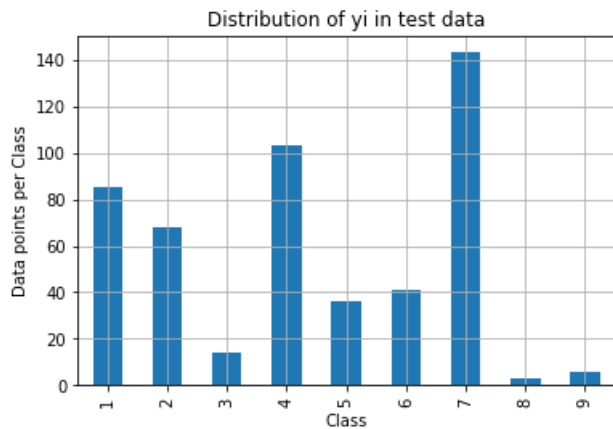
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

```

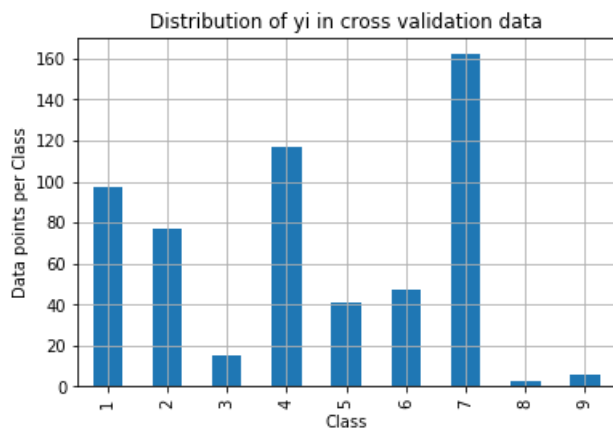




Number of data points in class 7 : 648 (28.711 %)
 Number of data points in class 4 : 466 (20.647 %)
 Number of data points in class 1 : 386 (17.102 %)
 Number of data points in class 2 : 307 (13.602 %)
 Number of data points in class 6 : 187 (8.285 %)
 Number of data points in class 5 : 165 (7.311 %)
 Number of data points in class 3 : 60 (2.658 %)
 Number of data points in class 9 : 25 (1.108 %)
 Number of data points in class 8 : 13 (0.576 %)



Number of data points in class 7 : 143 (28.657 %)
 Number of data points in class 4 : 103 (20.641 %)
 Number of data points in class 1 : 85 (17.034 %)
 Number of data points in class 2 : 68 (13.627 %)
 Number of data points in class 6 : 41 (8.216 %)
 Number of data points in class 5 : 36 (7.214 %)
 Number of data points in class 3 : 14 (2.806 %)
 Number of data points in class 9 : 6 (1.202 %)
 Number of data points in class 8 : 3 (0.601 %)



Number of data points in class 7 : 162 (28.673 %)
 Number of data points in class 4 : 117 (20.708 %)
 Number of data points in class 1 : 97 (17.168 %)
 Number of data points in class 2 : 77 (13.628 %)
 Number of data points in class 6 : 47 (8.319 %)
 Number of data points in class 5 : 41 (7.257 %)
 Number of data points in class 3 : 15 (2.655 %)
 Number of data points in class 9 : 6 (1.062 %)
 Number of data points in class 8 : 3 (0.531 %)

3.2 Prediction using a 'Random' Model

9.2 Prediction using a Random Model

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

In [151]:

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

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

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

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

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

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

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

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

In [152]:

```
# 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(v cv.cv predicted v. 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))

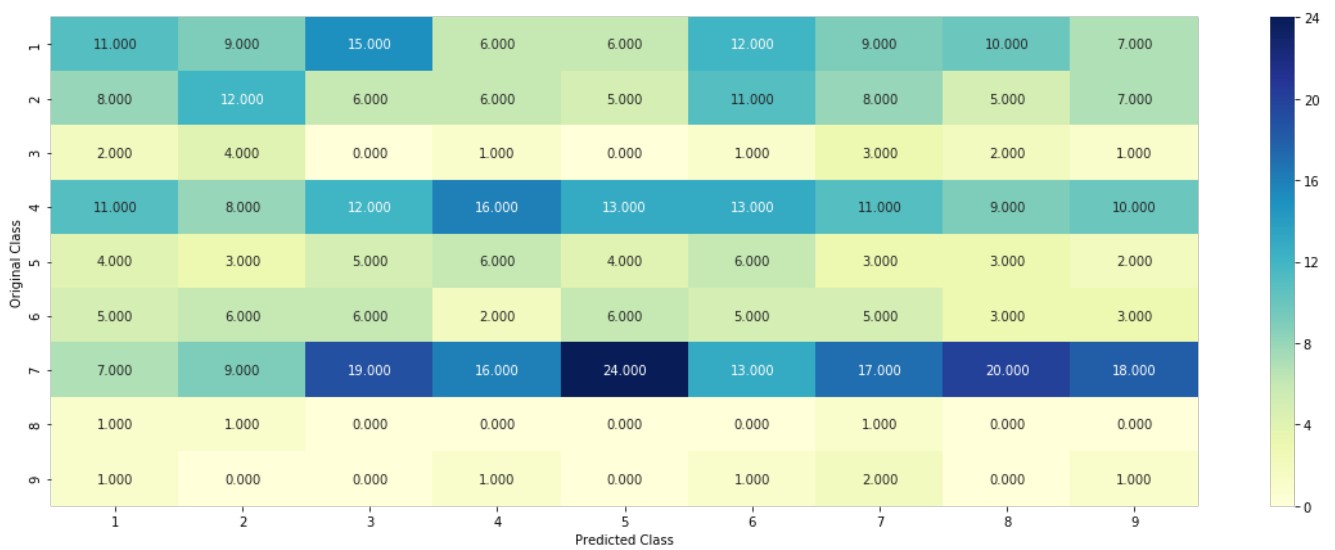
predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

```

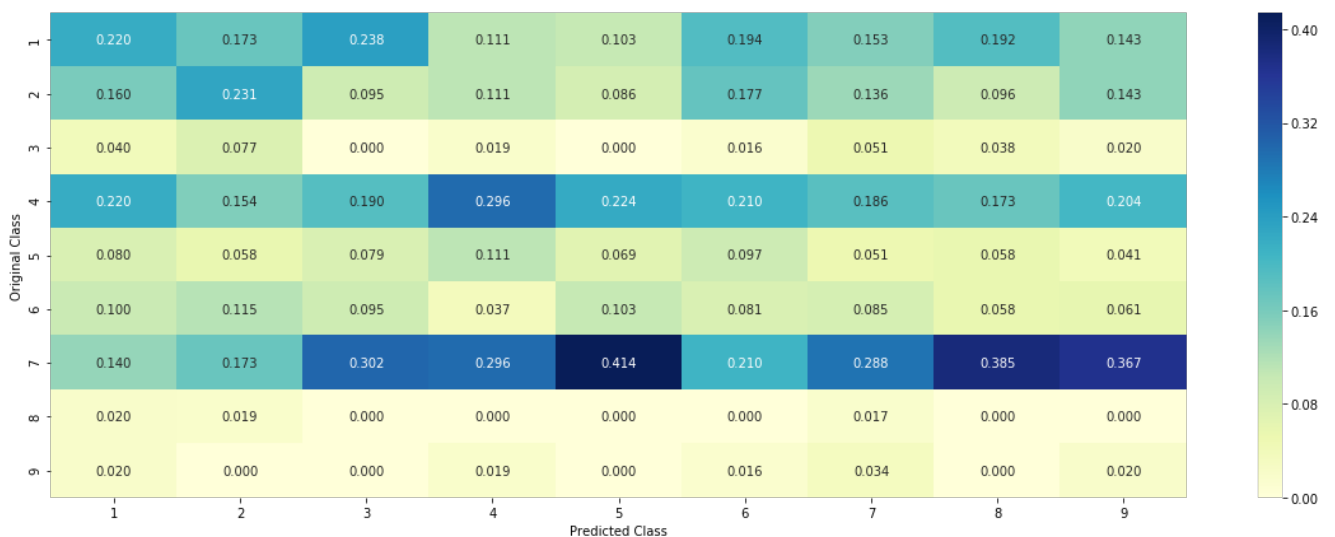
Log loss on Cross Validation Data using Random Model 2.4168088076223566

Log loss on Test Data using Random Model 2.488997692134992

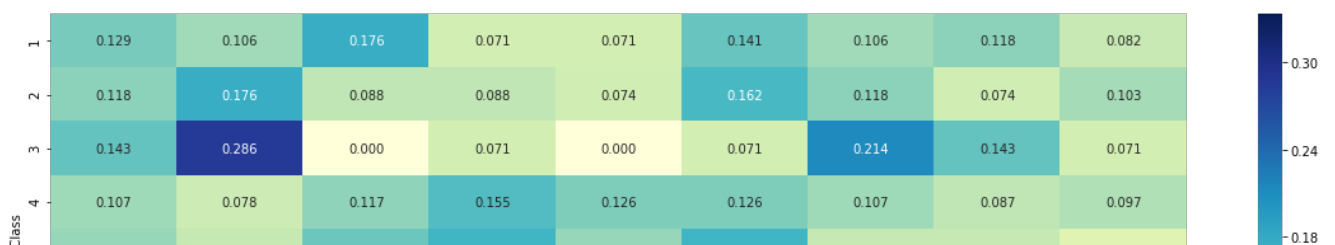
----- Confusion matrix -----

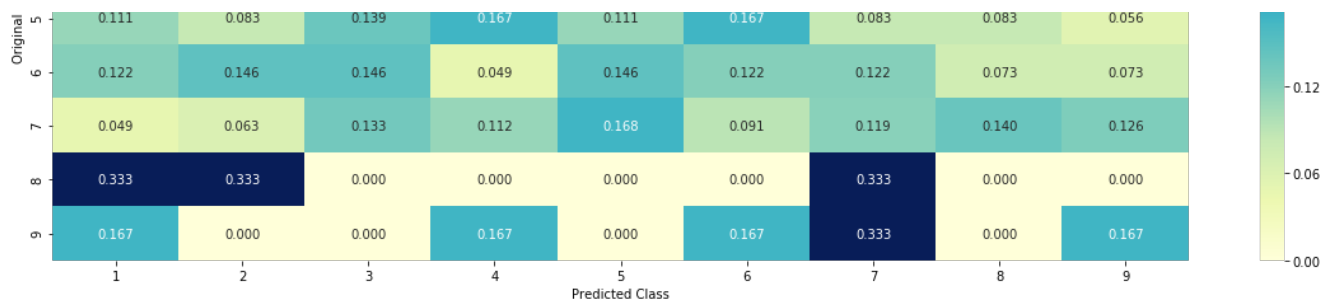


----- Precision matrix (Column Sum=1) -----



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





3.3 Univariate Analysis

In [153]:

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

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

```

# print(train_df.loc[train_df['Class']==1] & (train_df['Gene']=='BRCA1'])
#
# ID      Gene      Variation  Class
# 2470    2470    BRCA1          S1715C      1
# 2486    2486    BRCA1          S1841R      1
# 2614    2614    BRCA1          M1R        1
# 2432    2432    BRCA1          L1657P      1
# 2567    2567    BRCA1          T1685A      1
# 2583    2583    BRCA1          E1660G      1
# 2634    2634    BRCA1          W1718L      1
# cls_cnt.shape[0] will return the number of rows

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

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):
    # print(gv_dict)
    # {'BRCA1': [0.20075757575757575, 0.03787878787878788, 0.06818181818181817, 0.13636363636363635, 0.25, 0.19318181818181818, 0.03787878787878788, 0.03787878787878788, 0.03787878787878788],
    #      'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366, 0.27040816326530615, 0.061224489795918366, 0.066326530612244902, 0.051020408163265307, 0.051020408163265307, 0.056122448979591837],
    #      'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.06818181818181817, 0.06818181818181817, 0.0625, 0.34659090909090912, 0.0625, 0.056818181818181816],
    #      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608, 0.07878787878787878, 0.13939393939393939, 0.34545454545454546, 0.060606060606060608, 0.060606060606060608, 0.060606060606060608],
    #      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917, 0.46540880503144655, 0.075471698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081761006289, 0.062893081761006289],
    #      'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295, 0.072847682119205295, 0.066225165562913912, 0.066225165562913912, 0.27152317880794702, 0.066225165562913912, 0.066225165562913912],
    #      'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333333333333334, 0.073333333333333334, 0.09333333333333333, 0.080000000000000002, 0.29999999999999999, 0.066666666666666666, 0.066666666666666666],
    #      ...
    #    }
    gv_dict = get_gv_fea_dict(alpha, feature, df)
    # value_count is similar in get_gv_fea_dict
    value_count = train_df[feature].value_counts()

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

```

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

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

3.2.1 Univariate Analysis on Gene Feature

Q1. Gene, What type of feature it is ?

Ans. Gene is a categorical variable

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

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

In [154]:

```
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      170
TP53       112
EGFR        90
BRCA2       89
PTEN        83
KIT          63
BRAF         61
ERBB2        50
PDGFRA       49
ALK          48
Name: Gene, dtype: int64
```

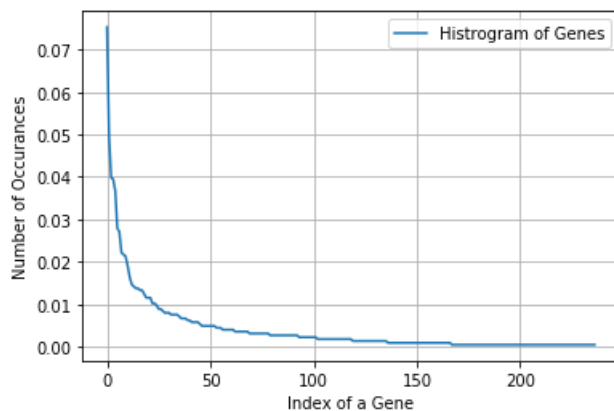
In [155]:

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

Ans: There are 237 different categories of genes in the train data, and they are distributed as follows

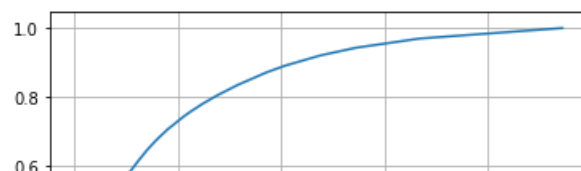
In [156]:

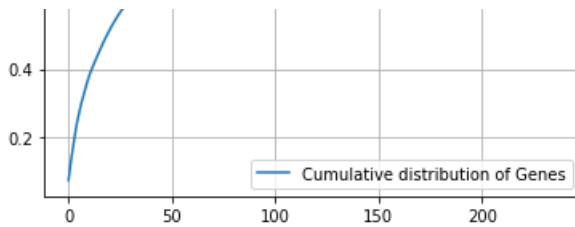
```
s = sum(unique_genes.values);
h = unique_genes.values/s;
plt.plot(h, label="Histogram of Genes")
plt.xlabel('Index of a Gene')
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```



In [157]:

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





Q3. How to featurize this Gene feature ?

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

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

1. One hot Encoding
2. Response coding

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

In [158]:

```
#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 [159]:

```
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: (2257, 9)

In [160]:

```
# one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer(ngram_range=(1, 5), min_df=5, max_features=2000)
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 [161]:

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

Out[161]:

```
2118      CCND1
1172      PIK3CA
2531      BRCA1
1319      MLH1
1587      FOXO1
Name: Gene, dtype: object
```

In [162]:

```
print(len(gene_vectorizer.get_feature_names()))
```

In [163]:

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

train_gene_feature is converted feature using one-hot encoding method. The shape of gene feature: (2257, 102)

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

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

In [164]:

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

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

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

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

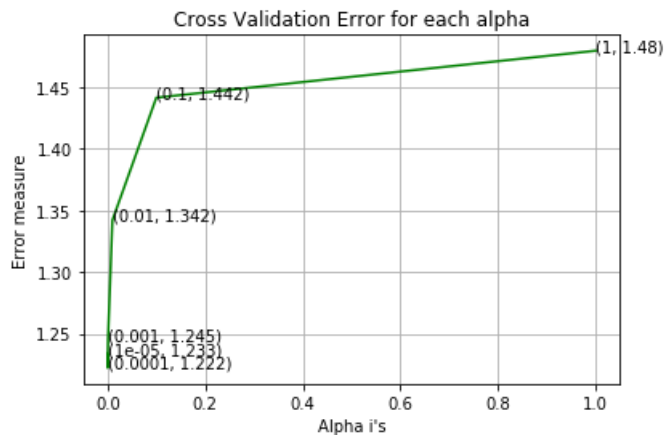
cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_gene_feature_onehotCoding, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_gene_feature_onehotCoding, y_train)
    predict_y =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

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

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

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

For values of alpha = 1e-05 The log loss is: 1.233082559800587
 For values of alpha = 0.0001 The log loss is: 1.2221940840152774
 For values of alpha = 0.001 The log loss is: 1.2450805756151184
 For values of alpha = 0.01 The log loss is: 1.342117024105676
 For values of alpha = 0.1 The log loss is: 1.4416912929150347
 For values of alpha = 1 The log loss is: 1.479821596737203



For values of best alpha = 0.0001 The train log loss is: 1.0976991693026268
 For values of best alpha = 0.0001 The cross validation log loss is: 1.2221940840152774
 For values of best alpha = 0.0001 The test log loss is: 1.2435308521810813

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

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

In [165]:

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

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

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

Q6. How many data points in Test and CV datasets are covered by the 237 genes in train dataset?

Ans

1. In test data 482 out of 499 : 96.59318637274549
2. In cross validation data 547 out of 565 : 96.8141592920354

3.2.2 Univariate Analysis on Variation Feature

Q7. Variation, What type of feature is it ?

Ans. Variation is a categorical variable

Q8. How many categories are there?

In [166]:

```
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 : 2051
 Truncating_Mutations 63
 Deletion 54


```

Deletion          31
Amplification     42
Fusions           23
Overexpression    4
Q61R              3
I31M              2
Q61H              2
T286A             2
Y64A              2
Name: Variation, dtype: int64

```

In [167]:

```

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

```

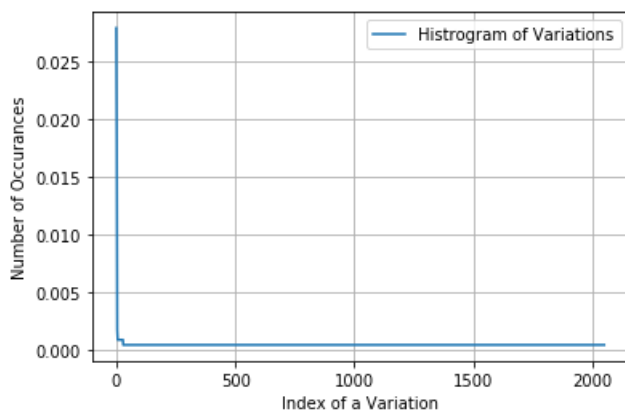
Ans: There are 2051 different categories of variations in the train data, and they are distributed as follows

In [168]:

```

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

```



In [169]:

```

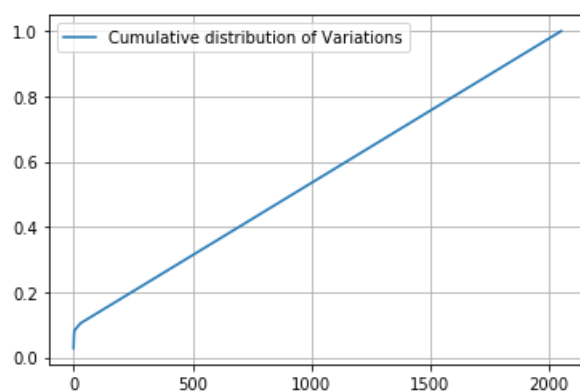
c = np.cumsum(h)
print(c)
plt.plot(c, label='Cumulative distribution of Variations')
plt.grid()
plt.legend()
plt.show()

```

```

[0.02791316 0.05183872 0.0704475 ... 0.99911387 0.99955693 1.
]

```



Q9. How to featurize this Variation feature ?

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

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

1. One hot Encoding
2. Response coding

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

In [170]:

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

In [171]:

```
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: (2257, 9)

In [172]:

```
# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer(ngram_range=(1, 5), min_df=5, max_features=2000)
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 [173]:

```
print(len(variation_vectorizer.get_feature_names()))
```

11

In [174]:

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

train_variation_feature_onehotEncoded is converted feature using the one-hot encoding method. The shape of Variation feature: (2257, 11)

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

Let's build a model just like the earlier!

In [175]:

```
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=1000, tol=1e-05)
```

```

clf=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.clas
ses_, eps=1e-15))

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

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

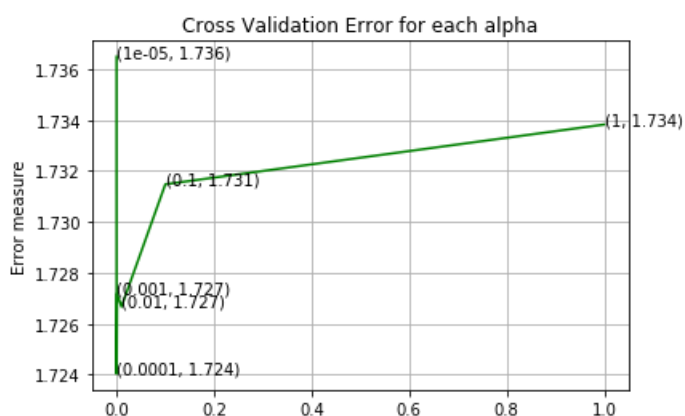
predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
print('For values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

```

```

For values of alpha = 1e-05 The log loss is: 1.7364856255686358
For values of alpha = 0.0001 The log loss is: 1.724036523659092
For values of alpha = 0.001 The log loss is: 1.7271889292854585
For values of alpha = 0.01 The log loss is: 1.7266729063976909
For values of alpha = 0.1 The log loss is: 1.7314793293247905
For values of alpha = 1 The log loss is: 1.733823844396466

```



For values of best alpha = 0.0001 The train log loss is: 1.704089251690859
 For values of best alpha = 0.0001 The cross validation log loss is: 1.724036523659092
 For values of best alpha = 0.0001 The test log loss is: 1.7198633130580128

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

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

In [176]:

```
print("Q12. How many data points are covered by total ", unique_variations.shape[0], "variations i
n 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.s
hape[0])*100)
```

Q12. How many data points are covered by total 2051 variations in test and cross validation data sets?

Ans

1. In test data 59 out of 499 : 11.823647294589177
2. In cross validation data 53 out of 565 : 9.380530973451327

3.2.3 Univariate Analysis on Text Feature

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

In [177]:

```
# 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 [178]:

```
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 [179]:

```
# building a TfidfVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer = TfidfVectorizer(ngram_range=(1, 5), min_df=5, max_features=2000)
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 : 2000

In [180]:

```
print(len(text_vectorizer.get_feature_names()))
```

2000

In [181]:

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

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

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

In [182]:

```
#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 [183]:

```
# 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 [184]:

```
# 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 [185]:

```
#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 [186]:

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

```
Counter({8.594153366286104: 10, 12.891230049429144: 5, 6.29229028774205: 4, 6.501033802781964: 3,
10.43695953849777: 2, 10.125724625153829: 2, 9.252686666484626: 2, 8.938097459995912: 2,
7.839321115457525: 2, 7.830031886716227: 2, 7.692142720084124: 2, 7.613830451842098: 2,
7.309696865785778: 2, 6.986811562377829: 2, 6.56221890614874: 2, 6.55752003649352: 2,
4.888952654439375: 2, 225.42948355392306: 1, 149.66390074149146: 1, 119.58323799868397: 1,
113.42884656382635: 1, 108.29047521155961: 1, 102.79411351698076: 1, 98.04651807497045: 1,
95.94233220515461: 1, 95.59428654218644: 1, 95.56880264630162: 1, 94.04934953104438: 1,
92.46502236926922: 1, 80.3897877782465: 1, 77.02857303330072: 1, 75.40752555013351: 1,
72.99049976098365: 1, 72.9079752783215: 1, 71.08319234152768: 1, 70.1077576731767: 1,
69.72327060475827: 1, 67.5668681315225: 1, 67.5052473390158: 1, 64.97326880471562: 1,
61.386117653738076: 1, 61.120336263173186: 1, 60.9058592966349: 1, 59.11763474030046: 1,
57.4597626894257: 1, 57.146900110091586: 1, 56.2412403834435: 1, 56.216897363398225: 1,
55.911968457134904: 1, 55.639389949686425: 1, 54.649852904359506: 1, 54.024669254980786: 1,
52.997642800262376: 1, 51.71456016847814: 1, 51.29253509811004: 1, 50.45035342704496: 1, 49.5077930
7361812: 1, 49.39853754442469: 1, 48.911284911172125: 1, 48.90625858797139: 1, 46.45134271446248:
1, 44.25426785401855: 1, 43.63584841527159: 1, 43.005137142204255: 1, 42.54781442049824: 1,
41.46714481887317: 1, 41.17796154822636: 1, 40.90767554314691: 1, 40.00092611981866: 1,
39.297472159950125: 1, 38.91126694990929: 1, 38.36159527149141: 1, 38.26827481468472: 1, 38.2219365
2816724: 1, 38.13178128216346: 1, 38.11997521163218: 1, 37.769199531726436: 1, 37.3952120880463:
1, 37.0539937513755: 1, 36.94872642312225: 1, 36.79874949315312: 1, 36.722884855915474: 1,
36.462319214835084: 1, 36.24937281837738: 1, 36.17247134487735: 1, 36.03614972356889: 1, 35.9650802
8669837: 1, 35.79961916251964: 1, 35.558069918272274: 1, 35.01731366293294: 1, 34.865821663637796:
1, 34.82488799094069: 1, 34.28798708836485: 1, 33.58092922550729: 1, 33.16279337003154: 1,
32.63311054262648: 1, 32.22301908506014: 1, 31.87571132514676: 1, 31.428536690412134: 1,
31.080811895820005: 1, 30.909166025879877: 1, 30.86053505466067: 1, 30.844132497041617: 1,
30.826597199869656: 1, 30.77601405778812: 1, 30.734549870714027: 1, 30.703447812559958: 1,
30.420905208132407: 1, 30.154764206081317: 1, 30.060305877260024: 1, 29.94692126236269: 1,
29.84226184218965: 1, 29.83841705956542: 1, 29.50903161136105: 1, 29.39027796186567: 1,
29.24441114933914: 1, 28.801405656020567: 1, 28.715214439309115: 1, 28.612844138008594: 1,
28.50852432654113: 1, 28.306832137796306: 1, 28.200470552059084: 1, 28.072248611917885: 1,
28.030002251168845: 1, 27.91880346122593: 1, 27.825368441717742: 1, 27.776632403029964: 1,
27.7224434181262: 1, 27.643154573860485: 1, 27.628712682041915: 1, 27.5945169070444: 1,
27.569948277234587: 1, 27.444216560506973: 1, 27.249889584545702: 1, 26.795545922219834: 1,
26.657563368917682: 1, 26.623677309853694: 1, 26.405507352275034: 1, 26.354424453431296: 1,
26.17372788193746: 1, 26.112290717863065: 1, 26.092212517238558: 1, 25.98129838268241: 1,
25.78024598359976: 1, 25.722342135994857: 1, 25.647526136366213: 1, 25.62035925547689: 1,
25.500860283562627: 1, 25.414167945405808: 1, 25.330079674154042: 1, 25.232641029819074: 1,
25.12336518139868: 1, 25.078416706396077: 1, 24.96092866638543: 1, 24.91175363635633: 1,
24.800106068185863: 1, 24.713621807754816: 1, 24.20499521542631: 1, 24.103976421790747: 1,
23.907206564566064: 1, 23.900642526916457: 1, 23.87867822713171: 1, 23.845139094903306: 1,
23.766547112357795: 1, 23.759128191336842: 1, 23.571470112737313: 1, 23.34457613245488: 1,
23.28779645477189: 1, 23.12103761148571: 1, 23.0441728572191: 1, 22.93335236498018: 1,
22.88991600016289: 1, 22.85345700484406: 1, 22.668383908107728: 1, 22.532092885838615: 1,
22.50582860913214: 1, 22.49905370568828: 1, 22.370775901711408: 1, 22.366483218808153: 1,
22.192814445844515: 1, 22.14799002911057: 1, 22.02606754461766: 1, 21.95631268675962: 1, 21.8317041
7193158: 1, 21.784666452670336: 1, 21.774928033953593: 1, 21.731397036627175: 1,
21.711100427311514: 1, 21.620193779878253: 1, 21.56552304250454: 1, 21.55255456285572: 1,
21.5314311058076: 1, 21.477161540168634: 1, 21.360451533087694: 1, 21.313100424154253: 1,
21.272850849472274: 1, 21.030811941443027: 1, 20.94925884075451: 1, 20.94790439374481: 1,
20.89349061288042: 1, 20.866850039024886: 1, 20.76544961069521: 1, 20.70649364318755: 1,
20.61830355636278: 1, 20.57472649856754: 1, 20.564013257286078: 1, 20.371117928951445: 1,
20.343351896311056: 1, 20.33149993029813: 1, 20.179107828697834: 1, 20.16468241452966: 1,
20.14456546735738: 1, 20.086888555387885: 1, 19.96664776453087: 1, 19.84665267387858: 1,
19.726663555555555: 1, 19.606666666666666: 1, 19.486666666666666: 1, 19.366666666666666: 1,
19.246666666666666: 1, 19.126666666666666: 1, 19.006666666666666: 1, 18.886666666666666: 1,
18.766666666666666: 1, 18.646666666666666: 1, 18.526666666666666: 1, 18.406666666666666: 1,
18.286666666666666: 1, 18.166666666666666: 1, 18.046666666666666: 1, 17.926666666666666: 1,
17.806666666666666: 1, 17.686666666666666: 1, 17.566666666666666: 1, 17.446666666666666: 1,
17.326666666666666: 1, 17.206666666666666: 1, 17.086666666666666: 1, 16.966666666666666: 1,
16.846666666666666: 1, 16.726666666666666: 1, 16.606666666666666: 1, 16.486666666666666: 1,
16.366666666666666: 1, 16.246666666666666: 1, 16.126666666666666: 1, 16.006666666666666: 1,
15.886666666666666: 1, 15.766666666666666: 1, 15.646666666666666: 1, 15.526666666666666: 1,
15.406666666666666: 1, 15.286666666666666: 1, 15.166666666666666: 1, 15.046666666666666: 1,
14.926666666666666: 1, 14.806666666666666: 1, 14.686666666666666: 1, 14.566666666666666: 1,
14.446666666666666: 1, 14.326666666666666: 1, 14.206666666666666: 1, 14.086666666666666: 1,
13.966666666666666: 1, 13.846666666666666: 1, 13.726666666666666: 1, 13.606666666666666: 1,
13.486666666666666: 1, 13.366666666666666: 1, 13.246666666666666: 1, 13.126666666666666: 1,
13.006666666666666: 1, 12.886666666666666: 1, 12.766666666666666: 1, 12.646666666666666: 1,
12.526666666666666: 1, 12.406666666666666: 1, 12.286666666666666: 1, 12.166666666666666: 1,
12.046666666666666: 1, 11.926666666666666: 1, 11.806666666666666: 1, 11.686666666666666: 1,
11.566666666666666: 1, 11.446666666666666: 1, 11.326666666666666: 1, 11.206666666666666: 1,
11.086666666666666: 1, 10.966666666666666: 1, 10.846666666666666: 1, 10.726666666666666: 1,
10.606666666666666: 1, 10.486666666666666: 1, 10.366666666666666: 1, 10.246666666666666: 1,
10.126666666666666: 1, 10.006666666666666: 1, 9.886666666666666: 1, 9.766666666666666: 1,
9.646666666666666: 1, 9.526666666666666: 1, 9.406666666666666: 1, 9.286666666666666: 1,
9.166666666666666: 1, 9.046666666666666: 1, 8.926666666666666: 1, 8.806666666666666: 1,
8.686666666666666: 1, 8.566666666666666: 1, 8.446666666666666: 1, 8.326666666666666: 1,
8.206666666666666: 1, 8.086666666666666: 1, 7.966666666666666: 1, 7.846666666666666: 1,
7.726666666666666: 1, 7.606666666666666: 1, 7.486666666666666: 1, 7.366666666666666: 1,
7.246666666666666: 1, 7.126666666666666: 1, 7.006666666666666: 1, 6.886666666666666: 1,
6.766666666666666: 1, 6.646666666666666: 1, 6.526666666666666: 1, 6.406666666666666: 1,
6.286666666666666: 1, 6.166666666666666: 1, 6.046666666666666: 1, 5.926666666666666: 1,
5.806666666666666: 1, 5.686666666666666: 1, 5.566666666666666: 1, 5.446666666666666: 1,
5.326666666666666: 1, 5.206666666666666: 1, 5.086666666666666: 1, 4.966666666666666: 1,
4.846666666666666: 1, 4.726666666666666: 1, 4.606666666666666: 1, 4.486666666666666: 1,
4.366666666666666: 1, 4.246666666666666: 1, 4.126666666666666: 1, 4.006666666666666: 1,
3.886666666666666: 1, 3.766666666666666: 1, 3.646666666666666: 1, 3.526666666666666: 1,
3.406666666666666: 1, 3.286666666666666: 1, 3.166666666666666: 1, 3.046666666666666: 1,
2.926666666666666: 1, 2.806666666666666: 1, 2.686666666666666: 1, 2.566666666666666: 1,
2.446666666666666: 1, 2.326666666666666: 1, 2.206666666666666: 1, 2.086666666666666: 1,
1.966666666666666: 1, 1.846666666666666: 1, 1.726666666666666: 1, 1.606666666666666: 1,
1.486666666666666: 1, 1.366666666666666: 1, 1.246666666666666: 1, 1.126666666666666: 1,
1.006666666666666: 1, 0.886666666666666: 1, 0.766666666666666: 1, 0.646666666666666: 1,
0.526666666666666: 1, 0.406666666666666: 1, 0.286666666666666: 1, 0.166666666666666: 1,
0.046666666666666: 1, 0.006666666666666: 1, 0.000000000000000: 1})
```

20.14405049735738: 1, 20.098209553817885: 1, 19.942694776453997: 1, 19.940953367037952: 1, 19.934182597067263: 1, 19.929560091469984: 1, 19.89968529870869: 1, 19.879716179482138: 1, 19.875510420043813: 1, 19.873150977454486: 1, 19.8559635444608382: 1, 19.79708292246175: 1, 19.7787386598821: 1, 19.776814380291665: 1, 19.760438689490027: 1, 19.70261671252069: 1, 19.698731405278156: 1, 19.677789131939488: 1, 19.666998112526553: 1, 19.573214403986274: 1, 19.344702759436274: 1, 19.330374515480106: 1, 19.27433821821475: 1, 19.175668650363857: 1, 19.15932898448668: 1, 19.12863230720731: 1, 19.12803797335822: 1, 19.0838270759914: 1, 19.016544583400183: 1, 18.986072763859827: 1, 18.908771186686852: 1, 18.89378294638672: 1, 18.883191566551798: 1, 18.852285839476867: 1, 18.8438958171887: 1, 18.826682937242616: 1, 18.78505692789075: 1, 18.773477884177538: 1, 18.769426915609216: 1, 18.736526191365044: 1, 18.728084662488474: 1, 18.699825774042715: 1, 18.65004227175297: 1, 18.546110453624635: 1, 18.536183924097244: 1, 18.488951673632716: 1, 18.458017549924758: 1, 18.443120809086416: 1, 18.42453208669081: 1, 18.332277953711024: 1, 18.2762490270142: 1, 18.234935922597774: 1, 18.19913343960927: 1, 18.14164026609705: 1, 18.13266793746642: 1, 18.107920112745333: 1, 18.081505419637374: 1, 18.022152414844587: 1, 18.01942500871764: 1, 17.990686936765822: 1, 17.98765868537402: 1, 17.96519974072061: 1, 17.91155429472569: 1, 17.879346860863144: 1, 17.87875462450218: 1, 17.867620595467965: 1, 17.852946480990298: 1, 17.82595482433231: 1, 17.825909646310613: 1, 17.77243049200857: 1, 17.700173893169946: 1, 17.698712434690933: 1, 17.69566076654217: 1, 17.655510974102693: 1, 17.56596221268967: 1, 17.514817082398: 1, 17.48296848212528: 1, 17.459355316382318: 1, 17.457965785840695: 1, 17.44630382288315: 1, 17.291606287680434: 1, 17.270218002290772: 1, 17.24153043158105: 1, 17.22152703785632: 1, 17.1550789959142: 1, 17.06296179563358: 1, 17.048322471846724: 1, 16.994868014433695: 1, 16.97467816072813: 1, 16.935727004246836: 1, 16.879536162599003: 1, 16.878335812156998: 1, 16.839713732808974: 1, 16.81042986513155: 1, 16.78165846991575: 1, 16.776957911687106: 1, 16.762487543706218: 1, 16.740006260333324: 1, 16.689124416942157: 1, 16.66585892508806: 1, 16.663027280882254: 1, 16.661340142674867: 1, 16.655425247637037: 1, 16.636104706907656: 1, 16.62188357925837: 1, 16.496582391710767: 1, 16.489617081304594: 1, 16.347786947672304: 1, 16.339522055711647: 1, 16.319477526940858: 1, 16.297367878602678: 1, 16.281316203710634: 1, 16.25371946097071: 1, 16.149608928903707: 1, 16.091868056523133: 1, 16.01720623824177: 1, 16.01400385996727: 1, 16.011228967039393: 1, 15.949643624693266: 1, 15.939883223965753: 1, 15.934080202684832: 1, 15.912957445636602: 1, 15.896512624529295: 1, 15.883356019670623: 1, 15.831770460926672: 1, 15.829083342691836: 1, 15.82672079008296: 1, 15.798314786223061: 1, 15.776953729043221: 1, 15.733507632478199: 1, 15.70697016512009: 1, 15.670737043324898: 1, 15.593111661835072: 1, 15.584318410943137: 1, 15.5466880964452: 1, 15.539632688418841: 1, 15.491741851039226: 1, 15.463009105871818: 1, 15.446553159062953: 1, 15.423433771191775: 1, 15.4145200145347: 1, 15.387962098586874: 1, 15.386048022042518: 1, 15.316451591993989: 1, 15.27722610457979: 1, 15.277109817367167: 1, 15.239295601373653: 1, 15.195188010888979: 1, 15.179528535597568: 1, 15.17588320121404: 1, 15.16638102871181: 1, 15.158039461305188: 1, 15.128146397412351: 1, 15.115594081116315: 1, 15.080520753418813: 1, 15.068133300043888: 1, 15.061371581453454: 1, 15.053579321281708: 1, 15.024014435879833: 1, 15.017314135819637: 1, 15.009129675755139: 1, 14.9888361273568: 1, 14.973734689415318: 1, 14.97199915557904: 1, 14.9447823247878565: 1, 14.919209387642011: 1, 14.90180905165812: 1, 14.891171151019655: 1, 14.890251981409603: 1, 14.867158573720241: 1, 14.847825971922893: 1, 14.837915174718018: 1, 14.83047033912047: 1, 14.810746399316656: 1, 14.782095749135438: 1, 14.764595332494414: 1, 14.702728982041009: 1, 14.685578975323017: 1, 14.606870896828049: 1, 14.585265151652143: 1, 14.5641425698146

12.354913909038611: 1, 12.33892066372216: 1, 12.323621542690718: 1, 12.322213972494978: 1,
12.322017328379372: 1, 12.31211989528383: 1, 12.302288983614062: 1, 12.301825407160884: 1,
12.294894168536713: 1, 12.274600211369945: 1, 12.264659730692827: 1, 12.257576329421221: 1,
12.253696264166674: 1, 12.238185523423978: 1, 12.218961340908367: 1, 12.214110743369943: 1,
12.204983213154136: 1, 12.130945189283137: 1, 12.124776244399495: 1, 12.11954222781882: 1,
12.110670033913356: 1, 12.067225395786815: 1, 12.062301361762275: 1, 12.048943686007645: 1,
12.027870405701638: 1, 12.024566856755069: 1, 12.00790548023223: 1, 12.007593399583648: 1,
12.007289918348926: 1, 11.959402674085666: 1, 11.900128815318997: 1, 11.89366615940601: 1,
11.87599559740856: 1, 11.875790399355443: 1, 11.863216840849255: 1, 11.816895228516955: 1,
11.81242511823557: 1, 11.789919178256339: 1, 11.766602111793588: 1, 11.725740070129122: 1,
11.719096913630295: 1, 11.700323585498243: 1, 11.668023464933803: 1, 11.666201029534255: 1,
11.643453239318369: 1, 11.642827805252223: 1, 11.637285657140518: 1, 11.632116811571205: 1,
11.628547752488904: 1, 11.623547199292853: 1, 11.6229327393979: 1, 11.610588907892252: 1,
11.577007398685593: 1, 11.568170329334484: 1, 11.566588844756016: 1, 11.513783209570853: 1,
11.509748543181345: 1, 11.500713504878616: 1, 11.480641646049605: 1, 11.458638191006802: 1,
11.440614292933859: 1, 11.424143191073206: 1, 11.41366177165945: 1, 11.39282385098021: 1,
11.379273758839753: 1, 11.374778100215297: 1, 11.352829254002948: 1, 11.351312383565197: 1,
11.350443326735604: 1, 11.348281378676926: 1, 11.340337666508349: 1, 11.263823968964907: 1,
11.262704572022118: 1, 11.254780345881203: 1, 11.240825131484494: 1, 11.21783410195861: 1,
11.212513074508507: 1, 11.17447603080127: 1, 11.166788424705214: 1, 11.161836010781942: 1,
11.131766207343993: 1, 11.130880256730864: 1, 11.09229706296338: 1, 11.083801878128236: 1,
11.083138515072521: 1, 11.06978745301316: 1, 11.035501244497983: 1, 10.964072716963956: 1,
10.94735792093913: 1, 10.941370856845394: 1, 10.939356317452695: 1, 10.937758125638712: 1,
10.934704534172688: 1, 10.928965711625292: 1, 10.91889717803017: 1, 10.90988489271989: 1,
10.886446793705202: 1, 10.881137552221864: 1, 10.871282154589768: 1, 10.85047185872942: 1,
10.850347754466611: 1, 10.842162252147176: 1, 10.834797832495244: 1, 10.819415031753959: 1,
10.808786524434085: 1, 10.783645699512569: 1, 10.779849163756134: 1, 10.752267194069445: 1,
10.747717439010865: 1, 10.692965962166504: 1, 10.67071627758446: 1, 10.652902132818468: 1,
10.643970775975129: 1, 10.630380704291166: 1, 10.623553211582722: 1, 10.61829001531341: 1,
10.59296782865387: 1, 10.588608064778597: 1, 10.567218273002215: 1, 10.551467496463607: 1,
10.527176222075827: 1, 10.525103514266023: 1, 10.522612980587052: 1, 10.493135593881718: 1,
10.471232955680788: 1, 10.446932853031798: 1, 10.435767940303519: 1, 10.42372747822137: 1,
10.39754748094216: 1, 10.38605592680967: 1, 10.379078868861487: 1, 10.377310571210723: 1,
10.372468484977981: 1, 10.367949537261667: 1, 10.366240518826618: 1, 10.36424092808549: 1,
10.351894722738113: 1, 10.351348280782341: 1, 10.349182709176288: 1, 10.332604780045548: 1,
10.322465855217436: 1, 10.312350693071844: 1, 10.301469069944424: 1, 10.296337873182159: 1,
10.259515812986434: 1, 10.24139266450659: 1, 10.235593710301467: 1, 10.232964351426789: 1,
10.220956106593611: 1, 10.21721154817877: 1, 10.211246684685062: 1, 10.200351612992785: 1,
10.17490423970536: 1, 10.150440949688289: 1, 10.150254895762625: 1, 10.127811807477899: 1,
10.118829560759353: 1, 10.115828596010862: 1, 10.090793172400687: 1, 10.080648822977256: 1,
10.069629558336915: 1, 10.067347111143286: 1, 10.057819000832422: 1, 10.052413108778088: 1,
10.040839554136548: 1, 10.037319432475098: 1, 10.02839855295126: 1, 10.022761177346581: 1,
10.006786888860724: 1, 10.000832491371755: 1, 9.992935701117489: 1, 9.984636588814787: 1,
9.965139644245939: 1, 9.961274278395297: 1, 9.949866057492384: 1, 9.93944584550195: 1,
9.908087340379496: 1, 9.889202401875282: 1, 9.8873471884605: 1, 9.886992114633346: 1,
9.885869212451084: 1, 9.880026109822403: 1, 9.86385820065295: 1, 9.86152103493594: 1,
9.85218867552253: 1, 9.835272929732453: 1, 9.829904157084348: 1, 9.825592633942108: 1,
9.821429315973305: 1, 9.796905197750661: 1, 9.791821142493468: 1, 9.790332934509987: 1,
9.783834876494128: 1, 9.777183116769871: 1, 9.773317881812462: 1, 9.759019242032732: 1,
9.758260603043302: 1, 9.737436341144846: 1, 9.73302271426518: 1, 9.71831163468062: 1,
9.714296105545642: 1, 9.700945743242228: 1, 9.700125738977253: 1, 9.694959658485974: 1,
9.68855940465038: 1, 9.685440499259295: 1, 9.683519615517984: 1, 9.681569389119622: 1,
9.678116869837588: 1, 9.673219072389145: 1, 9.660584583982057: 1, 9.659993208989851: 1,
9.659116323045597: 1, 9.653680841174907: 1, 9.64529993625919: 1, 9.620099291703376: 1,
9.61732038099862: 1, 9.6147618263851: 1, 9.609657049366456: 1, 9.60694684620856: 1,
9.575011545685298: 1, 9.57361579410812: 1, 9.568920380495019: 1, 9.53285600764345: 1,
9.518042290179988: 1, 9.5157096450277: 1, 9.514498183529206: 1, 9.506503903530067: 1,
9.48927949772396: 1, 9.447513574961725: 1, 9.443054933982978: 1, 9.430259339463452: 1,
9.421512827615302: 1, 9.41430722291173: 1, 9.39626417182657: 1, 9.393670646215519: 1,
9.391303937993055: 1, 9.390462242100595: 1, 9.383479096209797: 1, 9.378828792311108: 1,
9.36908687058936: 1, 9.366998913912884: 1, 9.347173299516976: 1, 9.337381767057531: 1,
9.332928184431463: 1, 9.32725167870737: 1, 9.313404480861658: 1, 9.298826043017943: 1,
9.297902915937103: 1, 9.293430273943653: 1, 9.291597237506428: 1, 9.29090841691673: 1,
9.287008101471095: 1, 9.275497121337454: 1, 9.26831823985191: 1, 9.235757543285892: 1,
9.23189952056598: 1, 9.211923392372686: 1, 9.202681645354094: 1, 9.1966115708858495: 1,
9.195957821233787: 1, 9.194423591555944: 1, 9.190722901343149: 1, 9.184722345196644: 1,
9.1791161277807: 1, 9.167471087813695: 1, 9.144655109721146: 1, 9.13335482354787: 1,
9.131712789674616: 1, 9.12959688308878: 1, 9.125101780559653: 1, 9.120143746604304: 1,
9.113973324964634: 1, 9.09066603503584: 1, 9.073191138820269: 1, 9.072877897454797: 1,
9.072406755339612: 1, 9.0666182348071: 1, 9.05377476255861: 1, 9.050881862848795: 1,
9.046855612702055: 1, 9.04434662208998: 1, 9.034086411463383: 1, 9.021434787140276: 1,
9.012969462512826: 1, 9.010091557140862: 1, 9.009082779812626: 1, 8.998724697603144: 1,
8.998691374750173: 1, 8.989874630445787: 1, 8.975217614814326: 1, 8.974426766837292: 1,
8.970594745683046: 1, 8.968094548635426: 1, 8.965664799575197: 1, 8.960838972316827: 1,
8.955598675450888: 1, 8.952282414416459: 1, 8.942053914635915: 1, 8.939327415823218: 1,
8.93923251749711: 1, 8.926707920156973: 1, 8.923795906901494: 1, 8.922501611248485: 1,
8.91910219458822: 1, 8.912561114912187: 1, 8.903188049195338: 1, 8.885935015698255: 1,

8.879953236558999: 1, 8.877839785205724: 1, 8.875602904599917: 1, 8.872702400855948: 1, 8.86674852221078: 1, 8.856737609285082: 1, 8.856668648528037: 1, 8.855322415774785: 1, 8.85203328866514: 1, 8.837115090727815: 1, 8.836891877217191: 1, 8.813707823316781: 1, 8.80596461888137: 1, 8.793180787829272: 1, 8.791319395402246: 1, 8.785606338072675: 1, 8.7779115578022797: 1, 8.772119725567197: 1, 8.770457735152666: 1, 8.761663510319172: 1, 8.746003550275796: 1, 8.743812517009367: 1, 8.740568884576845: 1, 8.723842998819247: 1, 8.72176492004877: 1, 8.721069602012244: 1, 8.717316764148814: 1, 8.712781674782422: 1, 8.710541281052281: 1, 8.708609034460054: 1, 8.699224340266872: 1, 8.693024712651003: 1, 8.692460883949844: 1, 8.689110125885296: 1, 8.685566667340357: 1, 8.678599155013536: 1, 8.678497538943127: 1, 8.66837083920941: 1, 8.668222397096624: 1, 8.66050447844875: 1, 8.656411928759521: 1, 8.64830673045765: 1, 8.647337754439533: 1, 8.643203127375074: 1, 8.634386724019558: 1, 8.623527655341634: 1, 8.622133888392822: 1, 8.610495727194428: 1, 8.606436400135243: 1, 8.58450929506422: 1, 8.581150098828836: 1, 8.576520486226263: 1, 8.573162226987396: 1, 8.57266881689705: 1, 8.57172875501637: 1, 8.568230499623418: 1, 8.56344372299174: 1, 8.560469645684707: 1, 8.548101538251798: 1, 8.546615028901723: 1, 8.544607786473247: 1, 8.532078643297954: 1, 8.501953298466749: 1, 8.497302220543183: 1, 8.483637938246906: 1, 8.469387858951215: 1, 8.460793602911409: 1, 8.446186596187607: 1, 8.4271835461032: 1, 8.42647134015228: 1, 8.42006251693999: 1, 8.40660858162082: 1, 8.40202581802553: 1, 8.396003501094352: 1, 8.37809030313599: 1, 8.374996223967461: 1, 8.374570592757538: 1, 8.37429021165667: 1, 8.368324553902045: 1, 8.36425139276093: 1, 8.35304289940751: 1, 8.348894865913362: 1, 8.341855103451856: 1, 8.335137371057808: 1, 8.307242416668224: 1, 8.301590329491585: 1, 8.276361086327956: 1, 8.270413711263233: 1, 8.269481919864718: 1, 8.26536067965857: 1, 8.264347987717846: 1, 8.26322600737903: 1, 8.258035548358304: 1, 8.257416473183033: 1, 8.2550469327277: 1, 8.2543756461577: 1, 8.253211042381528: 1, 8.2493593291608: 1, 8.246610973468803: 1, 8.238817526131811: 1, 8.236419746311727: 1, 8.23504708502384: 1, 8.230882539620815: 1, 8.22735841390181: 1, 8.218323360606608: 1, 8.216805535815924: 1, 8.212532699779873: 1, 8.209877103163858: 1, 8.204775374305537: 1, 8.198955401378269: 1, 8.196392723345054: 1, 8.191596641664109: 1, 8.186318817892854: 1, 8.183233673088207: 1, 8.17776329425532: 1, 8.176507123529374: 1, 8.163115404089003: 1, 8.161086196896829: 1, 8.152744093536647: 1, 8.152605010858807: 1, 8.144504156949315: 1, 8.134252167616475: 1, 8.123846506704117: 1, 8.117915027816649: 1, 8.102867971266058: 1, 8.101377995426141: 1, 8.101035884302613: 1, 8.095423369483779: 1, 8.088117471137625: 1, 8.085653889562542: 1, 8.077540963058153: 1, 8.077415604441288: 1, 8.076646936807592: 1, 8.062914629072885: 1, 8.051156876267715: 1, 8.045043259303942: 1, 8.04428545170823: 1, 8.042405581939892: 1, 8.030870010848677: 1, 8.026539564132943: 1, 8.025664690498742: 1, 8.02501346335687: 1, 8.023046656787272: 1, 8.012677791617715: 1, 8.006676372285714: 1, 8.003449766197804: 1, 7.998543550644615: 1, 7.994264948238377: 1, 7.994038095405836: 1, 7.977765734941353: 1, 7.958969055009495: 1, 7.955860953881023: 1, 7.945534728498097: 1, 7.935703773684051: 1, 7.935642906850508: 1, 7.935403178623025: 1, 7.934581892705588: 1, 7.933674400038745: 1, 7.927780983089533: 1, 7.9262365623300965: 1, 7.924885843056973: 1, 7.924824142091652: 1, 7.922612787595712: 1, 7.914693917881262: 1, 7.897339482790411: 1, 7.894867227445524: 1, 7.885143773232172: 1, 7.884398621826785: 1, 7.879042037195777: 1, 7.87466601996211: 1, 7.870771321804749: 1, 7.869092190287327: 1, 7.866753910424794: 1, 7.860153493225763: 1, 7.856756928807752: 1, 7.853411739493612: 1, 7.8518696806552715: 1, 7.8371884844920805: 1, 7.836310429130546: 1, 7.835264961315315: 1, 7.832252544938002: 1, 7.832104403848898: 1, 7.826622742331712: 1, 7.826573819190918: 1, 7.822695876670623: 1, 7.818505538666986: 1, 7.813808165882252: 1, 7.8118982482848125: 1, 7.805710326290275: 1, 7.804304730832059: 1, 7.770226072057995: 1, 7.766754009003501: 1, 7.75761038934458: 1, 7.743720055820509: 1, 7.742690384069222: 1, 7.737761818552241: 1, 7.737079046190304: 1, 7.7349667996925575: 1, 7.734509306012274: 1, 7.726708922228827: 1, 7.7258059385901054: 1, 7.71540712845422: 1, 7.7145328575148495: 1, 7.7075454631818605: 1, 7.706801279750411: 1, 7.697392068512475: 1, 7.692960638226273: 1, 7.689163649229882: 1, 7.687252634847213: 1, 7.686830487553939: 1, 7.678353544443433: 1, 7.6760169645492615: 1, 7.6641140061339295: 1, 7.660601633126618: 1, 7.656537197229839: 1, 7.656167542769154: 1, 7.636194092961562: 1, 7.633378373998138: 1, 7.628507250969874: 1, 7.613999909558325: 1, 7.608972866230454: 1, 7.608659100799882: 1, 7.599941953748322: 1, 7.591207617119702: 1, 7.591118481707564: 1, 7.59040422211946: 1, 7.589385854180043: 1, 7.586332185485165: 1, 7.580286624594976: 1, 7.571765661519544: 1, 7.571182455859327: 1, 7.570576984661961: 1, 7.563818036880219: 1, 7.5618734400535: 1, 7.550676122243403: 1, 7.54867820264301: 1, 7.5486906793936175: 1, 7.537718700555438: 1, 7.535758302412769: 1, 7.5277614118499425: 1, 7.5201601751168: 1, 7.5184098345291845: 1, 7.504723812686514: 1, 7.503201630426188: 1, 7.495959630522043: 1, 7.494907130134344: 1, 7.493214040524039: 1, 7.492874843696073: 1, 7.489052288445717: 1, 7.484732554576982: 1, 7.4792760532168225: 1, 7.477584756851264: 1, 7.470626983053536: 1, 7.468580120016989: 1, 7.466978543967108: 1, 7.463247355597147: 1, 7.457584690078454: 1, 7.456448862215589: 1, 7.451608366574558: 1, 7.451370017881523: 1, 7.447989349703241: 1, 7.445658310499931: 1, 7.442975899595548: 1, 7.437378038570513: 1, 7.435745238245782: 1, 7.433399127922853: 1, 7.431380655617274: 1, 7.428769816111975: 1, 7.426434446914093: 1, 7.423700151428918: 1, 7.398165628982593: 1, 7.396313146915249: 1, 7.393733343138338: 1, 7.391307321703137: 1, 7.388773675018825: 1, 7.386246207080015: 1, 7.382847093149148: 1, 7.3741581463878605: 1, 7.368711476455816: 1, 7.368584371216464: 1, 7.367195262020458: 1, 7.365733183568286: 1, 7.365084320528162: 1, 7.363122915547791: 1, 7.359632245676502: 1, 7.343231603477289: 1, 7.338399179081608: 1, 7.33322241624004: 1, 7.322182181053062: 1, 7.320912058002222: 1, 7.316101759323428: 1, 7.311393619352961: 1, 7.308289543185395: 1, 7.3082158547502925: 1, 7.306064834486312: 1, 7.305816797600795: 1, 7.305632694094798: 1, 7.30478429274362: 1, 7.301024325046365: 1, 7.2972379315946805: 1, 7.280243157101305: 1, 7.25817291359616: 1, 7.250449461395435: 1, 7.246284173064676: 1, 7.240983795292997: 1, 7.237942848791883: 1, 7.224062827762056: 1, 7.219434753869385: 1, 7.2154909001483: 1, 7.212526988129468: 1, 7.209304171386223: 1,

7.195653119973359: 1, 7.188771110323447: 1, 7.187735549372542: 1, 7.1869693979516205: 1, 7.177471382451386: 1, 7.172188214541874: 1, 7.16511350732615: 1, 7.162218004531782: 1, 7.160566989386235: 1, 7.145070113495332: 1, 7.144349108778011: 1, 7.144100420337684: 1, 7.136945272520762: 1, 7.134533102460122: 1, 7.122825217959747: 1, 7.113921091029487: 1, 7.113808644704322: 1, 7.111510774301491: 1, 7.111404993423469: 1, 7.108781873795506: 1, 7.105278222858686: 1, 7.094517499503336: 1, 7.089834289157117: 1, 7.085786570221349: 1, 7.073723811944147: 1, 7.069569282716057: 1, 7.05469606049467: 1, 7.0459792640855206: 1, 7.043250929750974: 1, 7.042701719855126: 1, 7.041323100996833: 1, 7.028492568309781: 1, 7.028292837296942: 1, 7.0223160038039305: 1, 7.020617769941924: 1, 7.017038414506008: 1, 7.013359582971792: 1, 7.011932110244325: 1, 6.993745098077591: 1, 6.990127727594754: 1, 6.983798298705794: 1, 6.9830154097519666: 1, 6.964211492317896: 1, 6.963632615532775: 1, 6.952727773590034: 1, 6.947423108971404: 1, 6.9393128090823: 1, 6.9378924654330225: 1, 6.9334568128262175: 1, 6.931768063121236: 1, 6.917420880549754: 1, 6.915649331839087: 1, 6.91485439 9083428: 1, 6.9123713030370935: 1, 6.911829742223225: 1, 6.907614670030558: 1, 6.906863852905828: 1, 6.905548695067743: 1, 6.903550615478182: 1, 6.8979586567896565: 1, 6.883067931165554: 1, 6.876794538567564: 1, 6.875210979615283: 1, 6.870332697806177: 1, 6.8656043526828014: 1, 6.8619347996338815: 1, 6.8554962929285015: 1, 6.847182948925575: 1, 6.845308173107574: 1, 6.843028194773864: 1, 6.839350634464762: 1, 6.837930588893324: 1, 6.8373303933126: 1, 6.836355665373743: 1, 6.833370017258292: 1, 6.826035983938348: 1, 6.824090794041156: 1, 6.823292432188462: 1, 6.821670752194509: 1, 6.818922383120029: 1, 6.818625952310405: 1, 6.815098687980464: 1, 6.803998048353619: 1, 6.803098355102456: 1, 6.791583362687245: 1, 6.7899116573950415: 1, 6.784548730191906: 1, 6.779589432177469: 1, 6.779387075254446: 1, 6.77816038 51708: 1, 6.767470173812131: 1, 6.765539643527911: 1, 6.762549426070973: 1, 6.7476755211048145: 1, 6.741482195830864: 1, 6.736993009732557: 1, 6.734113523297366: 1, 6.7338411998704135: 1, 6.716813974902041: 1, 6.71651462086752: 1, 6.70726980126582: 1, 6.703009515288239: 1, 6.701056189218819: 1, 6.694521288578774: 1, 6.694006805422901: 1, 6.688858576083831: 1, 6.68417976680613: 1, 6.68059454567234: 1, 6.675262242338963: 1, 6.6626830367889545: 1, 6.656843567524824: 1, 6.645791264488515: 1, 6.635563629892964: 1, 6.62708496949618: 1, 6.623656732642094: 1, 6.619995262011232: 1, 6.6185691175536725: 1, 6.613835510244248: 1, 6.60623139428177: 1, 6.605676555639996: 1, 6.603011292544142: 1, 6.596512258407412: 1, 6.594078800752774: 1, 6.591125138379063: 1, 6.589981593737998: 1, 6.586562470043666: 1, 6.5806085963037235: 1, 6.578859607893112: 1, 6.578492688948223: 1, 6.568792837101024: 1, 6.56822629 1293677: 1, 6.559056297668872: 1, 6.556842039554975: 1, 6.548637271016234: 1, 6.543871833565101: 1, 6.541184489778629: 1, 6.5394990279402965: 1, 6.537822578343027: 1, 6.536045465073636: 1, 6.520120831281223: 1, 6.515563587382454: 1, 6.508528851463983: 1, 6.501161586160881: 1, 6.500473050337444: 1, 6.499642595654314: 1, 6.496230104522725: 1, 6.492360312896272: 1, 6.49004421377781: 1, 6.478824425144189: 1, 6.471136372097758: 1, 6.469548644106322: 1, 6.4585338833875054: 1, 6.456115772157545: 1, 6.453298995259167: 1, 6.448604576854004: 1, 6.44853079 8319813: 1, 6.44837206803887: 1, 6.445831277075485: 1, 6.444874720208956: 1, 6.4447045843905135: 1, 6.438168077437756: 1, 6.433121750068218: 1, 6.432780373300176: 1, 6.431872201168065: 1, 6.430815728361049: 1, 6.430685335827576: 1, 6.428334812711625: 1, 6.422725593079599: 1, 6.4215721471720055: 1, 6.420956328169909: 1, 6.418893723941316: 1, 6.414600314354635: 1, 6.40494268 1368538: 1, 6.402462170958604: 1, 6.401132770476779: 1, 6.388359701479123: 1, 6.3804019185990315: 1, 6.373740388579233: 1, 6.361060136561711: 1, 6.360333364806624: 1, 6.35995548007729: 1, 6.356961582054231: 1, 6.348590728338855: 1, 6.34696724923216: 1, 6.3429607596683: 1, 6.339625868143785: 1, 6.332322521184014: 1, 6.33225630542078: 1, 6.328625114452751: 1, 6.325260758014128: 1, 6.324255823682791: 1, 6.32161763078224: 1, 6.318686821956279: 1, 6.311145448637443: 1, 6.296596256410152: 1, 6.293996671142392: 1, 6.29254387868259: 1, 6.291233640936824: 1, 6.287552191947521: 1, 6.2871701781458125: 1, 6.2857156053795515: 1, 6.279052753571865: 1, 6.275999299160987: 1, 6.272628327921942: 1, 6.265662027714928: 1, 6.262568733200062: 1, 6.261191564236765: 1, 6.257025982360714: 1, 6.255916912598764: 1, 6.255456707565467: 1, 6.245775986103721: 1, 6.239842907846048: 1, 6.231364950262502: 1, 6.230924808224252: 1, 6.229147356409668: 1, 6.225091491299142: 1, 6.223518104048368: 1, 6.214618651688099: 1, 6.2084340479487095: 1, 6.205482891693868: 1, 6.199257136971378: 1, 6.197491980588748: 1, 6.185756299907896: 1, 6.183202640884544: 1, 6.182887600283459: 1, 6.18249596499441: 1, 6.177190548997374: 1, 6.17293024603745: 1, 6.169842552391914: 1, 6.152937879512105: 1, 6.152684002241841: 1, 6.152018381039757: 1, 6.146916372466631: 1, 6.13935761736804: 1, 6.132960517539995: 1, 6.130742614698839: 1, 6.115061705536334: 1, 6.113644060209828: 1, 6.109546899396782: 1, 6.108507498638792: 1, 6.102216169097589: 1, 6.091316776808786: 1, 6.0894766060596295: 1, 6.08027869657122: 1, 6.078909209148451: 1, 6.065324128859916: 1, 6.050029546765223: 1, 6.048916040708434: 1, 6.045593158150724: 1, 6.041384217800808: 1, 6.029282550132979: 1, 6.015590781841139: 1, 6.015097051664544: 1, 6.012145544105281: 1, 6.003706445067376: 1, 6.001342144259179: 1, 6.000642762956984: 1, 5.999283674037049: 1, 5.99788812025839: 1, 5.99724469177686: 1, 5.991407973700849: 1, 5.98986096752403: 1, 5.986410135263609: 1, 5.983086990730731: 1, 5.9817520047108745: 1, 5.98077649380523: 1, 5.980024362192422: 1, 5.971947351818514: 1, 5.971017583753861: 1, 5.970520029430779: 1, 5.967977208793935: 1, 5.96426127067981: 1, 5.963068896879114: 1, 5.956585079818992: 1, 5.95460925623019: 1, 5.95458694900644: 1, 5.954236662018317: 1, 5.953953270725752: 1, 5.944267899294497: 1, 5.942825374563659: 1, 5.93813076465409: 1, 5.935407778054932: 1, 5.935024267231281: 1, 5.9313453205280835: 1, 5.922721446311886: 1, 5.922198701779073: 1, 5.9164650956252025: 1, 5.912027871535016: 1, 5.904189193287237: 1, 5.903963206482505: 1, 5.903931641736175: 1, 5.903600652245732: 1, 5.899384025321669: 1, 5.891174836943688: 1, 5.8872427192521615: 1, 5.886480278190184: 1, 5.884625427687373: 1, 5.884341282196662: 1, 5.881098602013199: 1, 5.879258707906551: 1, 5.879122132890086: 1, 5.877944666980077: 1, 5.873989570804453: 1, 5.870331414393747: 1, 5.8658785747231255: 1, 5.859256858972805: 1, 5.857130968952953: 1, 5.854351903456167: 1, 5.852273052093466: 1, 5.852093358559783: 1, 5.846552174618787: 1, 5.845935915389509: 1, 5.845607555615441: 1,

5.84188447563457: 1, 5.840566774911274: 1, 5.826507630737109: 1, 5.824149680161676: 1, 5.823399185297636: 1, 5.812087757518777: 1, 5.809933838376124: 1, 5.8022564714937115: 1, 5.8020523939930495: 1, 5.7970717038996575: 1, 5.793458415878454: 1, 5.789151239690762: 1, 5.787567715829301: 1, 5.787059620534036: 1, 5.780065742687914: 1, 5.778733479094732: 1, 5.772298568279479: 1, 5.766483953191144: 1, 5.760229302244786: 1, 5.760191867484349: 1, 5.760113707863088: 1, 5.749791163922584: 1, 5.746239310864933: 1, 5.745244803817244: 1, 5.740884960604078: 1, 5.735623277497041: 1, 5.731858338656977: 1, 5.730314146254625: 1, 5.729837628743183: 1, 5.72044654637108: 1, 5.719606282653876: 1, 5.71862460662191: 1, 5.717036040561051: 1, 5.716638433084852: 1, 5.7161944184873565: 1, 5.713446245922326: 1, 5.712361865118431: 1, 5.711928454773895: 1, 5.704215641007188: 1, 5.703918689846563: 1, 5.702269454132776: 1, 5.698122925081619: 1, 5.695772193352759: 1, 5.693397835711586: 1, 5.693245699345299: 1, 5.691712905217233: 1, 5.682455756703673: 1, 5.679933227345937: 1, 5.678974608130793: 1, 5.678823345377231: 1, 5.677708558082862: 1, 5.676129679945563: 1, 5.674180127870725: 1, 5.670723446498605: 1, 5.669337876864157: 1, 5.667214483904998: 1, 5.667130481322348: 1, 5.664368223162768: 1, 5.65753488330356: 1, 5.657519717893916: 1, 5.650086039721582: 1, 5.647039673023564: 1, 5.646988968579963: 1, 5.64313827772232: 1, 5.642893234067691: 1, 5.636635270155471: 1, 5.635767457153182: 1, 5.635149726182266: 1, 5.633992816758303: 1, 5.626372426089908: 1, 5.625643201123235: 1, 5.625566360150697: 1, 5.623877687641618: 1, 5.623304745071088: 1, 5.6228333876608385: 1, 5.617846103445427: 1, 5.60902321443792: 1, 5.593931376250802: 1, 5.573874552093727: 1, 5.573262322172215: 1, 5.573140476520437: 1, 5.5714032907097595: 1, 5.571027528710737: 1, 5.569174628419467: 1, 5.5665423841794235: 1, 5.565095277615214: 1, 5.560264010179829: 1, 5.560156411489854: 1, 5.559196420952133: 1, 5.557526002293113: 1, 5.545502686901957: 1, 5.5446318489024335: 1, 5.542476153385423: 1, 5.54210474481343: 1, 5.538431148798349: 1, 5.5371588383018215: 1, 5.534919211073624: 1, 5.534352804471436: 1, 5.533001371327809: 1, 5.531339905133814: 1, 5.522233013812014: 1, 5.520432180149745: 1, 5.520390163256907: 1, 5.515819605514254: 1, 5.509592938710696: 1, 5.508318692191006: 1, 5.507323853337976: 1, 5.499898998821875: 1, 5.4984559317987385: 1, 5.497617122735511: 1, 5.496691929239375: 1, 5.495825486371091: 1, 5.492772074616687: 1, 5.490958116872524: 1, 5.490583044589588: 1, 5.49057536846779: 1, 5.490095738792401: 1, 5.48627888533167: 1, 5.474045830703931: 1, 5.473314826799138: 1, 5.467501057992931: 1, 5.462433097033034: 1, 5.441299132395783: 1, 5.43999698588953: 1, 5.43648123503099: 1, 5.4357269410979026: 1, 5.432614075222547: 1, 5.431289674153378: 1, 5.430057556387989: 1, 5.4293324306804385: 1, 5.426896397097763: 1, 5.423010849191527: 1, 5.422947871701359: 1, 5.419329461741959: 1, 5.418659216579986: 1, 5.4151712264468355: 1, 5.410213679410685: 1, 5.405885330729384: 1, 5.405718995476058: 1, 5.400946348201658: 1, 5.398757007329241: 1, 5.390936284745154: 1, 5.39065694932394: 1, 5.387175768684874: 1, 5.385894691034405: 1, 5.373969681773698: 1, 5.369978763022761: 1, 5.369774015302776: 1, 5.367800033163312: 1, 5.356750090120179: 1, 5.353039331468642: 1, 5.341069746225395: 1, 5.337071732567945: 1, 5.3336418425238135: 1, 5.332787192701245: 1, 5.331626062086764: 1, 5.32848414209773: 1, 5.327743400390499: 1, 5.327423191734625: 1, 5.3270596418335225: 1, 5.326297133501019: 1, 5.325657995594006: 1, 5.324647083237166: 1, 5.324160606276318: 1, 5.317361328869733: 1, 5.300749175389959: 1, 5.299290246028076: 1, 5.280110938775182: 1, 5.279940073768986: 1, 5.276165551072847: 1, 5.27525370573545: 1, 5.274617284944648: 1, 5.270843403579468: 1, 5.266957987557796: 1, 5.256313522565246: 1, 5.25455184901875: 1, 5.249424656866847: 1, 5.248883378852159: 1, 5.2459141683089205: 1, 5.239787659868004: 1, 5.234797200726986: 1, 5.23341434157335: 1, 5.2300204694097605: 1, 5.229648715023363: 1, 5.228503331917899: 1, 5.224388254567361: 1, 5.221945544030882: 1, 5.212597656220152: 1, 5.206664800107267: 1, 5.202808684093221: 1, 5.2023389502257755: 1, 5.198810101385718: 1, 5.198479053394528: 1, 5.195478198788059: 1, 5.1896028810873895: 1, 5.186636861671483: 1, 5.183987423072675: 1, 5.183190708598058: 1, 5.180551950514975: 1, 5.169076948229175: 1, 5.168391984632888: 1, 5.167962881765782: 1, 5.1661461613611035: 1, 5.1649919430933: 1, 5.163642683423685: 1, 5.1636310268744205: 1, 5.159029724868976: 1, 5.155903978906707: 1, 5.1552158300625495: 1, 5.154624076886591: 1, 5.151201396790296: 1, 5.150654495871013: 1, 5.145109237456816: 1, 5.143139772223639: 1, 5.13341411627616: 1, 5.1324259737629685: 1, 5.130199123198474: 1, 5.126284556860417: 1, 5.1250859890438205: 1, 5.124418260614058: 1, 5.118117266730573: 1, 5.116606046138037: 1, 5.112434594376135: 1, 5.107363676930964: 1, 5.107079070286368: 1, 5.1049246908328545: 1, 5.102734582841756: 1, 5.102494580593238: 1, 5.101500734984361: 1, 5.100550102770743: 1, 5.099962944129771: 1, 5.098759633664765: 1, 5.098653483862214: 1, 5.090147107679556: 1, 5.089140829070071: 1, 5.086827238439607: 1, 5.0808883883207345: 1, 5.079935095694573: 1, 5.079518018342138: 1, 5.077712904446601: 1, 5.072779190709165: 1, 5.071160425592302: 1, 5.070384647418978: 1, 5.069335633741044: 1, 5.067546595289422: 1, 5.064784408318373: 1, 5.0635590008162055: 1, 5.061539293343713: 1, 5.061355011461879: 1, 5.0596175725236305: 1, 5.055167041588525: 1, 5.054337113202608: 1, 5.054277960510644: 1, 5.052899612184039: 1, 5.050562304509274: 1, 5.047806746280948: 1, 5.0303505040652015: 1, 5.026550630104197: 1, 5.025277404546997: 1, 5.023499293036635: 1, 5.023334605874647: 1, 5.023211695063804: 1, 5.02092143464149: 1, 5.015252983327759: 1, 5.01075022176069: 1, 5.006187413418719: 1, 5.004998945299513: 1, 5.0048214402912174: 1, 4.994340779566953: 1, 4.99083727355279: 1, 4.989524211511532: 1, 4.984937553461114: 1, 4.980902326661787: 1, 4.9718638910154676: 1, 4.970820183341754: 1, 4.969668774551273: 1, 4.967406787962255: 1, 4.960788867083377: 1, 4.958503361907825: 1, 4.95342070387335: 1, 4.952772264381587: 1, 4.952501461309308: 1, 4.95209196735262: 1, 4.9427597865245945: 1, 4.939190661345618: 1, 4.937124354244177: 1, 4.936305516291757: 1, 4.933994514963243: 1, 4.918549704539593: 1, 4.917984214208678: 1, 4.914463269165306: 1, 4.914059370390886: 1, 4.912127259246129: 1, 4.9095770588197425: 1, 4.905866802477527: 1, 4.900341324366064: 1, 4.897015287546051: 1, 4.892070038301602: 1, 4.871591971485748: 1, 4.867792543968598: 1, 4.864396393078039: 1, 4.852483897575169: 1, 4.850689570539308: 1, 4.846439082030693: 1, 4.839956897649465: 1, 4.839671597626838: 1, 4.8382963686773826: 1, 4.8366043288962235: 1, 4.83624941602738: 1, 4.831407385382671: 1, 4.830774035201668: 1, 4.827987348027044: 1,

4.821435600154054: 1, 4.816773715941438: 1, 4.814878162433417: 1, 4.810291208409349: 1, 4.8095007979740885: 1, 4.808855468378069: 1, 4.802194882791266: 1, 4.791785253929847: 1, 4.782044364592723: 1, 4.775509332057427: 1, 4.772564867994075: 1, 4.77149198298563: 1, 4.769903885545292: 1, 4.769666481333969: 1, 4.765841338877689: 1, 4.7640042818393065: 1, 4.762961577150079: 1, 4.748638914895689: 1, 4.747304269527383: 1, 4.7410667577633525: 1, 4.740729321101958: 1, 4.737333718590375: 1, 4.7278358970013645: 1, 4.7264335510807465: 1, 4.718003755578054: 1, 4.714119341230794: 1, 4.712530044926843: 1, 4.705105561397923: 1, 4.702754464157433: 1, 4.702332869732983: 1, 4.699683104589926: 1, 4.698846763508681: 1, 4.691135408237548: 1, 4.689087517946064: 1, 4.68416760047677: 1, 4.67009971654441: 1, 4.669515193688528: 1, 4.667471620951465: 1, 4.6669323594242895: 1, 4.6661947241444235: 1, 4.6655196825069245: 1, 4.663742271897914: 1, 4.66267072657829: 1, 4.662305002977375: 1, 4.661861987595252: 1, 4.66122517481151: 1, 4.65732687507335: 1, 4.654846746206334: 1, 4.650187611816545: 1, 4.650150751387292: 1, 4.643394803177734: 1, 4.6431110577068395: 1, 4.6415273954209955: 1, 4.636675469750966: 1, 4.635842292729796: 1, 4.635146780468495: 1, 4.631347263785195: 1, 4.621990122686766: 1, 4.6174567256014125: 1, 4.605392196663065: 1, 4.600002482085826: 1, 4.599394559509389: 1, 4.597341663593558: 1, 4.579663744864896: 1, 4.578994037180124: 1, 4.5778633085018345: 1, 4.57694056900522: 1, 4.5724167745890725: 1, 4.567572144627855: 1, 4.567273199551433: 1, 4.566160139565958: 1, 4.558820965288001: 1, 4.549249691177796: 1, 4.544523192263898: 1, 4.541095686799693: 1, 4.540712628110586: 1, 4.528043385575052: 1, 4.527961953309776: 1, 4.526078971287783: 1, 4.524661157765575: 1, 4.524447415872944: 1, 4.502693144621285: 1, 4.4988505931157325: 1, 4.492235803169127: 1, 4.483624675523727: 1, 4.481451675006587: 1, 4.481273236367334: 1, 4.479692059338242: 1, 4.478849865356813: 1, 4.47740192143835: 1, 4.475778857594541: 1, 4.4744185854538765: 1, 4.472463734037876: 1, 4.46998395792946: 1, 4.462885710783098: 1, 4.4608419491657605: 1, 4.452524231565765: 1, 4.450465811405836: 1, 4.445410354108077: 1, 4.443848180668747: 1, 4.442855668328901: 1, 4.434986803736098: 1, 4.43386156446043: 1, 4.429814733961859: 1, 4.421792642054135: 1, 4.4212901825133475: 1, 4.41847421479602: 1, 4.401611611712928: 1, 4.4012193052268715: 1, 4.399938473683351: 1, 4.399637033598226: 1, 4.398903570597419: 1, 4.396898050653742: 1, 4.396833019915102: 1, 4.395848327409134: 1, 4.392958304133408: 1, 4.39081064447686: 1, 4.39016995189121: 1, 4.385630537377727: 1, 4.383715841036224: 1, 4.374630336315908: 1, 4.372215364664976: 1, 4.371459991874058: 1, 4.370745731560836: 1, 4.365348985148766: 1, 4.358519001105127: 1, 4.355914615902342: 1, 4.346501394713578: 1, 4.345954001854997: 1, 4.340939689200948: 1, 4.33406840881724: 1, 4.330951233585537: 1, 4.32813318880459: 1, 4.32682018182838: 1, 4.3240804268877655: 1, 4.315473607985715: 1, 4.3121654024575475: 1, 4.307690825054319: 1, 4.306926123119908: 1, 4.3027655407489025: 1, 4.290836732813504: 1, 4.290491937261216: 1, 4.2842317152847995: 1, 4.282412078333787: 1, 4.281289235266931: 1, 4.270962929179206: 1, 4.256623844559316: 1, 4.255196757850798: 1, 4.245327912538191: 1, 4.235267555890625: 1, 4.22318547445586: 1, 4.209316264277829: 1, 4.202741118986923: 1, 4.202541839470851: 1, 4.200270615041777: 1, 4.197084630639547: 1, 4.193248599170173: 1, 4.188681601791211: 1, 4.18372987865737: 1, 4.18371977299696: 1, 4.182441658909898: 1, 4.172915180987755: 1, 4.172159864454462: 1, 4.169356261236256: 1, 4.169318873724207: 1, 4.165076365003772: 1, 4.164947067006613: 1, 4.147238190240023: 1, 4.146785618646595: 1, 4.14516543099556: 1, 4.141298819897186: 1, 4.1329706363712395: 1, 4.126063016327842: 1, 4.125261186481366: 1, 4.111899587718293: 1, 4.080676833643397: 1, 4.080373031304546: 1, 4.080216145825188: 1, 4.075851170590296: 1, 4.070361747541808: 1, 4.066416364723072: 1, 4.0471869873007265: 1, 4.044531302735939: 1, 4.039979419513664: 1, 4.039207035153472: 1, 4.038237982342366: 1, 4.0289754273929645: 1, 4.026217583758714: 1, 4.013788346661682: 1, 4.006606993210097: 1, 4.002324666613009: 1, 3.99705385862014: 1, 3.992124809091589: 1, 3.991301410025376: 1, 3.936203253710087: 1, 3.9060426159707693: 1, 3.8843400419887772: 1, 3.881400177613667: 1, 3.867181084660802: 1, 3.8342641565346005: 1, 3.8183871530880853: 1, 3.679591064543468: 1, 3.6605740227682784: 1, 3.6218019265279917: 1, 3.598734753115727: 1, 3.0003010393288387: 1}})

In [187]:

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

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

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

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

```

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

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

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

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

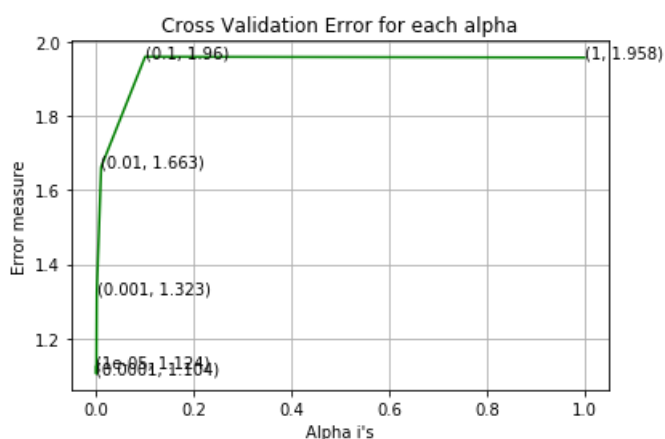
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best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

```

```

For values of alpha = 1e-05 The log loss is: 1.123612945244866
For values of alpha = 0.0001 The log loss is: 1.104219261562574
For values of alpha = 0.001 The log loss is: 1.3232031883263133
For values of alpha = 0.01 The log loss is: 1.6626205135664829
For values of alpha = 0.1 The log loss is: 1.9599373914282405
For values of alpha = 1 The log loss is: 1.9577474000242514

```



```

For values of best alpha = 0.0001 The train log loss is: 0.7221439195449428
For values of best alpha = 0.0001 The cross validation log loss is: 1.104219261562574
For values of best alpha = 0.0001 The test log loss is: 1.227753379264546

```

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

Ans. Yes, it seems like!

In [188]:

```
def get_intersec_text(df):
    df_text_vec = TfidfVectorizer(ngram_range=(1, 5),min_df=5,max_features=2000)
    df_text_fea = df_text_vec.fit_transform(df['TEXT'])
    df_text_features = df_text_vec.get_feature_names()

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

In [189]:

```
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")
```

92.55 % of word of test data appeared in train data
92.6 % of word of Cross Validation appeared in train data

3.3. Feature engineering

In [259]:

```
def getTextLength(textdata):
    length = []
    for text in textdata:
        length.append(len(text))

    return length
```

3.3.1 Take text length as new feature

In [260]:

```
train_text_len = getTextLength(train_df['TEXT'])
cv_text_len = getTextLength(cv_df['TEXT'])
test_text_len = getTextLength(test_df['TEXT'])
```

In [261]:

```
# reshape
train_text_len = np.reshape(train_text_len,(len(train_text_len),1))
test_text_len = np.reshape(test_text_len,(len(test_text_len),1))
cv_text_len = np.reshape(cv_text_len,(len(cv_text_len),1))
```

In [262]:

```
# Normalization
from sklearn.preprocessing import normalize
train_text_len = normalize(train_text_len,axis=0)
test_text_len = normalize(test_text_len,axis=0)
cv_text_len = normalize(cv_text_len,axis=0)
```

In [263]:

```
print(train_text_len.shape)
print(cv_text_len.shape)
print(test_text_len.shape)
```

(2124, 1)
(532, 1)
(665, 1)

3.3.2 Gene Length

In [264]:

```
train_gene_len = getTextLength(train_df['Gene'])
cv_gene_len = getTextLength(cv_df['Gene'])
test_gene_len = getTextLength(test_df['Gene'])
```

In [265]:

```
# reshape
train_gene_len = np.reshape(train_gene_len, (len(train_gene_len), 1))
test_gene_len = np.reshape(test_gene_len, (len(test_gene_len), 1))
cv_gene_len = np.reshape(cv_gene_len, (len(cv_gene_len), 1))
```

In [266]:

```
# Normalization
train_gene_len = normalize(train_gene_len, axis=0)
test_gene_len = normalize(test_gene_len, axis=0)
cv_gene_len = normalize(cv_gene_len, axis=0)
```

In [267]:

```
print(train_gene_len.shape)
print(cv_gene_len.shape)
print(test_gene_len.shape)
```

```
(2124, 1)
(532, 1)
(665, 1)
```

3.3.3. Variation Length

In [268]:

```
train_variation_len = getTextLength(train_df['Variation'])
cv_variation_len = getTextLength(cv_df['Variation'])
test_variation_len = getTextLength(test_df['Variation'])
```

In [269]:

```
# reshape
train_variation_len = np.reshape(train_variation_len, (len(train_variation_len), 1))
test_variation_len = np.reshape(test_variation_len, (len(test_variation_len), 1))
cv_variation_len = np.reshape(cv_variation_len, (len(cv_variation_len), 1))
```

In [270]:

```
# Normalization
from sklearn.preprocessing import normalize
train_variation_len = normalize(train_variation_len, axis=0)
test_variation_len = normalize(test_variation_len, axis=0)
cv_variation_len = normalize(cv_variation_len, axis=0)
```

In [271]:

```
print(train_variation_len.shape)
print(cv_variation_len.shape)
print(test_variation_len.shape)
```

```
(2124, 1)
(532, 1)
(665, 1)
```

4. Machine Learning Models

In [203]:

```
#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 [204]:

```
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 [205]:

```
# this function will be used just for naive bayes
# for the given indices, we will print the name of the features
# and we will check whether the feature present in the test point text or not
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = TfidfVectorizer(ngram_range=(1, 5), min_df=5, max_features=2000)
    var_count_vec = TfidfVectorizer(ngram_range=(1, 5), min_df=5, max_features=2000)
    text_count_vec = TfidfVectorizer(ngram_range=(1, 5), min_df=5, max_features=2000)

    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())
    fea3_len = len(text_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))
        elif (v < fea1_len + fea2_len + fea3_len):
            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")

def get_impfeature_names_LR(indices, text, gene, var, no_features):
```



```

gene_count_vec = CountVectorizer(ngram_range=(1, 5),min_df=5,max_features=2000)
var_count_vec = CountVectorizer(ngram_range=(1, 5),min_df=5,max_features=2000)
text_count_vec = CountVectorizer(ngram_range=(1, 5),min_df=5,max_features=2000)

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())
fea3_len = len(text_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_r
o))
    elif(v<fea1_len+fea2_len+fea3_len):
        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 [206]:

```

# merging gene, variance and text features

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

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

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

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

cv_x_onehotCoding = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding))
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 [207]:

```
# Stack Text length
train_x_onehotCoding = hstack((train_x_onehotCoding,train_text_len))
test_x_onehotCoding = hstack((test_x_onehotCoding,test_text_len))
cv_x_onehotCoding = hstack((cv_x_onehotCoding,cv_text_len))
```

In [208]:

```
# Stack Gene length
train_x_onehotCoding = hstack((train_x_onehotCoding,train_gene_len))
test_x_onehotCoding = hstack((test_x_onehotCoding,test_gene_len))
cv_x_onehotCoding = hstack((cv_x_onehotCoding,cv_gene_len))
```

In [209]:

```
# Stack Variation length
train_x_onehotCoding = hstack((train_x_onehotCoding,train_variation_len)).tocsr()
test_x_onehotCoding = hstack((test_x_onehotCoding,test_variation_len)).tocsr()
cv_x_onehotCoding = hstack((cv_x_onehotCoding,cv_variation_len)).tocsr()
```

In [210]:

```
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 = (2257, 2116)
(number of data points * number of features) in test data = (499, 2116)
(number of data points * number of features) in cross validation data = (565, 2116)
```

In [211]:

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

```
Response encoding features :
(number of data points * number of features) in train data = (2257, 27)
(number of data points * number of features) in test data = (499, 27)
(number of data points * number of features) in cross validation data = (565, 27)
```

4.1. Base Line Model

4.1.1. Naive Bayes

4.1.1.1. Hyper parameter tuning

In [75]:

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

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

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

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

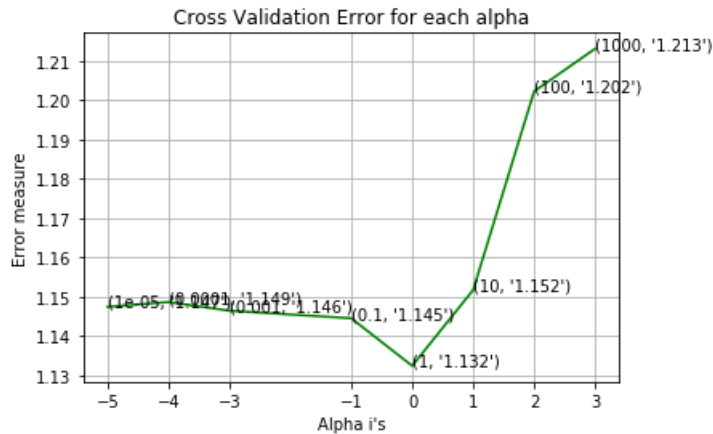
for alpha = 1e-05
Log Loss : 1.1474366145627115
for alpha = 0.0001
Log Loss : 1.148641360129883

```

```

Log Loss : 1.1464474057225402
for alpha = 0.001
Log Loss : 1.1445225217205688
for alpha = 1
Log Loss : 1.132406099663628
for alpha = 10
Log Loss : 1.1517320333344243
for alpha = 100
Log Loss : 1.2022759527726867
for alpha = 1000
Log Loss : 1.213134467353376

```



For values of best alpha = 1 The train log loss is: 1.0735679457662106
 For values of best alpha = 1 The cross validation log loss is: 1.132406099663628
 For values of best alpha = 1 The test log loss is: 1.2197509929678902

4.1.1.2. Testing the model with best hyper paramters

In [76]:

```

# 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 : 1.1495018441215494

Number of missclassified point : 0.37593984962406013

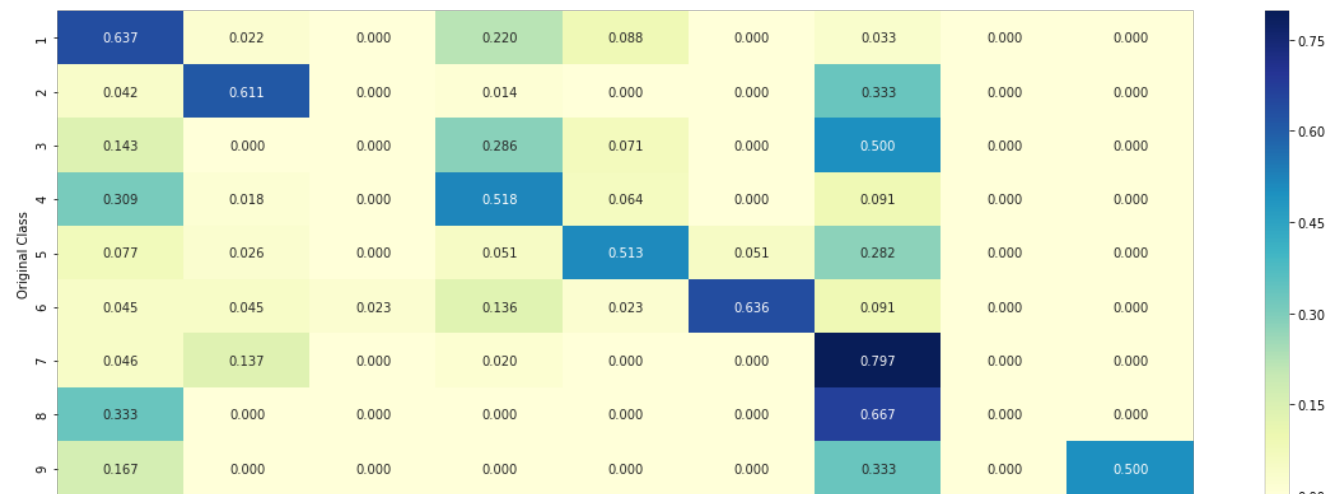
----- Confusion matrix -----



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



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



4.1.1.3. Feature Importance, Correctly classified point

In [77]:

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

Predicted Class : 4

Predicted Class Probabilities: [[0.0559 0.0516 0.0193 0.6884 0.0411 0.0384 0.0948 0.0053 0.0052]]

Actual Class : 1

```
-----
14 Text feature [activity] present in test data point [True]
15 Text feature [protein] present in test data point [True]
16 Text feature [proteins] present in test data point [True]
17 Text feature [pten] present in test data point [True]
19 Text feature [function] present in test data point [True]
20 Text feature [experiments] present in test data point [True]
24 Text feature [whereas] present in test data point [True]
26 Text feature [results] present in test data point [True]
27 Text feature [functional] present in test data point [True]
28 Text feature [indicated] present in test data point [True]
29 Text feature [shown] present in test data point [True]
30 Text feature [determined] present in test data point [True]
31 Text feature [type] present in test data point [True]
33 Text feature [acid] present in test data point [True]
34 Text feature [described] present in test data point [True]
35 Text feature [important] present in test data point [True]
36 Text feature [amino] present in test data point [True]
39 Text feature [may] present in test data point [True]
40 Text feature [ability] present in test data point [True]
41 Text feature [bind] present in test data point [True]
42 Text feature [wild] present in test data point [True]
43 Text feature [also] present in test data point [True]
44 Text feature [related] present in test data point [True]
45 Text feature [two] present in test data point [True]
46 Text feature [whether] present in test data point [True]
47 Text feature [mutations] present in test data point [True]
49 Text feature [vitro] present in test data point [True]
50 Text feature [reduced] present in test data point [True]
52 Text feature [catalytic] present in test data point [True]
53 Text feature [tagged] present in test data point [True]
58 Text feature [levels] present in test data point [True]
60 Text feature [three] present in test data point [True]
61 Text feature [thus] present in test data point [True]
62 Text feature [retained] present in test data point [True]
63 Text feature [either] present in test data point [True]
64 Text feature [30] present in test data point [True]
65 Text feature [associated] present in test data point [True]
66 Text feature [indicate] present in test data point [True]
67 Text feature [discussion] present in test data point [True]
68 Text feature [containing] present in test data point [True]
69 Text feature [lower] present in test data point [True]
70 Text feature [purified] present in test data point [True]
71 Text feature [effects] present in test data point [True]
72 Text feature [previously] present in test data point [True]
74 Text feature [although] present in test data point [True]
75 Text feature [phosphatase] present in test data point [True]
76 Text feature [functions] present in test data point [True]
77 Text feature [effect] present in test data point [True]
78 Text feature [transfected] present in test data point [True]
79 Text feature [introduction] present in test data point [True]
80 Text feature [about] present in test data point [True]
```

```

80 Text feature [show] present in test data point [True]
81 Text feature [determine] present in test data point [True]
82 Text feature [critical] present in test data point [True]
83 Text feature [suggest] present in test data point [True]
85 Text feature [therefore] present in test data point [True]
86 Text feature [transfection] present in test data point [True]
87 Text feature [see] present in test data point [True]
88 Text feature [buffer] present in test data point [True]
89 Text feature [stability] present in test data point [True]
90 Text feature [terminal] present in test data point [True]
92 Text feature [used] present in test data point [True]
93 Text feature [assay] present in test data point [True]
94 Text feature [using] present in test data point [True]
95 Text feature [however] present in test data point [True]
96 Text feature [analysis] present in test data point [True]
97 Text feature [loss] present in test data point [True]
Out of the top 100 features 66 are present in query point

```

4.1.1.4. Feature Importance, Incorrectly classified point

In [78]:

```

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

```

Predicted Class : 7

Predicted Class Probabilities: [[0.0507 0.0582 0.0143 0.0657 0.0446 0.0417 0.7141 0.0056 0.0051]]

Actual Class : 7

```

-----
14 Text feature [activation] present in test data point [True]
19 Text feature [activated] present in test data point [True]
20 Text feature [kinase] present in test data point [True]
21 Text feature [cells] present in test data point [True]
23 Text feature [downstream] present in test data point [True]
24 Text feature [expressing] present in test data point [True]
26 Text feature [also] present in test data point [True]
27 Text feature [contrast] present in test data point [True]
29 Text feature [factor] present in test data point [True]
30 Text feature [cell] present in test data point [True]
31 Text feature [however] present in test data point [True]
32 Text feature [compared] present in test data point [True]
33 Text feature [similar] present in test data point [True]
35 Text feature [independent] present in test data point [True]
36 Text feature [increased] present in test data point [True]
37 Text feature [10] present in test data point [True]
38 Text feature [shown] present in test data point [True]
39 Text feature [signaling] present in test data point [True]
40 Text feature [addition] present in test data point [True]
41 Text feature [growth] present in test data point [True]
42 Text feature [well] present in test data point [True]
43 Text feature [previously] present in test data point [True]
44 Text feature [mutations] present in test data point [True]
48 Text feature [suggest] present in test data point [True]
52 Text feature [found] present in test data point [True]
54 Text feature [may] present in test data point [True]
56 Text feature [phosphorylation] present in test data point [True]
57 Text feature [presence] present in test data point [True]
60 Text feature [recently] present in test data point [True]
62 Text feature [potential] present in test data point [True]
65 Text feature [showed] present in test data point [True]
66 Text feature [pathways] present in test data point [True]
68 Text feature [inhibition] present in test data point [True]
69 Text feature [increase] present in test data point [True]
71 Text feature [although] present in test data point [True]
74 Text feature [mechanism] present in test data point [True]

```

```

74 Text feature [mechanism] present in test data point [True]
75 Text feature [serum] present in test data point [True]
77 Text feature [mutation] present in test data point [True]
78 Text feature [12] present in test data point [True]
80 Text feature [observed] present in test data point [True]
81 Text feature [without] present in test data point [True]
82 Text feature [results] present in test data point [True]
83 Text feature [fig] present in test data point [True]
85 Text feature [total] present in test data point [True]
86 Text feature [described] present in test data point [True]
87 Text feature [studies] present in test data point [True]
88 Text feature [interestingly] present in test data point [True]
90 Text feature [oncogenic] present in test data point [True]
91 Text feature [mutant] present in test data point [True]
92 Text feature [using] present in test data point [True]
95 Text feature [either] present in test data point [True]
96 Text feature [different] present in test data point [True]
97 Text feature [including] present in test data point [True]
98 Text feature [examined] present in test data point [True]
Out of the top 100 features 54 are present in query point

```

4.2. K Nearest Neighbour Classification

4.2.1. Hyper parameter tuning

In [232]:

```

# 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 = [2,3,4,5, 11,12,13,14,15, 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 probabilties we use log-probability estimates
    print("Log Loss :",log_loss(cv_y, sig_clf_probs))

fig_ax = plt.subplots()

```



```

fig, ax = plt.subplots(1)
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],str(txt)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

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

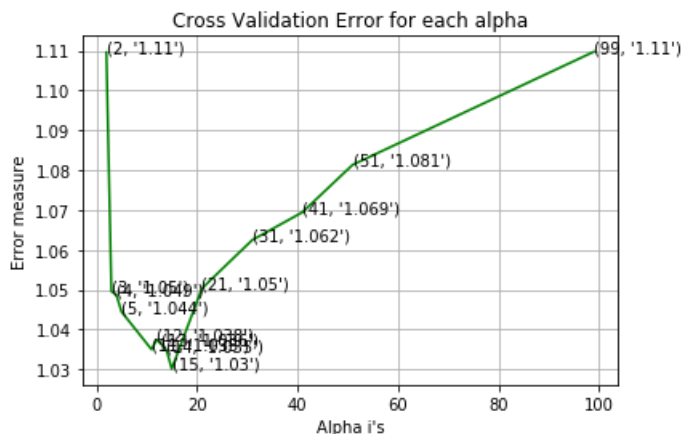
predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for alpha = 2
Log Loss : 1.1095710974165804
for alpha = 3
Log Loss : 1.0495845344930375
for alpha = 4
Log Loss : 1.0486939859397164
for alpha = 5
Log Loss : 1.0444886732604612
for alpha = 11
Log Loss : 1.0349895993225595
for alpha = 12
Log Loss : 1.0375090484316702
for alpha = 13
Log Loss : 1.036260293489855
for alpha = 14
Log Loss : 1.0347179542108307
for alpha = 15
Log Loss : 1.0301496794535479
for alpha = 21
Log Loss : 1.0503188634887977
for alpha = 31
Log Loss : 1.062459112077698
for alpha = 41
Log Loss : 1.0694396671982722
for alpha = 51
Log Loss : 1.0813570992508348
for alpha = 99
Log Loss : 1.1097205709637812

```



```

For values of best alpha = 15 The train log loss is: 0.6891548154821295
For values of best alpha = 15 The cross validation log loss is: 1.0301496794535479
For values of best alpha = 15 The test log loss is: 1.0401504132222222

```

For values of best alpha = 15 The test log loss is: 1.049150413303032

4.2.2. Testing the model with best hyper paramters

In [80]:

```
# 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)
```

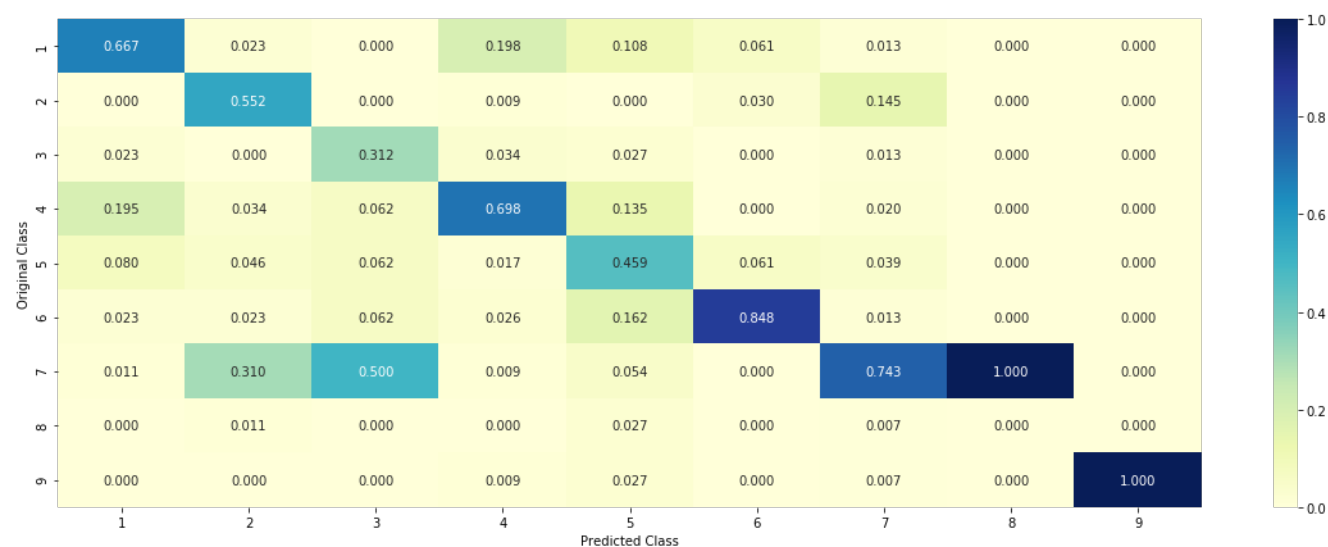
Log loss : 0.9667620353376408

Number of mis-classified points : 0.33646616541353386

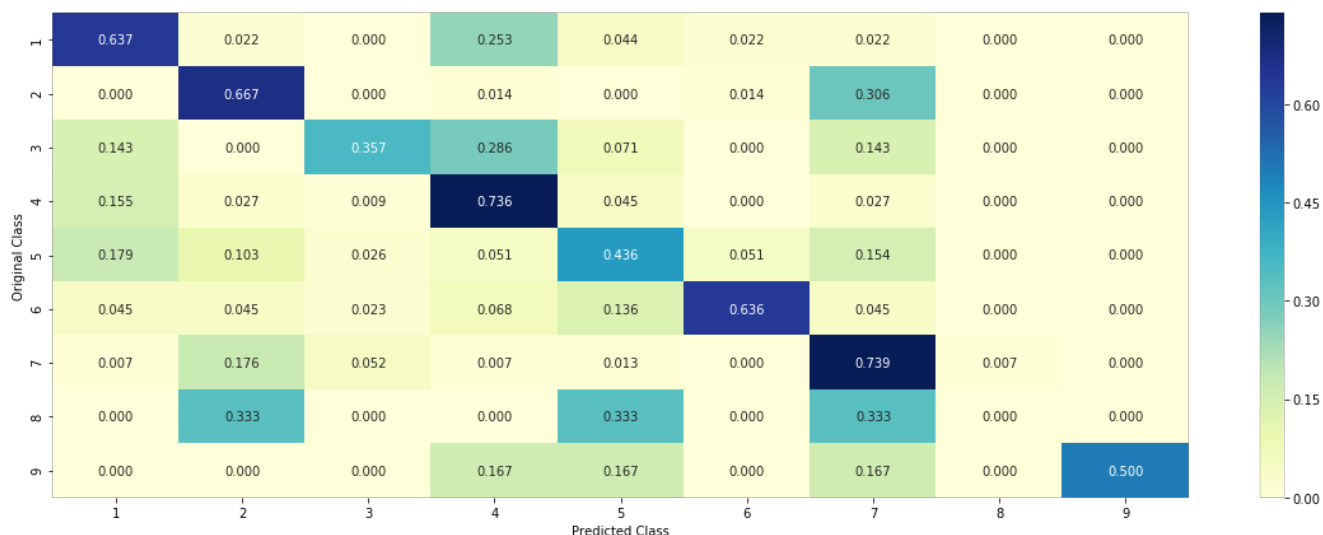
----- Confusion matrix -----



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



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



4.2.3. Sample Query point -1

In [81]:

```
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 : 2

Actual Class : 1

The 13 nearest neighbours of the test points belongs to classes [4 4 4 1 1 4 4 4 4 4 4 4 4]

Frequency of nearest points : Counter({4: 11, 1: 2})

4.2.4. Sample Query Point-2

In [82]:

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

test_point_index = 100

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

Predicted Class : 7

Actual Class : 7

the k value for knn is 13 and the nearest neighbours of the test points belongs to classes [7 7 7 7 7 7 7 7 7 7 7 7 7]

```
/ / / 2 / / / / 2]
Fequency of nearest points : Counter({7: 11, 2: 2})
```

4.3. Logistic Regression

In [254]:

```
# one-hot encoding of Gene feature.
gene_vectorizer_LR = CountVectorizer(ngram_range=(1, 4),min_df=4,max_features=1500)
train_gene_feature_onehotCoding_LR = gene_vectorizer_LR.fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding_LR = gene_vectorizer_LR.transform(test_df['Gene'])
cv_gene_feature_onehotCoding_LR = gene_vectorizer_LR.transform(cv_df['Gene'])
```

In [255]:

```
# one-hot encoding of variation feature.
variation_vectorizer_LR = CountVectorizer(ngram_range=(1, 4),min_df=4,max_features=1500)
train_variation_feature_onehotCoding_LR = variation_vectorizer_LR.fit_transform(train_df['Variation'])
test_variation_feature_onehotCoding_LR = variation_vectorizer_LR.transform(test_df['Variation'])
cv_variation_feature_onehotCoding_LR = variation_vectorizer_LR.transform(cv_df['Variation'])
```

In [256]:

```
# building a CountVectorizer with all the words that occurred minimum 3 times in train data
text_vectorizer_LR = CountVectorizer(ngram_range=(1, 4),min_df=4,max_features=1500)
train_text_feature_onehotCoding_LR = text_vectorizer_LR.fit_transform(train_df['TEXT'])
# don't forget to normalize every feature
train_text_feature_onehotCoding_LR = normalize(train_text_feature_onehotCoding_LR, axis=0)

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

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

# getting all the feature names (words)
train_text_features_LR= text_vectorizer_LR.get_feature_names()
```

In [257]:

```
# Stacking all 3 features

train_gene_var_onehotCoding =
hstack((train_gene_feature_onehotCoding_LR,train_variation_feature_onehotCoding_LR))
test_gene_var_onehotCoding =
hstack((test_gene_feature_onehotCoding_LR,test_variation_feature_onehotCoding_LR))
cv_gene_var_onehotCoding =
hstack((cv_gene_feature_onehotCoding_LR,cv_variation_feature_onehotCoding_LR))

train_x_onehotCoding_LR = hstack((train_gene_var_onehotCoding, train_text_feature_onehotCoding_LR))
test_x_onehotCoding_LR = hstack((test_gene_var_onehotCoding, test_text_feature_onehotCoding_LR))
cv_x_onehotCoding_LR = hstack((cv_gene_var_onehotCoding, cv_text_feature_onehotCoding_LR))
```

In [272]:

```
# Stack Text length
train_x_onehotCoding_LR = hstack((train_x_onehotCoding_LR,train_text_len))
test_x_onehotCoding_LR = hstack((test_x_onehotCoding_LR,test_text_len))
cv_x_onehotCoding_LR = hstack((cv_x_onehotCoding_LR,cv_text_len))
```

In [273]:

```
# Stack Gene length
train_x_onehotCoding_LR = hstack((train_x_onehotCoding_LR,train_gene_len))
```

```
test_x_onehotCoding_LR = hstack((test_x_onehotCoding_LR, test_gene_len))
cv_x_onehotCoding_LR = hstack((cv_x_onehotCoding_LR, cv_gene_len))
```

In [274]:

```
# Stack Variation length
train_x_onehotCoding_LR = hstack((train_x_onehotCoding_LR, train_variation_len)).tocsr()
test_x_onehotCoding_LR = hstack((test_x_onehotCoding_LR, test_variation_len)).tocsr()
cv_x_onehotCoding_LR = hstack((cv_x_onehotCoding_LR, cv_variation_len)).tocsr()
```

In [275]:

```
print(train_x_onehotCoding_LR.shape)
print(test_x_onehotCoding_LR.shape)
print(cv_x_onehotCoding_LR.shape)
```

```
(2124, 1633)
(665, 1633)
(532, 1633)
```

In [277]:

```
train_y = np.array(list(train_df['Class']))
cv_y = np.array(list(cv_df['Class']))
test_y = np.array(list(test_df['Class']))
```

4.3.1. With Class balancing

4.3.1.1. Hyper parameter tuning

In [278]:

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

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

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

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

alpha = [10 ** x for x in range(-8, 3)]
cv_log_error_array = []
for i in alpha:
```

```

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

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

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

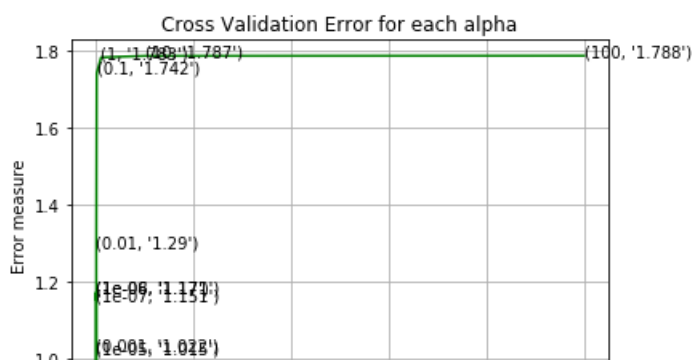
predict_y = sig_clf.predict_proba(train_x_onehotCoding_LR)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding_LR)
print('For values of best alpha = ', alpha[best_alpha],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

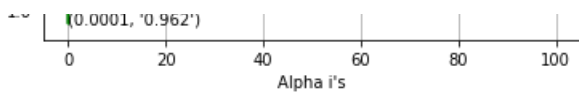
```

```

for alpha = 1e-08
Log Loss : 1.1700231813366946
for alpha = 1e-07
Log Loss : 1.1509879511769372
for alpha = 1e-06
Log Loss : 1.170522922706586
for alpha = 1e-05
Log Loss : 1.0147744547347068
for alpha = 0.0001
Log Loss : 0.9622898946365958
for alpha = 0.001
Log Loss : 1.0219106694444027
for alpha = 0.01
Log Loss : 1.2899484887301387
for alpha = 0.1
Log Loss : 1.742398763582952
for alpha = 1
Log Loss : 1.7832622173155708
for alpha = 10
Log Loss : 1.7871339048521875
for alpha = 100
Log Loss : 1.7875402989565086

```





For values of best alpha = 0.0001 The train log loss is: 0.7187138401196952
 For values of best alpha = 0.0001 The cross validation log loss is: 0.9622898946365958
 For values of best alpha = 0.0001 The test log loss is: 0.9425295409189242

4.3.1.2. Testing the model with best hyper paramters

In [279]:

```
# 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)

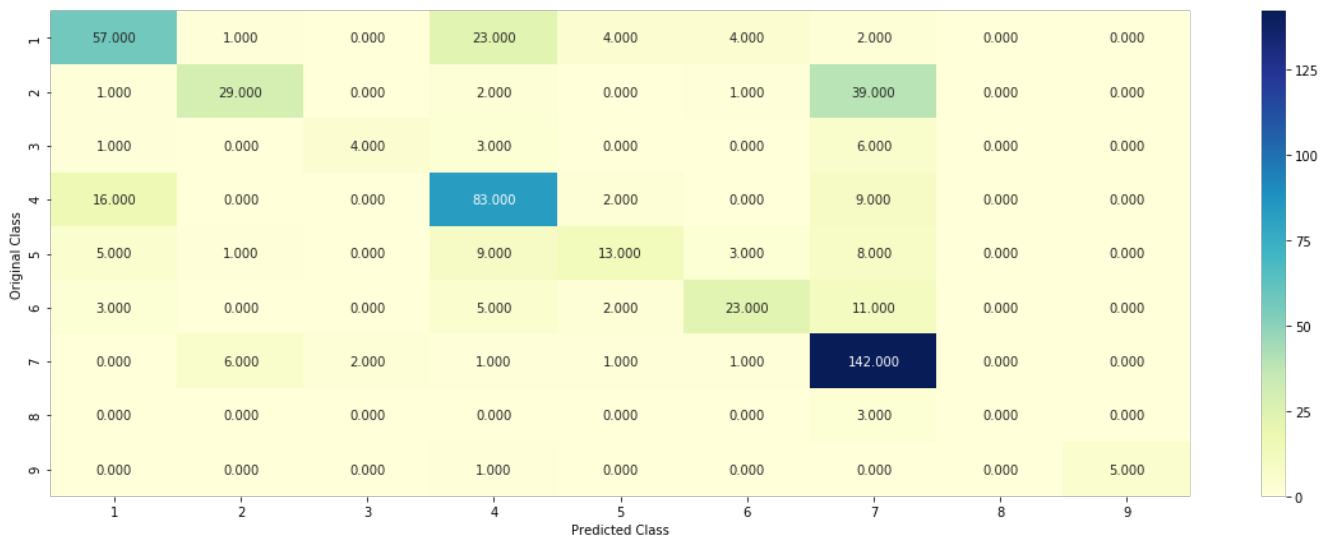
# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# 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_LR, train_y, cv_x_onehotCoding_LR, cv_y, clf)
```

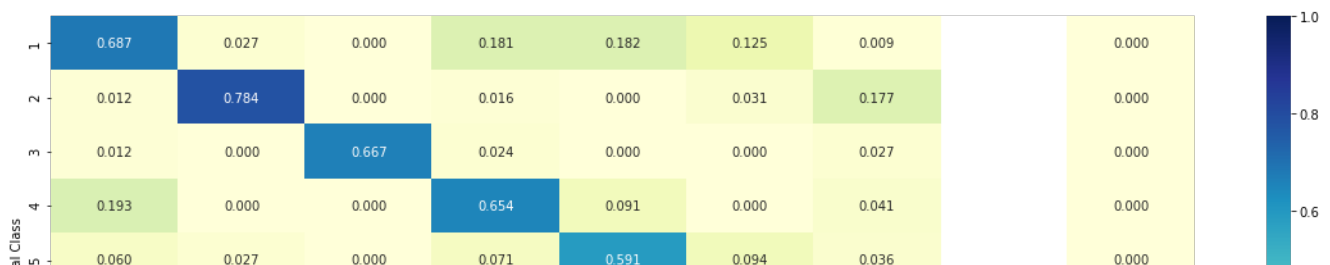
Log loss : 0.9622898946365958

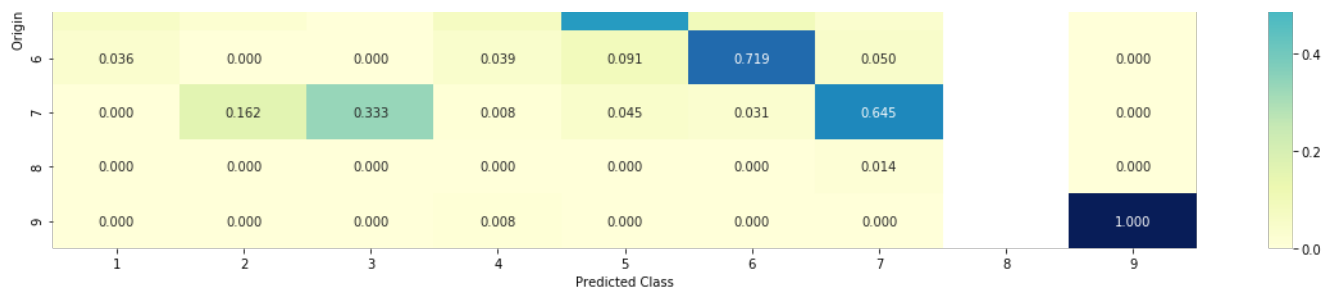
Number of mis-classified points : 0.3308270676691729

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

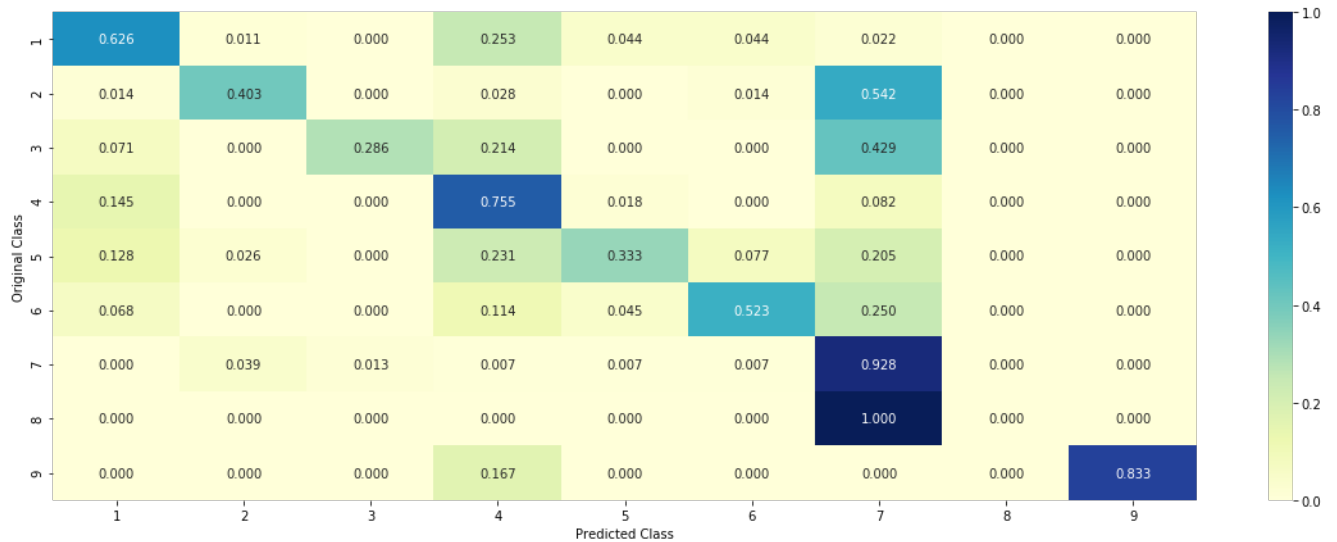


----- Precision matrix (Column Sum=1) -----





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



4.3.1.3. Feature Importance

In [93]:

```
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"]))

def get_imp_feature_names_LR(text, indices, removed_ind = []):
    word_present = 0
    tabulte_list = []
    incresingorder_ind = 0
    for i in indices:
        if i < train_gene_feature_onehotCoding_LR.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_LR[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_LR[i], yes_no])
        incresingorder_ind += 1
    print(word_present, "most important features are present in our query point")
    print("-"*50)
    print("The features that are most important of the ",predicted_cls[0]," class:")
    print (tabulate(tabulte_list, headers=["Index",'Feature name', 'Present or Not']))

```

4.3.1.3.1. Correctly Classified point

In [94]:

```

# from tabulate import tabulate
clf = SGDCClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', ran
dom_state=42)
clf.fit(train_x_onehotCoding_LR,train_y)
test_point_index = 1
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding_LR[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding_LR[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_) [predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names_LR(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: [[1.690e-02 4.900e-03 2.000e-04 8.073e-01 1.800e-03 3.000e-04 1.635e-01

4.700e-03 5.000e-04]]

Actual Class : 1

```

-----
16 Text feature [beads] present in test data point [True]
17 Text feature [material] present in test data point [True]
23 Text feature [iii] present in test data point [True]
32 Text feature [defect] present in test data point [True]
33 Text feature [washed] present in test data point [True]
34 Text feature [nucleus] present in test data point [True]
35 Text feature [motifs] present in test data point [True]
37 Text feature [tagged] present in test data point [True]
38 Text feature [loss] present in test data point [True]
40 Text feature [2a] present in test data point [True]
46 Text feature [sigma] present in test data point [True]
47 Text feature [endogenous] present in test data point [True]
48 Text feature [encodes] present in test data point [True]
49 Text feature [pcdna3] present in test data point [True]
54 Text feature [recombination] present in test data point [True]
58 Text feature [show] present in test data point [True]
69 Text feature [pbs] present in test data point [True]
76 Text feature [experimental] present in test data point [True]
77 Text feature [homozygous] present in test data point [True]
79 Text feature [changes] present in test data point [True]
80 Text feature [sds] present in test data point [True]
81 Text feature [lacking] present in test data point [True]
82 Text feature [transfected] present in test data point [True]
83 Text feature [functionally] present in test data point [True]
84 Text feature [nacl] present in test data point [True]
86 Text feature [deleted] present in test data point [True]
87 Text feature [density] present in test data point [True]
88 Text feature [phosphatase] present in test data point [True]
90 Text feature [mm] present in test data point [True]
91 Text feature [lysis] present in test data point [True]
93 Text feature [online] present in test data point [True]
94 Text feature [contribute] present in test data point [True]
100 Text feature [linker] present in test data point [True]
105 Text feature [phenotype] present in test data point [True]
109 Text feature [s4] present in test data point [True]
110 Text feature [heterozygous] present in test data point [True]
112 Text feature [domains] present in test data point [True]
114 Text feature [see] present in test data point [True]
118 Text feature [cannot] present in test data point [True]

```

119 Text feature [biotechnology] present in test data point [True]
123 Text feature [1998] present in test data point [True]
125 Text feature [procedures] present in test data point [True]
126 Text feature [representative] present in test data point [True]
129 Text feature [suggesting] present in test data point [True]
130 Text feature [risk] present in test data point [True]
131 Text feature [plates] present in test data point [True]
132 Text feature [assay] present in test data point [True]
133 Text feature [modified] present in test data point [True]
134 Text feature [29] present in test data point [True]
135 Text feature [min] present in test data point [True]
136 Text feature [assays] present in test data point [True]
137 Text feature [due] present in test data point [True]
138 Text feature [cohort] present in test data point [True]
143 Text feature [bound] present in test data point [True]
147 Text feature [express] present in test data point [True]
149 Text feature [stability] present in test data point [True]
150 Text feature [2010] present in test data point [True]
154 Text feature [rhoa] present in test data point [True]
155 Text feature [defective] present in test data point [True]
157 Text feature [motif] present in test data point [True]
158 Text feature [depletion] present in test data point [True]
159 Text feature [isoform] present in test data point [True]
160 Text feature [genetic] present in test data point [True]
162 Text feature [detailed] present in test data point [True]
163 Text feature [regulates] present in test data point [True]
166 Text feature [functional] present in test data point [True]
169 Text feature [times] present in test data point [True]
171 Text feature [2b] present in test data point [True]
172 Text feature [knockdown] present in test data point [True]
173 Text feature [endometrial] present in test data point [True]
174 Text feature [lower] present in test data point [True]
176 Text feature [right] present in test data point [True]
177 Text feature [supporting] present in test data point [True]
178 Text feature [primarily] present in test data point [True]
180 Text feature [product] present in test data point [True]
182 Text feature [co] present in test data point [True]
183 Text feature [database] present in test data point [True]
184 Text feature [elevated] present in test data point [True]
185 Text feature [low] present in test data point [True]
186 Text feature [experiment] present in test data point [True]
190 Text feature [tris] present in test data point [True]
192 Text feature [induced] present in test data point [True]
194 Text feature [observation] present in test data point [True]
195 Text feature [santa] present in test data point [True]
196 Text feature [culture] present in test data point [True]
197 Text feature [bars] present in test data point [True]
198 Text feature [s2] present in test data point [True]
202 Text feature [fetal] present in test data point [True]
203 Text feature [mice] present in test data point [True]
204 Text feature [buffer] present in test data point [True]
206 Text feature [thus] present in test data point [True]
207 Text feature [localized] present in test data point [True]
210 Text feature [protein] present in test data point [True]
211 Text feature [tested] present in test data point [True]
215 Text feature [incidence] present in test data point [True]
217 Text feature [proportion] present in test data point [True]
220 Text feature [conserved] present in test data point [True]
222 Text feature [western] present in test data point [True]
223 Text feature [strongly] present in test data point [True]
225 Text feature [differences] present in test data point [True]
226 Text feature [p85] present in test data point [True]
227 Text feature [indicate] present in test data point [True]
228 Text feature [rabbit] present in test data point [True]
230 Text feature [nuclear] present in test data point [True]
236 Text feature [cruz] present in test data point [True]
237 Text feature [comparison] present in test data point [True]
238 Text feature [poor] present in test data point [True]
239 Text feature [light] present in test data point [True]
240 Text feature [affected] present in test data point [True]
242 Text feature [caused] present in test data point [True]
243 Text feature [examine] present in test data point [True]
244 Text feature [investigated] present in test data point [True]
245 Text feature [component] present in test data point [True]
246 Text feature [function] present in test data point [True]
255 Text feature [3b] present in test data point [True]
258 Text feature [purified] present in test data point [True]

259 Text feature [several] present in test data point [True]
261 Text feature [systems] present in test data point [True]
263 Text feature [since] present in test data point [True]
264 Text feature [likely] present in test data point [True]
265 Text feature [following] present in test data point [True]
267 Text feature [homologous] present in test data point [True]
268 Text feature [validated] present in test data point [True]
270 Text feature [similarly] present in test data point [True]
271 Text feature [strong] present in test data point [True]
274 Text feature [particular] present in test data point [True]
276 Text feature [note] present in test data point [True]
277 Text feature [control] present in test data point [True]
278 Text feature [lack] present in test data point [True]
279 Text feature [account] present in test data point [True]
282 Text feature [hypothesis] present in test data point [True]
283 Text feature [free] present in test data point [True]
287 Text feature [incubated] present in test data point [True]
289 Text feature [last] present in test data point [True]
290 Text feature [regulate] present in test data point [True]
293 Text feature [bind] present in test data point [True]
296 Text feature [reduced] present in test data point [True]
299 Text feature [26] present in test data point [True]
300 Text feature [overexpressed] present in test data point [True]
301 Text feature [decreased] present in test data point [True]
303 Text feature [high] present in test data point [True]
304 Text feature [suggest] present in test data point [True]
305 Text feature [individuals] present in test data point [True]
307 Text feature [significantly] present in test data point [True]
308 Text feature [localization] present in test data point [True]
309 Text feature [large] present in test data point [True]
310 Text feature [therefore] present in test data point [True]
311 Text feature [leads] present in test data point [True]
312 Text feature [part] present in test data point [True]
314 Text feature [green] present in test data point [True]
316 Text feature [severe] present in test data point [True]
318 Text feature [substrate] present in test data point [True]
319 Text feature [article] present in test data point [True]
321 Text feature [transfection] present in test data point [True]
322 Text feature [responsible] present in test data point [True]
326 Text feature [significant] present in test data point [True]
328 Text feature [enriched] present in test data point [True]
329 Text feature [experiments] present in test data point [True]
330 Text feature [activity] present in test data point [True]
331 Text feature [family] present in test data point [True]
334 Text feature [agarose] present in test data point [True]
336 Text feature [antibody] present in test data point [True]
337 Text feature [splicing] present in test data point [True]
340 Text feature [performed] present in test data point [True]
341 Text feature [classes] present in test data point [True]
342 Text feature [residue] present in test data point [True]
343 Text feature [4a] present in test data point [True]
349 Text feature [used] present in test data point [True]
350 Text feature [promotes] present in test data point [True]
367 Text feature [indeed] present in test data point [True]
370 Text feature [scale] present in test data point [True]
371 Text feature [near] present in test data point [True]
372 Text feature [majority] present in test data point [True]
376 Text feature [predicted] present in test data point [True]
377 Text feature [suggested] present in test data point [True]
379 Text feature [immunoprecipitated] present in test data point [True]
383 Text feature [thought] present in test data point [True]
384 Text feature [indicating] present in test data point [True]
385 Text feature [locus] present in test data point [True]
386 Text feature [ref] present in test data point [True]
388 Text feature [exome] present in test data point [True]
389 Text feature [little] present in test data point [True]
390 Text feature [lanes] present in test data point [True]
391 Text feature [sites] present in test data point [True]
392 Text feature [mutants] present in test data point [True]
394 Text feature [key] present in test data point [True]
395 Text feature [maintained] present in test data point [True]
397 Text feature [specifically] present in test data point [True]
398 Text feature [short] present in test data point [True]
399 Text feature [plasma] present in test data point [True]
400 Text feature [red] present in test data point [True]
401 Text feature [ml] present in test data point [True]
407 Text feature [1996] present in test data point [True]

```

409 Text feature [pten] present in test data point [True]
411 Text feature [complete] present in test data point [True]
412 Text feature [occurs] present in test data point [True]
413 Text feature [fbs] present in test data point [True]
414 Text feature [effect] present in test data point [True]
416 Text feature [open] present in test data point [True]
417 Text feature [materials] present in test data point [True]
420 Text feature [consequences] present in test data point [True]
421 Text feature [catalytic] present in test data point [True]
423 Text feature [explain] present in test data point [True]
424 Text feature [30] present in test data point [True]
425 Text feature [presented] present in test data point [True]
426 Text feature [disorders] present in test data point [True]
428 Text feature [via] present in test data point [True]
429 Text feature [relevant] present in test data point [True]
430 Text feature [baf3] present in test data point [True]
431 Text feature [groups] present in test data point [True]
433 Text feature [toward] present in test data point [True]
434 Text feature [monoclonal] present in test data point [True]
436 Text feature [interact] present in test data point [True]
437 Text feature [location] present in test data point [True]
438 Text feature [levels] present in test data point [True]
439 Text feature [domain] present in test data point [True]
441 Text feature [approximately] present in test data point [True]
443 Text feature [generated] present in test data point [True]
445 Text feature [critical] present in test data point [True]
446 Text feature [spectrum] present in test data point [True]
450 Text feature [initiation] present in test data point [True]
452 Text feature [rna] present in test data point [True]
453 Text feature [regulation] present in test data point [True]
454 Text feature [deletions] present in test data point [True]
455 Text feature [hr] present in test data point [True]
456 Text feature [mouse] present in test data point [True]
458 Text feature [15] present in test data point [True]
459 Text feature [supplemental] present in test data point [True]
461 Text feature [families] present in test data point [True]
465 Text feature [dna] present in test data point [True]
466 Text feature [means] present in test data point [True]
467 Text feature [characterized] present in test data point [True]
470 Text feature [1c] present in test data point [True]
471 Text feature [major] present in test data point [True]
472 Text feature [summary] present in test data point [True]
473 Text feature [involved] present in test data point [True]
476 Text feature [wt] present in test data point [True]
477 Text feature [profile] present in test data point [True]
478 Text feature [states] present in test data point [True]
480 Text feature [variants] present in test data point [True]
481 Text feature [tumorigenesis] present in test data point [True]
484 Text feature [medium] present in test data point [True]
487 Text feature [compound] present in test data point [True]
490 Text feature [seen] present in test data point [True]
491 Text feature [decrease] present in test data point [True]
492 Text feature [sensitivity] present in test data point [True]
494 Text feature [formation] present in test data point [True]
495 Text feature [laboratories] present in test data point [True]
498 Text feature [larger] present in test data point [True]
499 Text feature [pathogenesis] present in test data point [True]
Out of the top 500 features 250 are present in query point

```

4.3.1.3.2. Incorrectly Classified point

In [95]:

```

test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding_LR[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding_LR[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_) [predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names_LR(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: [[0.0174 0.1459 0.0024 0.0458 0.0318 0.009 0.7404 0.0055 0.0017]]

Actual Class : 7

13 Text feature [thyroid] present in test data point [True]
19 Text feature [overexpression] present in test data point [True]
26 Text feature [codon] present in test data point [True]
30 Text feature [clone] present in test data point [True]
32 Text feature [transforming] present in test data point [True]
34 Text feature [contained] present in test data point [True]
45 Text feature [cross] present in test data point [True]
51 Text feature [downstream] present in test data point [True]
54 Text feature [constitutively] present in test data point [True]
65 Text feature [malignant] present in test data point [True]
66 Text feature [activated] present in test data point [True]
67 Text feature [codons] present in test data point [True]
69 Text feature [500] present in test data point [True]
75 Text feature [expressing] present in test data point [True]
80 Text feature [log] present in test data point [True]
89 Text feature [activation] present in test data point [True]
93 Text feature [pathways] present in test data point [True]
101 Text feature [akt] present in test data point [True]
103 Text feature [transformed] present in test data point [True]
108 Text feature [nm] present in test data point [True]
119 Text feature [ligand] present in test data point [True]
121 Text feature [see] present in test data point [True]
123 Text feature [fold] present in test data point [True]
149 Text feature [12] present in test data point [True]
150 Text feature [commonly] present in test data point [True]
151 Text feature [increases] present in test data point [True]
152 Text feature [exposed] present in test data point [True]
158 Text feature [cos] present in test data point [True]
160 Text feature [transient] present in test data point [True]
161 Text feature [transformation] present in test data point [True]
162 Text feature [oncogene] present in test data point [True]
164 Text feature [order] present in test data point [True]
165 Text feature [cancers] present in test data point [True]
166 Text feature [occur] present in test data point [True]
169 Text feature [upstream] present in test data point [True]
173 Text feature [free] present in test data point [True]
177 Text feature [increased] present in test data point [True]
180 Text feature [university] present in test data point [True]
181 Text feature [high] present in test data point [True]
182 Text feature [phosphorylated] present in test data point [True]
185 Text feature [cytoplasmic] present in test data point [True]
191 Text feature [transduction] present in test data point [True]
192 Text feature [presence] present in test data point [True]
193 Text feature [similar] present in test data point [True]
195 Text feature [equivalent] present in test data point [True]
203 Text feature [1a] present in test data point [True]
205 Text feature [colonies] present in test data point [True]
208 Text feature [useful] present in test data point [True]
209 Text feature [grown] present in test data point [True]
211 Text feature [carcinomas] present in test data point [True]
217 Text feature [early] present in test data point [True]
218 Text feature [require] present in test data point [True]
226 Text feature [probe] present in test data point [True]
227 Text feature [g1] present in test data point [True]
231 Text feature [total] present in test data point [True]
232 Text feature [wide] present in test data point [True]
234 Text feature [right] present in test data point [True]
236 Text feature [vectors] present in test data point [True]
237 Text feature [rate] present in test data point [True]
240 Text feature [cyclin] present in test data point [True]
244 Text feature [event] present in test data point [True]
245 Text feature [regulated] present in test data point [True]
246 Text feature [carcinoma] present in test data point [True]
253 Text feature [formed] present in test data point [True]
259 Text feature [suggested] present in test data point [True]
274 Text feature [vector] present in test data point [True]
278 Text feature [thought] present in test data point [True]
282 Text feature [primarily] present in test data point [True]
283 Text feature [regulate] present in test data point [True]
284 Text feature [colony] present in test data point [True]
286 Text feature [parallel] present in test data point [True]

287 Text feature [molecules] present in test data point [True]
288 Text feature [density] present in test data point [True]
289 Text feature [per] present in test data point [True]
290 Text feature [phosphorylation] present in test data point [True]
292 Text feature [inhibition] present in test data point [True]
294 Text feature [advanced] present in test data point [True]
295 Text feature [000] present in test data point [True]
297 Text feature [five] present in test data point [True]
299 Text feature [55] present in test data point [True]
300 Text feature [shows] present in test data point [True]
306 Text feature [absence] present in test data point [True]
307 Text feature [phosphate] present in test data point [True]
310 Text feature [added] present in test data point [True]
312 Text feature [fig] present in test data point [True]
316 Text feature [screened] present in test data point [True]
317 Text feature [resulting] present in test data point [True]
320 Text feature [volume] present in test data point [True]
321 Text feature [tumors] present in test data point [True]
323 Text feature [strong] present in test data point [True]
325 Text feature [immunohistochemistry] present in test data point [True]
328 Text feature [spectrum] present in test data point [True]
329 Text feature [without] present in test data point [True]
333 Text feature [decrease] present in test data point [True]
334 Text feature [work] present in test data point [True]
336 Text feature [suggest] present in test data point [True]
343 Text feature [24] present in test data point [True]
346 Text feature [translocation] present in test data point [True]
347 Text feature [distinct] present in test data point [True]
348 Text feature [11] present in test data point [True]
349 Text feature [serum] present in test data point [True]
358 Text feature [linked] present in test data point [True]
364 Text feature [signaling] present in test data point [True]
365 Text feature [compared] present in test data point [True]
367 Text feature [difference] present in test data point [True]
369 Text feature [obtained] present in test data point [True]
378 Text feature [increase] present in test data point [True]
379 Text feature [exposure] present in test data point [True]
380 Text feature [crystal] present in test data point [True]
381 Text feature [factor] present in test data point [True]
383 Text feature [sequences] present in test data point [True]
385 Text feature [also] present in test data point [True]
386 Text feature [determine] present in test data point [True]
390 Text feature [test] present in test data point [True]
395 Text feature [previously] present in test data point [True]
398 Text feature [appears] present in test data point [True]
399 Text feature [kinases] present in test data point [True]
401 Text feature [water] present in test data point [True]
406 Text feature [plates] present in test data point [True]
414 Text feature [nih] present in test data point [True]
415 Text feature [interestingly] present in test data point [True]
416 Text feature [membrane] present in test data point [True]
417 Text feature [agar] present in test data point [True]
418 Text feature [oncogenic] present in test data point [True]
424 Text feature [kit] present in test data point [True]
426 Text feature [frequent] present in test data point [True]
427 Text feature [forms] present in test data point [True]
430 Text feature [possibility] present in test data point [True]
433 Text feature [cells] present in test data point [True]
435 Text feature [non] present in test data point [True]
437 Text feature [specific] present in test data point [True]
439 Text feature [cultured] present in test data point [True]
440 Text feature [structure] present in test data point [True]
442 Text feature [harboring] present in test data point [True]
443 Text feature [18] present in test data point [True]
452 Text feature [mutant] present in test data point [True]
453 Text feature [suggests] present in test data point [True]
457 Text feature [correlated] present in test data point [True]
460 Text feature [form] present in test data point [True]
466 Text feature [signal] present in test data point [True]
470 Text feature [17] present in test data point [True]
477 Text feature [could] present in test data point [True]
478 Text feature [extracts] present in test data point [True]
479 Text feature [dl] present in test data point [True]
480 Text feature [demonstrated] present in test data point [True]
481 Text feature [adjacent] present in test data point [True]
482 Text feature [carried] present in test data point [True]
485 Text feature [resistant] present in test data point [True]

487 Text feature [since] present in test data point [True]
497 Text feature [lymphoma] present in test data point [True]
Out of the top 500 features 150 are present in query point

4.3.2. Without Class balancing

4.3.2.1. Hyper paramter tuning

In [280]:

```
# 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(-8, 3)]
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_LR, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_onehotCoding_LR, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding_LR)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-10))
    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_LR, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding_LR, train_y)
```

```

sig_clf.fit(train_x_onehotCoding_LR, train_y)

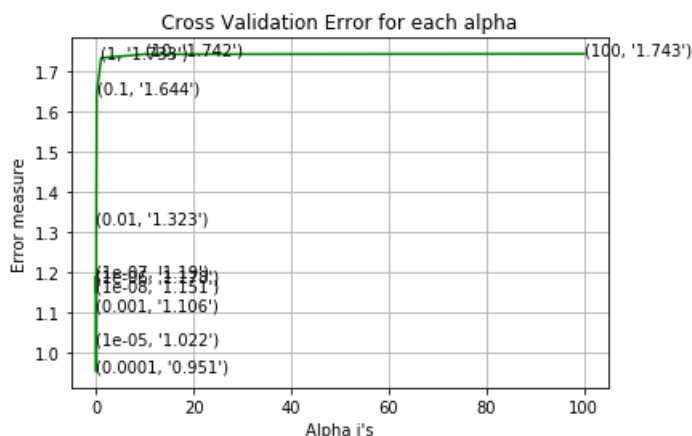
predict_y = sig_clf.predict_proba(train_x_onehotCoding_LR)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-10))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding_LR)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_lo
ss(y_cv, predict_y, labels=clf.classes_, eps=1e-10))
predict_y = sig_clf.predict_proba(test_x_onehotCoding_LR)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, p
redict_y, labels=clf.classes_, eps=1e-10))

```

```

for alpha = 1e-08
Log Loss : 1.1511800826167682
for alpha = 1e-07
Log Loss : 1.1896250192982063
for alpha = 1e-06
Log Loss : 1.1779188573085753
for alpha = 1e-05
Log Loss : 1.0218148529752458
for alpha = 0.0001
Log Loss : 0.9513510749847864
for alpha = 0.001
Log Loss : 1.1055812095324928
for alpha = 0.01
Log Loss : 1.3227600423889374
for alpha = 0.1
Log Loss : 1.6435369582934525
for alpha = 1
Log Loss : 1.7331501170330132
for alpha = 10
Log Loss : 1.7423153516832723
for alpha = 100
Log Loss : 1.7432966316294265

```



```

For values of best alpha = 0.0001 The train log loss is: 0.7118134230187894
For values of best alpha = 0.0001 The cross validation log loss is: 0.9513510749858702
For values of best alpha = 0.0001 The test log loss is: 0.9470475162688973

```

4.3.2.2. Testing model with best hyper parameters

In [282]:

```

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

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

```



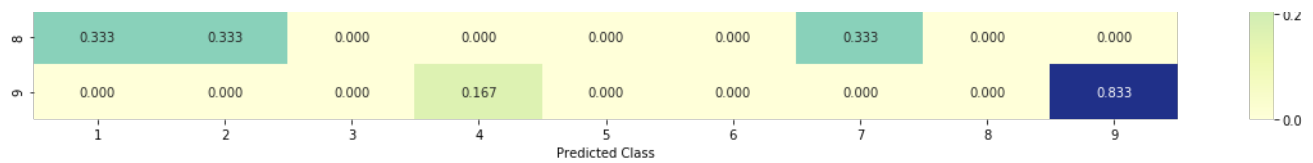
```
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding_LR, train_y, cv_x_onehotCoding_LR, cv_y, cl
f)
```

Number of mis-classified points : 0.32706766917293234

Original Class	1	2	3	4	5	6	7	8	9
1	58.000	1.000	0.000	22.000	4.000	4.000	2.000	0.000	0.000
2	1.000	31.000	0.000	2.000	0.000	1.000	37.000	0.000	0.000
3	1.000	0.000	1.000	4.000	0.000	0.000	8.000	0.000	0.000
4	15.000	0.000	0.000	85.000	1.000	0.000	9.000	0.000	0.000
5	8.000	1.000	0.000	7.000	13.000	3.000	7.000	0.000	0.000
6	5.000	0.000	0.000	5.000	2.000	23.000	9.000	0.000	0.000
7	0.000	8.000	0.000	1.000	1.000	1.000	142.000	0.000	0.000
8	1.000	1.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000
9	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	5.000

	1	2	3	4	5	6	7	8	9
1	0.652	0.024	0.000	0.173	0.190	0.125	0.009	0.000	0.000
2	0.011	0.738	0.000	0.016	0.000	0.031	0.172	0.000	0.000
3	0.011	0.000	1.000	0.031	0.000	0.000	0.037	0.000	0.000
4	0.169	0.000	0.000	0.669	0.048	0.000	0.042	0.000	0.000
5	0.090	0.024	0.000	0.055	0.619	0.094	0.033	0.000	0.000
6	0.056	0.000	0.000	0.039	0.095	0.719	0.042	0.000	0.000
7	0.000	0.190	0.000	0.008	0.048	0.031	0.660	0.000	0.000
8	0.011	0.024	0.000	0.000	0.000	0.000	0.005	0.000	0.000
9	0.000	0.000	0.000	0.008	0.000	0.000	0.000	0.000	1.000

Original Class	0	1	2	3	4	5	6	7
0	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
1	0.637	0.011	0.000	0.242	0.044	0.044	0.022	0.000
2	0.014	0.431	0.000	0.028	0.000	0.014	0.514	0.000
3	0.071	0.000	0.071	0.286	0.000	0.000	0.571	0.000
4	0.136	0.000	0.000	0.773	0.009	0.000	0.082	0.000
5	0.205	0.026	0.000	0.179	0.333	0.077	0.179	0.000
6	0.114	0.000	0.000	0.114	0.045	0.523	0.205	0.000
7	0.000	0.052	0.000	0.007	0.007	0.007	0.928	0.000



4.3.2.3. Feature Importance, Correctly Classified point

In [98]:

```
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding_LR,train_y)
test_point_index = 1
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding_LR[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding_LR[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,no_feature]
print("-"*50)
get_impfeature_names_LR(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: [[1.630e-02 5.900e-03 2.000e-04 8.046e-01 1.700e-03 3.000e-04 1.682e-01
2.800e-03 0.000e+00]]

Actual Class : 1

```
-----
15 Text feature [material] present in test data point [True]
18 Text feature [beads] present in test data point [True]
22 Text feature [iii] present in test data point [True]
29 Text feature [nucleus] present in test data point [True]
32 Text feature [washed] present in test data point [True]
34 Text feature [pcdna3] present in test data point [True]
35 Text feature [motifs] present in test data point [True]
38 Text feature [encodes] present in test data point [True]
40 Text feature [2a] present in test data point [True]
41 Text feature [defect] present in test data point [True]
48 Text feature [sigma] present in test data point [True]
51 Text feature [tagged] present in test data point [True]
53 Text feature [density] present in test data point [True]
54 Text feature [recombination] present in test data point [True]
56 Text feature [loss] present in test data point [True]
57 Text feature [show] present in test data point [True]
58 Text feature [lacking] present in test data point [True]
59 Text feature [endogenous] present in test data point [True]
62 Text feature [functionally] present in test data point [True]
63 Text feature [contribute] present in test data point [True]
65 Text feature [homozygous] present in test data point [True]
68 Text feature [experimental] present in test data point [True]
79 Text feature [sds] present in test data point [True]
80 Text feature [transfected] present in test data point [True]
83 Text feature [pbs] present in test data point [True]
85 Text feature [domains] present in test data point [True]
87 Text feature [cohort] present in test data point [True]
88 Text feature [deleted] present in test data point [True]
89 Text feature [linker] present in test data point [True]
91 Text feature [cannot] present in test data point [True]
93 Text feature [mm] present in test data point [True]
94 Text feature [see] present in test data point [True]
95 Text feature [biotechnology] present in test data point [True]
97 Text feature [changes] present in test data point [True]
99 Text feature [2010] present in test data point [True]
100 Text feature [phosphatase] present in test data point [True]
102 Text feature [suggesting] present in test data point [True]
103 Text feature [lysis] present in test data point [True]
106 Text feature [plates] present in test data point [True]
110 Text feature [due] present in test data point [True]
111 Text feature [online] present in test data point [True]
113 Text feature [representative] present in test data point [True]
115 Text feature [express] present in test data point [True]
```

115 Text feature [express] present in test data point [True]
117 Text feature [29] present in test data point [True]
119 Text feature [nacl] present in test data point [True]
120 Text feature [phenotype] present in test data point [True]
121 Text feature [s4] present in test data point [True]
124 Text feature [lower] present in test data point [True]
125 Text feature [1998] present in test data point [True]
126 Text feature [genetic] present in test data point [True]
130 Text feature [defective] present in test data point [True]
134 Text feature [motif] present in test data point [True]
135 Text feature [bound] present in test data point [True]
136 Text feature [min] present in test data point [True]
138 Text feature [detailed] present in test data point [True]
140 Text feature [right] present in test data point [True]
142 Text feature [isoform] present in test data point [True]
147 Text feature [assays] present in test data point [True]
150 Text feature [rhoa] present in test data point [True]
153 Text feature [product] present in test data point [True]
154 Text feature [regulates] present in test data point [True]
155 Text feature [risk] present in test data point [True]
159 Text feature [induced] present in test data point [True]
160 Text feature [santa] present in test data point [True]
163 Text feature [light] present in test data point [True]
164 Text feature [primarily] present in test data point [True]
167 Text feature [times] present in test data point [True]
168 Text feature [heterozygous] present in test data point [True]
170 Text feature [conserved] present in test data point [True]
171 Text feature [thus] present in test data point [True]
172 Text feature [depletion] present in test data point [True]
173 Text feature [observation] present in test data point [True]
174 Text feature [procedures] present in test data point [True]
175 Text feature [knockdown] present in test data point [True]
176 Text feature [functional] present in test data point [True]
177 Text feature [modified] present in test data point [True]
182 Text feature [assay] present in test data point [True]
184 Text feature [s2] present in test data point [True]
186 Text feature [2b] present in test data point [True]
187 Text feature [endometrial] present in test data point [True]
189 Text feature [tested] present in test data point [True]
191 Text feature [stability] present in test data point [True]
193 Text feature [strongly] present in test data point [True]
195 Text feature [several] present in test data point [True]
197 Text feature [localized] present in test data point [True]
198 Text feature [experiment] present in test data point [True]
200 Text feature [co] present in test data point [True]
202 Text feature [culture] present in test data point [True]
204 Text feature [elevated] present in test data point [True]
205 Text feature [fetal] present in test data point [True]
206 Text feature [following] present in test data point [True]
207 Text feature [bars] present in test data point [True]
208 Text feature [26] present in test data point [True]
209 Text feature [database] present in test data point [True]
210 Text feature [reduced] present in test data point [True]
212 Text feature [cruz] present in test data point [True]
214 Text feature [low] present in test data point [True]
218 Text feature [western] present in test data point [True]
221 Text feature [differences] present in test data point [True]
223 Text feature [incidence] present in test data point [True]
227 Text feature [caused] present in test data point [True]
231 Text feature [hypothesis] present in test data point [True]
234 Text feature [proportion] present in test data point [True]
239 Text feature [protein] present in test data point [True]
240 Text feature [tris] present in test data point [True]
241 Text feature [similarly] present in test data point [True]
243 Text feature [buffer] present in test data point [True]
245 Text feature [bind] present in test data point [True]
246 Text feature [supporting] present in test data point [True]
249 Text feature [last] present in test data point [True]
251 Text feature [note] present in test data point [True]
252 Text feature [indicate] present in test data point [True]
253 Text feature [particular] present in test data point [True]
254 Text feature [significantly] present in test data point [True]
257 Text feature [examine] present in test data point [True]
258 Text feature [regulate] present in test data point [True]
260 Text feature [systems] present in test data point [True]
261 Text feature [likely] present in test data point [True]
263 Text feature [investigated] present in test data point [True]
265 Text feature [mice] present in test data point [True]

266 Text feature [mice] present in test data point [True]
270 Text feature [leads] present in test data point [True]
273 Text feature [3b] present in test data point [True]
275 Text feature [lack] present in test data point [True]
277 Text feature [component] present in test data point [True]
278 Text feature [control] present in test data point [True]
280 Text feature [poor] present in test data point [True]
281 Text feature [affected] present in test data point [True]
282 Text feature [indeed] present in test data point [True]
284 Text feature [green] present in test data point [True]
285 Text feature [suggested] present in test data point [True]
286 Text feature [comparison] present in test data point [True]
287 Text feature [overexpressed] present in test data point [True]
288 Text feature [4a] present in test data point [True]
290 Text feature [rabbit] present in test data point [True]
291 Text feature [strong] present in test data point [True]
292 Text feature [function] present in test data point [True]
293 Text feature [nuclear] present in test data point [True]
295 Text feature [thought] present in test data point [True]
297 Text feature [high] present in test data point [True]
299 Text feature [open] present in test data point [True]
300 Text feature [large] present in test data point [True]
302 Text feature [short] present in test data point [True]
303 Text feature [transfection] present in test data point [True]
304 Text feature [near] present in test data point [True]
308 Text feature [suggest] present in test data point [True]
309 Text feature [part] present in test data point [True]
310 Text feature [since] present in test data point [True]
311 Text feature [individuals] present in test data point [True]
312 Text feature [therefore] present in test data point [True]
315 Text feature [account] present in test data point [True]
316 Text feature [p85] present in test data point [True]
317 Text feature [classes] present in test data point [True]
318 Text feature [significant] present in test data point [True]
320 Text feature [antibody] present in test data point [True]
322 Text feature [incubated] present in test data point [True]
323 Text feature [splicing] present in test data point [True]
326 Text feature [specifically] present in test data point [True]
330 Text feature [baf3] present in test data point [True]
334 Text feature [performed] present in test data point [True]
336 Text feature [free] present in test data point [True]
337 Text feature [validated] present in test data point [True]
338 Text feature [severe] present in test data point [True]
339 Text feature [decreased] present in test data point [True]
340 Text feature [majority] present in test data point [True]
341 Text feature [predicted] present in test data point [True]
342 Text feature [homologous] present in test data point [True]
344 Text feature [lanes] present in test data point [True]
345 Text feature [pathogenesis] present in test data point [True]
346 Text feature [disorders] present in test data point [True]
347 Text feature [enriched] present in test data point [True]
348 Text feature [substrate] present in test data point [True]
349 Text feature [article] present in test data point [True]
350 Text feature [toward] present in test data point [True]
351 Text feature [localization] present in test data point [True]
353 Text feature [domain] present in test data point [True]
355 Text feature [indicating] present in test data point [True]
358 Text feature [purified] present in test data point [True]
359 Text feature [deletions] present in test data point [True]
360 Text feature [residue] present in test data point [True]
362 Text feature [activity] present in test data point [True]
363 Text feature [used] present in test data point [True]
364 Text feature [experiments] present in test data point [True]
381 Text feature [plasma] present in test data point [True]
382 Text feature [onto] present in test data point [True]
384 Text feature [family] present in test data point [True]
387 Text feature [hr] present in test data point [True]
389 Text feature [consequences] present in test data point [True]
390 Text feature [materials] present in test data point [True]
391 Text feature [sensitivity] present in test data point [True]
392 Text feature [via] present in test data point [True]
393 Text feature [laboratories] present in test data point [True]
394 Text feature [little] present in test data point [True]
397 Text feature [sites] present in test data point [True]
398 Text feature [agarose] present in test data point [True]
399 Text feature [initiation] present in test data point [True]
400 Text feature [key] present in test data point [True]
402 Text feature [weak] present in test data point [True]

```

402 Text feature [red] present in test data point [True]
406 Text feature [mutants] present in test data point [True]
410 Text feature [30] present in test data point [True]
411 Text feature [inc] present in test data point [True]
412 Text feature [rna] present in test data point [True]
413 Text feature [promotes] present in test data point [True]
415 Text feature [immunoprecipitated] present in test data point [True]
416 Text feature [responsible] present in test data point [True]
420 Text feature [maintained] present in test data point [True]
423 Text feature [location] present in test data point [True]
424 Text feature [occurs] present in test data point [True]
426 Text feature [profile] present in test data point [True]
427 Text feature [complete] present in test data point [True]
428 Text feature [groups] present in test data point [True]
429 Text feature [exome] present in test data point [True]
432 Text feature [supplemental] present in test data point [True]
435 Text feature [explain] present in test data point [True]
437 Text feature [medium] present in test data point [True]
438 Text feature [locus] present in test data point [True]
439 Text feature [regulation] present in test data point [True]
440 Text feature [s1] present in test data point [True]
441 Text feature [relevant] present in test data point [True]
442 Text feature [instructions] present in test data point [True]
443 Text feature [characterized] present in test data point [True]
450 Text feature [tumorigenesis] present in test data point [True]
452 Text feature [scale] present in test data point [True]
453 Text feature [compound] present in test data point [True]
454 Text feature [summary] present in test data point [True]
455 Text feature [15] present in test data point [True]
457 Text feature [discovery] present in test data point [True]
458 Text feature [major] present in test data point [True]
461 Text feature [ref] present in test data point [True]
463 Text feature [spectrum] present in test data point [True]
464 Text feature [effect] present in test data point [True]
465 Text feature [fbs] present in test data point [True]
466 Text feature [prepared] present in test data point [True]
469 Text feature [ii] present in test data point [True]
470 Text feature [involved] present in test data point [True]
474 Text feature [monoclonal] present in test data point [True]
475 Text feature [293t] present in test data point [True]
476 Text feature [states] present in test data point [True]
477 Text feature [ml] present in test data point [True]
478 Text feature [recombinant] present in test data point [True]
481 Text feature [dna] present in test data point [True]
482 Text feature [altered] present in test data point [True]
483 Text feature [isolated] present in test data point [True]
487 Text feature [catalytic] present in test data point [True]
488 Text feature [generated] present in test data point [True]
489 Text feature [resulting] present in test data point [True]
492 Text feature [larger] present in test data point [True]
494 Text feature [collected] present in test data point [True]
499 Text feature [investigate] present in test data point [True]
Out of the top 500 features 248 are present in query point

```

4.3.2.4. Feature Importance, Inorrectly Classified point

In [99]:

```

test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding_LR[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
      np.round(sig_clf.predict_proba(test_x_onehotCoding_LR[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,no_feature]
print("-"*50)
get_impfeature_names_LR(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: [[0.0161 0.1602 0.003 0.0489 0.0332 0.0094 0.7225 0.0054 0.0012]]

Actual Class : 7

8 Text feature [thyroid] present in test data point [True]
10 Text feature [overexpression] present in test data point [True]
13 Text feature [codon] present in test data point [True]
26 Text feature [clone] present in test data point [True]
31 Text feature [cross] present in test data point [True]
35 Text feature [transforming] present in test data point [True]
36 Text feature [contained] present in test data point [True]
49 Text feature [downstream] present in test data point [True]
60 Text feature [constitutively] present in test data point [True]
76 Text feature [log] present in test data point [True]
80 Text feature [codons] present in test data point [True]
82 Text feature [see] present in test data point [True]
84 Text feature [malignant] present in test data point [True]
88 Text feature [cytoplasmic] present in test data point [True]
92 Text feature [500] present in test data point [True]
96 Text feature [expressing] present in test data point [True]
100 Text feature [activation] present in test data point [True]
102 Text feature [activated] present in test data point [True]
106 Text feature [upstream] present in test data point [True]
110 Text feature [fold] present in test data point [True]
111 Text feature [nm] present in test data point [True]
112 Text feature [increased] present in test data point [True]
118 Text feature [wide] present in test data point [True]
128 Text feature [increases] present in test data point [True]
130 Text feature [phosphorylated] present in test data point [True]
131 Text feature [occur] present in test data point [True]
132 Text feature [ligand] present in test data point [True]
143 Text feature [commonly] present in test data point [True]
147 Text feature [akt] present in test data point [True]
148 Text feature [pathways] present in test data point [True]
149 Text feature [la] present in test data point [True]
166 Text feature [shows] present in test data point [True]
167 Text feature [useful] present in test data point [True]
169 Text feature [right] present in test data point [True]
171 Text feature [event] present in test data point [True]
174 Text feature [formed] present in test data point [True]
176 Text feature [cos] present in test data point [True]
177 Text feature [parallel] present in test data point [True]
179 Text feature [high] present in test data point [True]
185 Text feature [l2] present in test data point [True]
189 Text feature [exposed] present in test data point [True]
190 Text feature [similar] present in test data point [True]
196 Text feature [strong] present in test data point [True]
198 Text feature [oncogene] present in test data point [True]
207 Text feature [transient] present in test data point [True]
209 Text feature [transformation] present in test data point [True]
210 Text feature [cancers] present in test data point [True]
212 Text feature [55] present in test data point [True]
213 Text feature [grown] present in test data point [True]
216 Text feature [presence] present in test data point [True]
218 Text feature [suggested] present in test data point [True]
220 Text feature [free] present in test data point [True]
221 Text feature [phosphorylation] present in test data point [True]
225 Text feature [advanced] present in test data point [True]
226 Text feature [transformed] present in test data point [True]
227 Text feature [total] present in test data point [True]
228 Text feature [regulated] present in test data point [True]
235 Text feature [early] present in test data point [True]
237 Text feature [g1] present in test data point [True]
238 Text feature [require] present in test data point [True]
239 Text feature [thought] present in test data point [True]
247 Text feature [linked] present in test data point [True]
249 Text feature [l1] present in test data point [True]
254 Text feature [probe] present in test data point [True]
255 Text feature [colonies] present in test data point [True]
257 Text feature [distinct] present in test data point [True]
261 Text feature [carcinomas] present in test data point [True]
267 Text feature [colony] present in test data point [True]
276 Text feature [primarily] present in test data point [True]
279 Text feature [fig] present in test data point [True]
280 Text feature [five] present in test data point [True]
281 Text feature [university] present in test data point [True]
283 Text feature [carcinoma] present in test data point [True]
284 Text feature [without] present in test data point [True]
286 Text feature [24] present in test data point [True]
291 Text feature [equivalent] present in test data point [True]
292 Text feature [spectrum] present in test data point [True]

292 Text feature [spectrum] present in test data point [True]
293 Text feature [resulting] present in test data point [True]
301 Text feature [order] present in test data point [True]
303 Text feature [test] present in test data point [True]
309 Text feature [cyclin] present in test data point [True]
310 Text feature [decrease] present in test data point [True]
313 Text feature [frequent] present in test data point [True]
315 Text feature [volume] present in test data point [True]
318 Text feature [vector] present in test data point [True]
321 Text feature [per] present in test data point [True]
324 Text feature [work] present in test data point [True]
325 Text feature [000] present in test data point [True]
327 Text feature [added] present in test data point [True]
328 Text feature [density] present in test data point [True]
330 Text feature [transduction] present in test data point [True]
334 Text feature [suggest] present in test data point [True]
337 Text feature [determine] present in test data point [True]
344 Text feature [rate] present in test data point [True]
347 Text feature [increase] present in test data point [True]
351 Text feature [regulate] present in test data point [True]
352 Text feature [specific] present in test data point [True]
353 Text feature [molecules] present in test data point [True]
354 Text feature [previously] present in test data point [True]
362 Text feature [kit] present in test data point [True]
364 Text feature [immunohistochemistry] present in test data point [True]
365 Text feature [vectors] present in test data point [True]
367 Text feature [agar] present in test data point [True]
369 Text feature [also] present in test data point [True]
370 Text feature [inhibition] present in test data point [True]
372 Text feature [signal] present in test data point [True]
376 Text feature [phosphate] present in test data point [True]
377 Text feature [lane] present in test data point [True]
378 Text feature [sequences] present in test data point [True]
382 Text feature [importance] present in test data point [True]
383 Text feature [94] present in test data point [True]
386 Text feature [crystal] present in test data point [True]
387 Text feature [compared] present in test data point [True]
388 Text feature [might] present in test data point [True]
389 Text feature [screened] present in test data point [True]
391 Text feature [harboring] present in test data point [True]
395 Text feature [17] present in test data point [True]
399 Text feature [possibility] present in test data point [True]
402 Text feature [suggests] present in test data point [True]
403 Text feature [canonical] present in test data point [True]
404 Text feature [experimental] present in test data point [True]
407 Text feature [obtained] present in test data point [True]
408 Text feature [signaling] present in test data point [True]
411 Text feature [non] present in test data point [True]
413 Text feature [factor] present in test data point [True]
416 Text feature [nih] present in test data point [True]
419 Text feature [serum] present in test data point [True]
423 Text feature [form] present in test data point [True]
425 Text feature [interestingly] present in test data point [True]
426 Text feature [lymphoma] present in test data point [True]
427 Text feature [oncogenic] present in test data point [True]
429 Text feature [translocation] present in test data point [True]
430 Text feature [water] present in test data point [True]
431 Text feature [37] present in test data point [True]
439 Text feature [absence] present in test data point [True]
441 Text feature [since] present in test data point [True]
442 Text feature [whereas] present in test data point [True]
445 Text feature [structure] present in test data point [True]
452 Text feature [profiles] present in test data point [True]
455 Text feature [cells] present in test data point [True]
456 Text feature [gastric] present in test data point [True]
459 Text feature [recently] present in test data point [True]
461 Text feature [selected] present in test data point [True]
462 Text feature [mutation] present in test data point [True]
463 Text feature [could] present in test data point [True]
470 Text feature [appears] present in test data point [True]
471 Text feature [residue] present in test data point [True]
473 Text feature [mutant] present in test data point [True]
476 Text feature [kinases] present in test data point [True]
481 Text feature [18] present in test data point [True]
484 Text feature [allowed] present in test data point [True]
485 Text feature [difference] present in test data point [True]
486 Text feature [exposure] present in test data point [True]
487 Text feature [formed] present in test data point [True]

```

48/ Text feature [forms] present in test data point [True]
491 Text feature [membrane] present in test data point [True]
496 Text feature [colon] present in test data point [True]
498 Text feature [carried] present in test data point [True]
Out of the top 500 features 157 are present in query point

```

4.4. Linear Support Vector Machines

4.4.1. Hyper paramter tuning

In [100]:

```

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

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

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

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

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

fig, ax = plt.subplots()
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=best_alpha, kernel='linear', probability=True, class_weight='balanced')

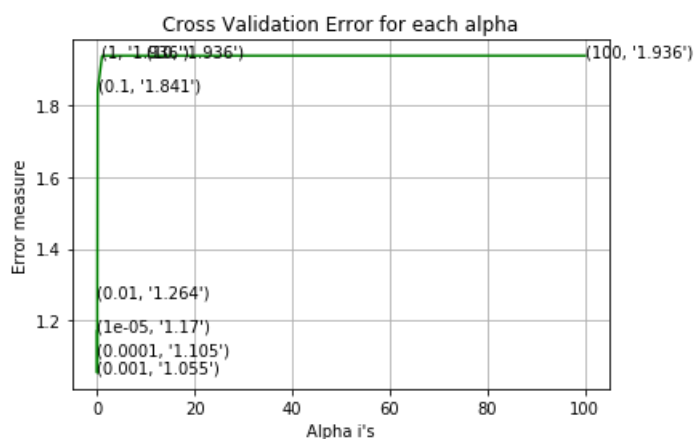
```



```
# clf = SVC(C=1, 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 : 1.1704842903220716
for C = 0.0001
Log Loss : 1.104636993729758
for C = 0.001
Log Loss : 1.0548220249478817
for C = 0.01
Log Loss : 1.2635906904398093
for C = 0.1
Log Loss : 1.8412825872146064
for C = 1
Log Loss : 1.9363885628244581
for C = 10
Log Loss : 1.9363882046423924
for C = 100
Log Loss : 1.936388172268587
```



```
For values of best alpha = 0.001 The train log loss is: 0.9280684335507319
For values of best alpha = 0.001 The cross validation log loss is: 1.0548220249478817
For values of best alpha = 0.001 The test log loss is: 1.1325369474956772
```

4.4.2. Testing model with best hyper parameters

In [101]:

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

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

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

```
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/mathematical-derivation-copy-8/
# -----

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

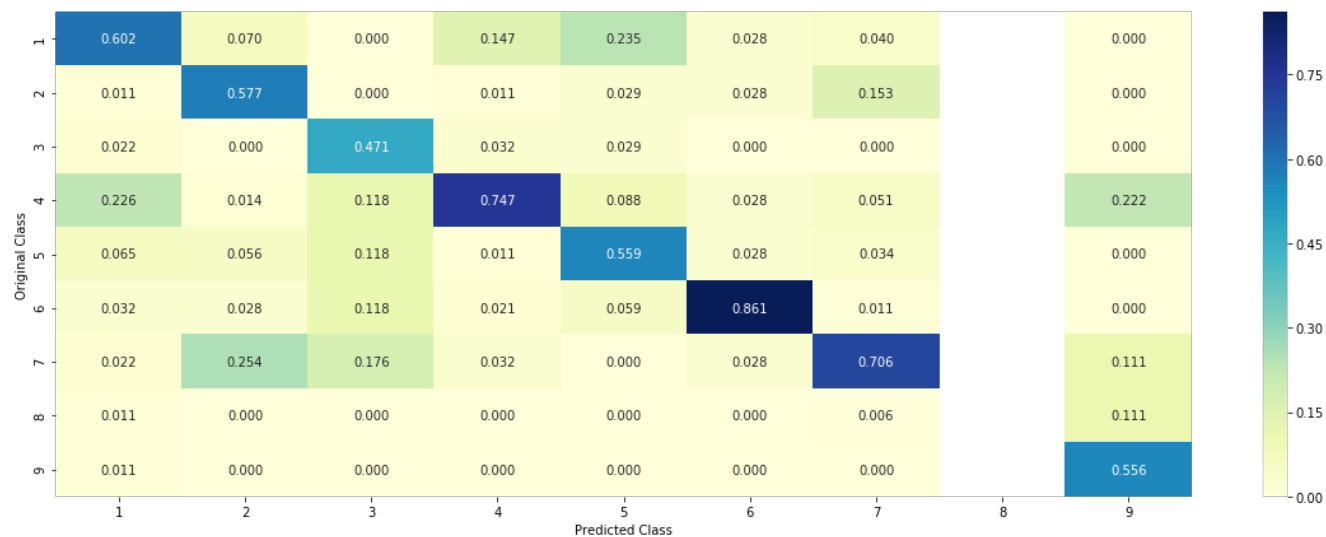
Log loss : 1.0548220249478817

Number of mis-classified points : 0.3308270676691729

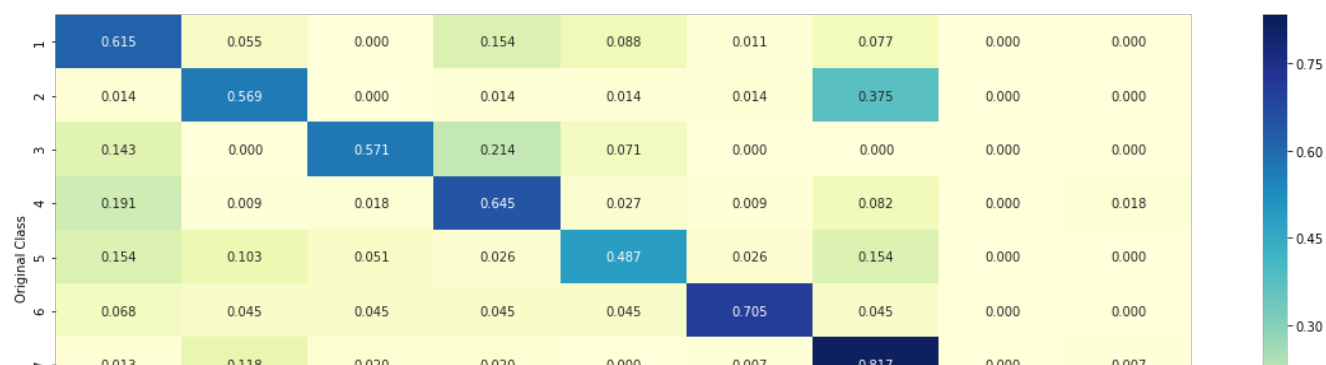
----- Confusion matrix -----

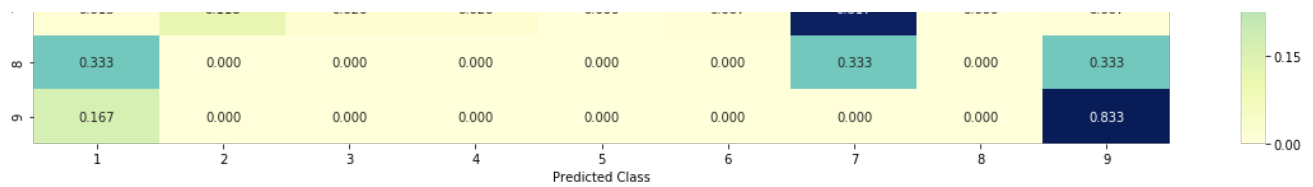


----- Precision matrix (Column Sum=1) -----



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





4.3.3. Feature Importance

4.3.3.1. For Correctly classified point

In [102]:

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

Predicted Class : 4

Predicted Class Probabilities: [[0.0675 0.0216 0.007 0.7767 0.0163 0.0057 0.0987 0.0047 0.0017]]

Actual Class : 1

```
-----
18 Text feature [iii] present in test data point [True]
27 Text feature [encodes] present in test data point [True]
29 Text feature [beads] present in test data point [True]
31 Text feature [transfected] present in test data point [True]
39 Text feature [show] present in test data point [True]
40 Text feature [material] present in test data point [True]
41 Text feature [homozygous] present in test data point [True]
45 Text feature [density] present in test data point [True]
47 Text feature [see] present in test data point [True]
49 Text feature [nucleus] present in test data point [True]
51 Text feature [suggesting] present in test data point [True]
61 Text feature [pbs] present in test data point [True]
63 Text feature [prepared] present in test data point [True]
64 Text feature [2a] present in test data point [True]
68 Text feature [help] present in test data point [True]
69 Text feature [proportion] present in test data point [True]
70 Text feature [truncated] present in test data point [True]
71 Text feature [high] present in test data point [True]
73 Text feature [washed] present in test data point [True]
74 Text feature [plates] present in test data point [True]
75 Text feature [western] present in test data point [True]
77 Text feature [short] present in test data point [True]
80 Text feature [online] present in test data point [True]
82 Text feature [s4] present in test data point [True]
83 Text feature [1998] present in test data point [True]
87 Text feature [plasma] present in test data point [True]
88 Text feature [ii] present in test data point [True]
89 Text feature [characterized] present in test data point [True]
91 Text feature [due] present in test data point [True]
94 Text feature [150] present in test data point [True]
104 Text feature [mm] present in test data point [True]
109 Text feature [tagged] present in test data point [True]
110 Text feature [functionally] present in test data point [True]
113 Text feature [incidence] present in test data point [True]
117 Text feature [times] present in test data point [True]
118 Text feature [affected] present in test data point [True]
119 Text feature [motif] present in test data point [True]
120 Text feature [supporting] present in test data point [True]
124 Text feature [lanes] present in test data point [True]
```

125 Text feature [involved] present in test data point [True]
127 Text feature [comparison] present in test data point [True]
129 Text feature [changes] present in test data point [True]
131 Text feature [motifs] present in test data point [True]
133 Text feature [particular] present in test data point [True]
134 Text feature [localization] present in test data point [True]
135 Text feature [onto] present in test data point [True]
136 Text feature [agarose] present in test data point [True]
137 Text feature [s2] present in test data point [True]
138 Text feature [stability] present in test data point [True]
139 Text feature [via] present in test data point [True]
140 Text feature [explain] present in test data point [True]
145 Text feature [29] present in test data point [True]
146 Text feature [defect] present in test data point [True]
147 Text feature [pcdna3] present in test data point [True]
149 Text feature [loss] present in test data point [True]
150 Text feature [indicate] present in test data point [True]
151 Text feature [investigate] present in test data point [True]
152 Text feature [several] present in test data point [True]
153 Text feature [6a] present in test data point [True]
154 Text feature [express] present in test data point [True]
157 Text feature [exposure] present in test data point [True]
158 Text feature [lysis] present in test data point [True]
161 Text feature [right] present in test data point [True]
163 Text feature [light] present in test data point [True]
164 Text feature [recombination] present in test data point [True]
165 Text feature [families] present in test data point [True]
166 Text feature [p85] present in test data point [True]
167 Text feature [induced] present in test data point [True]
168 Text feature [catalytic] present in test data point [True]
169 Text feature [considered] present in test data point [True]
171 Text feature [interact] present in test data point [True]
173 Text feature [experimental] present in test data point [True]
174 Text feature [thought] present in test data point [True]
175 Text feature [presented] present in test data point [True]
176 Text feature [cohort] present in test data point [True]
178 Text feature [depletion] present in test data point [True]
180 Text feature [rabbit] present in test data point [True]
181 Text feature [sigma] present in test data point [True]
182 Text feature [26] present in test data point [True]
183 Text feature [following] present in test data point [True]
184 Text feature [pathogenesis] present in test data point [True]
185 Text feature [investigated] present in test data point [True]
188 Text feature [contribute] present in test data point [True]
191 Text feature [linker] present in test data point [True]
193 Text feature [observation] present in test data point [True]
195 Text feature [stably] present in test data point [True]
197 Text feature [lacking] present in test data point [True]
198 Text feature [compound] present in test data point [True]
199 Text feature [finally] present in test data point [True]
201 Text feature [s3] present in test data point [True]
206 Text feature [accumulation] present in test data point [True]
208 Text feature [nuclear] present in test data point [True]
211 Text feature [open] present in test data point [True]
212 Text feature [bound] present in test data point [True]
214 Text feature [conserved] present in test data point [True]
215 Text feature [failure] present in test data point [True]
219 Text feature [tumorigenesis] present in test data point [True]
220 Text feature [near] present in test data point [True]
222 Text feature [little] present in test data point [True]
223 Text feature [fetal] present in test data point [True]
224 Text feature [major] present in test data point [True]
225 Text feature [nacl] present in test data point [True]
226 Text feature [substrate] present in test data point [True]
227 Text feature [consequences] present in test data point [True]
228 Text feature [differences] present in test data point [True]
230 Text feature [predicted] present in test data point [True]
231 Text feature [rate] present in test data point [True]
232 Text feature [disorders] present in test data point [True]
233 Text feature [lack] present in test data point [True]
235 Text feature [relatively] present in test data point [True]
236 Text feature [functional] present in test data point [True]
237 Text feature [deleted] present in test data point [True]
238 Text feature [inc] present in test data point [True]
241 Text feature [tris] present in test data point [True]
242 Text feature [co] present in test data point [True]
243 Text feature [classes] present in test data point [True]
...

245 Text feature [intensity] present in test data point [True]
246 Text feature [heterozygous] present in test data point [True]
248 Text feature [locus] present in test data point [True]
249 Text feature [complex] present in test data point [True]
250 Text feature [primarily] present in test data point [True]
251 Text feature [caused] present in test data point [True]
252 Text feature [isolated] present in test data point [True]
254 Text feature [culture] present in test data point [True]
255 Text feature [absence] present in test data point [True]
256 Text feature [1996] present in test data point [True]
257 Text feature [control] present in test data point [True]
258 Text feature [grown] present in test data point [True]
261 Text feature [suggest] present in test data point [True]
264 Text feature [extracts] present in test data point [True]
265 Text feature [s6] present in test data point [True]
268 Text feature [sds] present in test data point [True]
269 Text feature [procedures] present in test data point [True]
271 Text feature [require] present in test data point [True]
273 Text feature [indicating] present in test data point [True]
274 Text feature [represented] present in test data point [True]
276 Text feature [note] present in test data point [True]
279 Text feature [standard] present in test data point [True]
280 Text feature [chromosome] present in test data point [True]
282 Text feature [act] present in test data point [True]
283 Text feature [significant] present in test data point [True]
286 Text feature [regulates] present in test data point [True]
287 Text feature [blotting] present in test data point [True]
291 Text feature [specifically] present in test data point [True]
292 Text feature [green] present in test data point [True]
294 Text feature [completely] present in test data point [True]
295 Text feature [lower] present in test data point [True]
297 Text feature [individuals] present in test data point [True]
302 Text feature [tissue] present in test data point [True]
303 Text feature [poor] present in test data point [True]
305 Text feature [linked] present in test data point [True]
306 Text feature [purified] present in test data point [True]
308 Text feature [alterations] present in test data point [True]
309 Text feature [elevated] present in test data point [True]
310 Text feature [since] present in test data point [True]
312 Text feature [defective] present in test data point [True]
314 Text feature [indeed] present in test data point [True]
317 Text feature [large] present in test data point [True]
318 Text feature [product] present in test data point [True]
321 Text feature [bind] present in test data point [True]
324 Text feature [representative] present in test data point [True]
325 Text feature [indicated] present in test data point [True]
327 Text feature [collected] present in test data point [True]
329 Text feature [play] present in test data point [True]
330 Text feature [red] present in test data point [True]
332 Text feature [exome] present in test data point [True]
333 Text feature [family] present in test data point [True]
336 Text feature [old] present in test data point [True]
340 Text feature [age] present in test data point [True]
341 Text feature [suggested] present in test data point [True]
342 Text feature [products] present in test data point [True]
345 Text feature [discovery] present in test data point [True]
346 Text feature [pattern] present in test data point [True]
349 Text feature [protein] present in test data point [True]
352 Text feature [directly] present in test data point [True]
353 Text feature [used] present in test data point [True]
355 Text feature [5a] present in test data point [True]
357 Text feature [described] present in test data point [True]
359 Text feature [approximately] present in test data point [True]
361 Text feature [medium] present in test data point [True]
363 Text feature [leads] present in test data point [True]
364 Text feature [mice] present in test data point [True]
365 Text feature [day] present in test data point [True]
366 Text feature [listed] present in test data point [True]
367 Text feature [band] present in test data point [True]
368 Text feature [respectively] present in test data point [True]
371 Text feature [deletions] present in test data point [True]
372 Text feature [experiments] present in test data point [True]
374 Text feature [free] present in test data point [True]
378 Text feature [materials] present in test data point [True]
380 Text feature [15] present in test data point [True]
381 Text feature [dimer] present in test data point [True]
382 Text feature [mutants] present in test data point [True]

```

384 Text feature [demonstrate] present in test data point [True]
387 Text feature [dominant] present in test data point [True]
388 Text feature [lysed] present in test data point [True]
390 Text feature [36] present in test data point [True]
392 Text feature [2b] present in test data point [True]
394 Text feature [pten] present in test data point [True]
396 Text feature [detailed] present in test data point [True]
397 Text feature [genetic] present in test data point [True]
399 Text feature [seen] present in test data point [True]
402 Text feature [ph] present in test data point [True]
403 Text feature [figure] present in test data point [True]
404 Text feature [current] present in test data point [True]
408 Text feature [sites] present in test data point [True]
410 Text feature [harvested] present in test data point [True]
411 Text feature [able] present in test data point [True]
412 Text feature [strong] present in test data point [True]
413 Text feature [detect] present in test data point [True]
415 Text feature [normal] present in test data point [True]
416 Text feature [dna] present in test data point [True]
417 Text feature [groups] present in test data point [True]
418 Text feature [endometrial] present in test data point [True]
419 Text feature [article] present in test data point [True]
420 Text feature [p110] present in test data point [True]
424 Text feature [disease] present in test data point [True]
425 Text feature [overexpressed] present in test data point [True]
430 Text feature [inhibited] present in test data point [True]
431 Text feature [sufficient] present in test data point [True]
432 Text feature [test] present in test data point [True]
436 Text feature [cannot] present in test data point [True]
438 Text feature [point] present in test data point [True]
440 Text feature [2011] present in test data point [True]
441 Text feature [responsible] present in test data point [True]
443 Text feature [necessary] present in test data point [True]
444 Text feature [2007] present in test data point [True]
446 Text feature [fbs] present in test data point [True]
447 Text feature [40] present in test data point [True]
448 Text feature [complete] present in test data point [True]
449 Text feature [despite] present in test data point [True]
453 Text feature [domains] present in test data point [True]
454 Text feature [3b] present in test data point [True]
459 Text feature [exons] present in test data point [True]
462 Text feature [mutated] present in test data point [True]
463 Text feature [effect] present in test data point [True]
464 Text feature [generated] present in test data point [True]
472 Text feature [12] present in test data point [True]
475 Text feature [according] present in test data point [True]
476 Text feature [constructs] present in test data point [True]
477 Text feature [2012] present in test data point [True]
478 Text feature [assayed] present in test data point [True]
479 Text feature [subset] present in test data point [True]
481 Text feature [latter] present in test data point [True]
482 Text feature [shown] present in test data point [True]
483 Text feature [alone] present in test data point [True]
484 Text feature [substrates] present in test data point [True]
485 Text feature [research] present in test data point [True]
488 Text feature [domain] present in test data point [True]
490 Text feature [wt] present in test data point [True]
493 Text feature [phenotype] present in test data point [True]
495 Text feature [component] present in test data point [True]
498 Text feature [thus] present in test data point [True]
Out of the top 500 features 253 are present in query point

```

4.3.3.2. For Incorrectly classified point

In [103]:

```

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

```

```
get_impfeature_names(indices[0],
test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation']
.iloc[test_point_index], no_feature)
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0216 0.0852 0.0048 0.0836 0.0296 0.0081 0.7618 0.0035 0.0018]]

Actual Class : 7

```
-----
11 Text feature [thyroid] present in test data point [True]
25 Text feature [activation] present in test data point [True]
26 Text feature [codon] present in test data point [True]
36 Text feature [increased] present in test data point [True]
42 Text feature [cos] present in test data point [True]
46 Text feature [overexpression] present in test data point [True]
48 Text feature [advanced] present in test data point [True]
51 Text feature [akt] present in test data point [True]
61 Text feature [pathways] present in test data point [True]
62 Text feature [cross] present in test data point [True]
67 Text feature [downstream] present in test data point [True]
77 Text feature [activated] present in test data point [True]
80 Text feature [increase] present in test data point [True]
86 Text feature [expressing] present in test data point [True]
88 Text feature [transforming] present in test data point [True]
90 Text feature [signaling] present in test data point [True]
107 Text feature [upstream] present in test data point [True]
111 Text feature [phosphorylated] present in test data point [True]
115 Text feature [phosphate] present in test data point [True]
116 Text feature [nm] present in test data point [True]
119 Text feature [strong] present in test data point [True]
122 Text feature [codons] present in test data point [True]
124 Text feature [oncogenic] present in test data point [True]
125 Text feature [malignant] present in test data point [True]
134 Text feature [formed] present in test data point [True]
139 Text feature [compared] present in test data point [True]
141 Text feature [form] present in test data point [True]
143 Text feature [mutants] present in test data point [True]
152 Text feature [dl] present in test data point [True]
157 Text feature [exposed] present in test data point [True]
159 Text feature [94] present in test data point [True]
163 Text feature [equivalent] present in test data point [True]
166 Text feature [fold] present in test data point [True]
168 Text feature [domain] present in test data point [True]
171 Text feature [useful] present in test data point [True]
172 Text feature [together] present in test data point [True]
174 Text feature [immunohistochemistry] present in test data point [True]
175 Text feature [000] present in test data point [True]
182 Text feature [event] present in test data point [True]
186 Text feature [phosphorylation] present in test data point [True]
191 Text feature [possibility] present in test data point [True]
202 Text feature [ligand] present in test data point [True]
203 Text feature [epithelial] present in test data point [True]
207 Text feature [exhibited] present in test data point [True]
217 Text feature [lesions] present in test data point [True]
218 Text feature [generated] present in test data point [True]
227 Text feature [factor] present in test data point [True]
233 Text feature [nih] present in test data point [True]
234 Text feature [cyclin] present in test data point [True]
236 Text feature [lane] present in test data point [True]
237 Text feature [per] present in test data point [True]
240 Text feature [high] present in test data point [True]
249 Text feature [500] present in test data point [True]
251 Text feature [lead] present in test data point [True]
257 Text feature [1a] present in test data point [True]
258 Text feature [contained] present in test data point [True]
259 Text feature [mechanism] present in test data point [True]
260 Text feature [carcinomas] present in test data point [True]
261 Text feature [download] present in test data point [True]
262 Text feature [adjacent] present in test data point [True]
263 Text feature [vector] present in test data point [True]
272 Text feature [oncogene] present in test data point [True]
274 Text feature [numbers] present in test data point [True]
275 Text feature [commonly] present in test data point [True]
278 Text feature [membrane] present in test data point [True]
279 Text feature [suggest] present in test data point [True]
281 Text feature [probe] present in test data point [True]
283 Text feature [resulting] present in test data point [True]
```

```

284 Text feature [clone] present in test data point [True]
286 Text feature [log] present in test data point [True]
288 Text feature [free] present in test data point [True]
290 Text feature [represented] present in test data point [True]
295 Text feature [structures] present in test data point [True]
300 Text feature [agar] present in test data point [True]
305 Text feature [serum] present in test data point [True]
308 Text feature [constitutively] present in test data point [True]
310 Text feature [structure] present in test data point [True]
311 Text feature [75] present in test data point [True]
312 Text feature [regulated] present in test data point [True]
316 Text feature [importance] present in test data point [True]
320 Text feature [carried] present in test data point [True]
321 Text feature [total] present in test data point [True]
326 Text feature [since] present in test data point [True]
327 Text feature [occur] present in test data point [True]
331 Text feature [conducted] present in test data point [True]
336 Text feature [fig] present in test data point [True]
337 Text feature [recently] present in test data point [True]
341 Text feature [university] present in test data point [True]
345 Text feature [similar] present in test data point [True]
346 Text feature [cancers] present in test data point [True]
347 Text feature [includes] present in test data point [True]
349 Text feature [frequent] present in test data point [True]
355 Text feature [previously] present in test data point [True]
356 Text feature [immunoprecipitated] present in test data point [True]
358 Text feature [software] present in test data point [True]
363 Text feature [increases] present in test data point [True]
370 Text feature [serine] present in test data point [True]
372 Text feature [12] present in test data point [True]
373 Text feature [made] present in test data point [True]
375 Text feature [transformation] present in test data point [True]
376 Text feature [found] present in test data point [True]
377 Text feature [potential] present in test data point [True]
378 Text feature [transduction] present in test data point [True]
379 Text feature [mouse] present in test data point [True]
386 Text feature [showed] present in test data point [True]
391 Text feature [decrease] present in test data point [True]
393 Text feature [spectrum] present in test data point [True]
394 Text feature [presence] present in test data point [True]
395 Text feature [signal] present in test data point [True]
399 Text feature [factors] present in test data point [True]
400 Text feature [parallel] present in test data point [True]
403 Text feature [residues] present in test data point [True]
408 Text feature [crystal] present in test data point [True]
411 Text feature [cultured] present in test data point [True]
416 Text feature [beads] present in test data point [True]
424 Text feature [exposure] present in test data point [True]
425 Text feature [obtained] present in test data point [True]
428 Text feature [appear] present in test data point [True]
431 Text feature [kit] present in test data point [True]
433 Text feature [24] present in test data point [True]
437 Text feature [inhibition] present in test data point [True]
442 Text feature [induced] present in test data point [True]
443 Text feature [examined] present in test data point [True]
450 Text feature [transient] present in test data point [True]
453 Text feature [canonical] present in test data point [True]
454 Text feature [addition] present in test data point [True]
455 Text feature [density] present in test data point [True]
457 Text feature [primarily] present in test data point [True]
462 Text feature [produced] present in test data point [True]
466 Text feature [non] present in test data point [True]
472 Text feature [lysed] present in test data point [True]
474 Text feature [carcinoma] present in test data point [True]
476 Text feature [extracts] present in test data point [True]
479 Text feature [later] present in test data point [True]
480 Text feature [prior] present in test data point [True]
482 Text feature [threonine] present in test data point [True]
487 Text feature [colony] present in test data point [True]
491 Text feature [five] present in test data point [True]
493 Text feature [fusion] present in test data point [True]
494 Text feature [area] present in test data point [True]
Out of the top 500 features 140 are present in query point

```

4.5 Random Forest Classifier

4.5.1. Hyper paramter tuning (With One hot Encoding)

In [104]:

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

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

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

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

# 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 : 1.201430845073338
for n_estimators = 100 and max depth = 10
Log Loss : 1.1784615237751113
for n_estimators = 200 and max depth = 5
Log Loss : 1.190829868140809
for n_estimators = 200 and max depth = 10
Log Loss : 1.170846385500613
for n_estimators = 500 and max depth = 5
Log Loss : 1.1826336301257254
for n_estimators = 500 and max depth = 10
Log Loss : 1.1687346025846104
for n_estimators = 1000 and max depth = 5
Log Loss : 1.1794638530131654
for n_estimators = 1000 and max depth = 10
Log Loss : 1.1660062367793034
for n_estimators = 2000 and max depth = 5
Log Loss : 1.1808116720972777
for n_estimators = 2000 and max depth = 10
Log Loss : 1.1670230838872564
For values of best estimator = 1000 The train log loss is: 0.5730534008903525
For values of best estimator = 1000 The cross validation log loss is: 1.1660062367793034
For values of best estimator = 1000 The test log loss is: 1.1825216759139787

```

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

In [105]:

```

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

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

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

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

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)

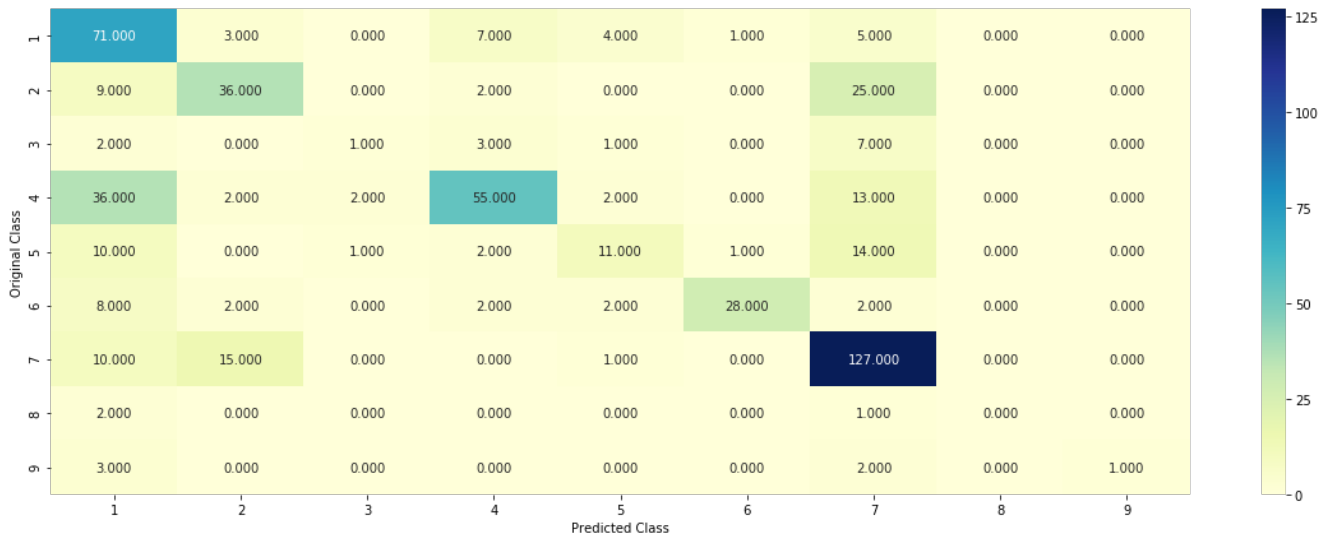
```

```

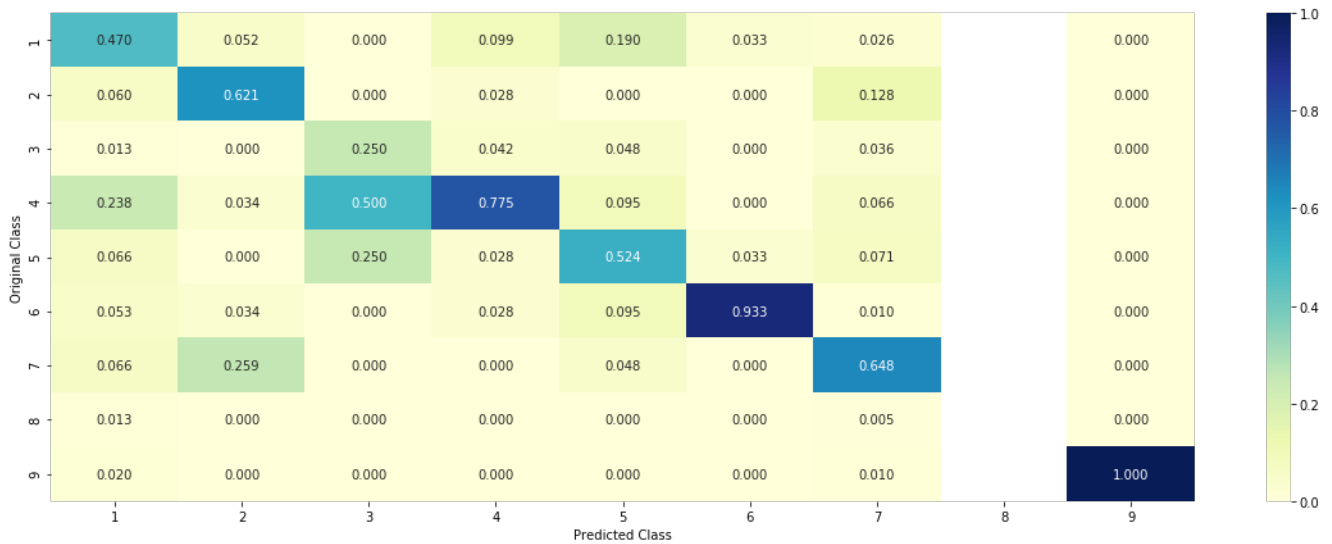
Log loss : 1.1660062367793034
Number of mis-classified points : 0.37969924812030076
Confusion matrix

```

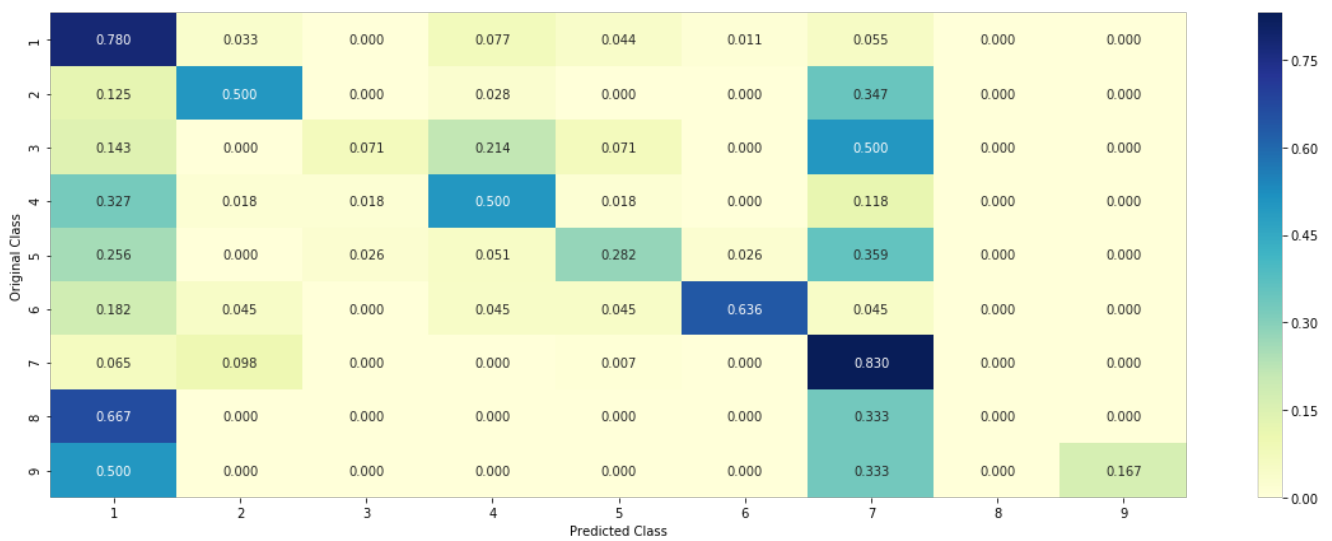
CONFUSION MATRIX



Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



4.5.3. Feature Importance

4.5.3.1. Correctly Classified point

In [106]:

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

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

Predicted Class : 7

Predicted Class Probabilities: [[0.1328 0.1717 0.0245 0.2005 0.0592 0.0527 0.3388 0.0098 0.01]]

Actual Class : 1

```
-----
0 Text feature [kinase] present in test data point [True]
1 Text feature [activating] present in test data point [True]
2 Text feature [activation] present in test data point [True]
4 Text feature [tyrosine] present in test data point [True]
5 Text feature [phosphorylation] present in test data point [True]
7 Text feature [activated] present in test data point [True]
9 Text feature [inhibitors] present in test data point [True]
13 Text feature [function] present in test data point [True]
15 Text feature [treatment] present in test data point [True]
16 Text feature [oncogenic] present in test data point [True]
17 Text feature [loss] present in test data point [True]
20 Text feature [constitutive] present in test data point [True]
21 Text feature [inhibitor] present in test data point [True]
22 Text feature [erk] present in test data point [True]
23 Text feature [signaling] present in test data point [True]
26 Text feature [receptor] present in test data point [True]
30 Text feature [therapy] present in test data point [True]
31 Text feature [therapeutic] present in test data point [True]
34 Text feature [pten] present in test data point [True]
36 Text feature [drug] present in test data point [True]
37 Text feature [expression] present in test data point [True]
38 Text feature [cells] present in test data point [True]
39 Text feature [activate] present in test data point [True]
43 Text feature [growth] present in test data point [True]
44 Text feature [trials] present in test data point [True]
46 Text feature [akt] present in test data point [True]
47 Text feature [variants] present in test data point [True]
48 Text feature [kinases] present in test data point [True]
49 Text feature [protein] present in test data point [True]
52 Text feature [constitutively] present in test data point [True]
54 Text feature [cell] present in test data point [True]
58 Text feature [functional] present in test data point [True]
59 Text feature [potential] present in test data point [True]
60 Text feature [inhibition] present in test data point [True]
65 Text feature [sensitivity] present in test data point [True]
67 Text feature [downstream] present in test data point [True]
68 Text feature [resistance] present in test data point [True]
69 Text feature [predicted] present in test data point [True]
70 Text feature [clinical] present in test data point [True]
72 Text feature [phosphatase] present in test data point [True]
73 Text feature [response] present in test data point [True]
75 Text feature [proteins] present in test data point [True]
76 Text feature [pathway] present in test data point [True]
78 Text feature [transfected] present in test data point [True]
79 Text feature [treated] present in test data point [True]
80 Text feature [human] present in test data point [True]
81 Text feature [factor] present in test data point [True]
82 Text feature [dna] present in test data point [True]
83 Text feature [carry] present in test data point [True]
85 Text feature [pathway] present in test data point [True]
```

```

85 Text feature [patients] present in test data point [True]
86 Text feature [expressing] present in test data point [True]
88 Text feature [transformation] present in test data point [True]
90 Text feature [affected] present in test data point [True]
91 Text feature [affect] present in test data point [True]
95 Text feature [expected] present in test data point [True]
97 Text feature [tagged] present in test data point [True]
98 Text feature [inhibited] present in test data point [True]
Out of the top 100 features 57 are present in query point

```

4.5.3.2. Inorrectly Classified point

In [107]:

```

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

```

```

Predicted Class : 7
Predicted Class Probabilities: [[0.1727 0.1577 0.0248 0.1803 0.0638 0.0547 0.3219 0.0106 0.0135]]
Actual Class : 7

```

```

-----
0 Text feature [kinase] present in test data point [True]
2 Text feature [activation] present in test data point [True]
5 Text feature [phosphorylation] present in test data point [True]
7 Text feature [activated] present in test data point [True]
13 Text feature [function] present in test data point [True]
16 Text feature [oncogenic] present in test data point [True]
17 Text feature [loss] present in test data point [True]
23 Text feature [signaling] present in test data point [True]
26 Text feature [receptor] present in test data point [True]
31 Text feature [therapeutic] present in test data point [True]
37 Text feature [expression] present in test data point [True]
38 Text feature [cells] present in test data point [True]
43 Text feature [growth] present in test data point [True]
46 Text feature [akt] present in test data point [True]
48 Text feature [kinases] present in test data point [True]
49 Text feature [protein] present in test data point [True]
52 Text feature [constitutively] present in test data point [True]
54 Text feature [cell] present in test data point [True]
58 Text feature [functional] present in test data point [True]
59 Text feature [potential] present in test data point [True]
60 Text feature [inhibition] present in test data point [True]
64 Text feature [transforming] present in test data point [True]
67 Text feature [downstream] present in test data point [True]
73 Text feature [response] present in test data point [True]
75 Text feature [proteins] present in test data point [True]
76 Text feature [pathway] present in test data point [True]
78 Text feature [transfected] present in test data point [True]
80 Text feature [human] present in test data point [True]
81 Text feature [factor] present in test data point [True]
82 Text feature [dna] present in test data point [True]
84 Text feature [oncogene] present in test data point [True]
86 Text feature [expressing] present in test data point [True]
88 Text feature [transformation] present in test data point [True]
91 Text feature [affect] present in test data point [True]
95 Text feature [expected] present in test data point [True]
97 Text feature [tagged] present in test data point [True]
Out of the top 100 features 36 are present in query point

```

4.5.3. Hyper paramter tuning (With Response Coding)

In [108]:

```

# -----
# 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 : 2.1435078707316584
for n_estimators = 10 and max depth = 3
Log Loss : 1.71794739847947
for n_estimators = 10 and max depth = 5
Log Loss : 1.4592693250292035
for n_estimators = 10 and max depth = 10
Log Loss : 1.6603329368618667
for n_estimators = 50 and max depth = 2
Log Loss : 1.7377277864367193
for n_estimators = 50 and max depth = 3
Log Loss : 1.42024496811053
for n_estimators = 50 and max depth = 5
Log Loss : 1.3262257208653996
for n_estimators = 50 and max depth = 10
Log Loss : 1.6242232985469247
for n_estimators = 100 and max depth = 2
Log Loss : 1.5865112668543828
for n_estimators = 100 and max depth = 3
Log Loss : 1.4233903619279067
for n_estimators = 100 and max depth = 5
Log Loss : 1.2516985817635862
for n_estimators = 100 and max depth = 10
Log Loss : 1.7135002102121826
for n_estimators = 200 and max depth = 2
Log Loss : 1.601326221388322
for n_estimators = 200 and max depth = 3
Log Loss : 1.4115232958853732
for n_estimators = 200 and max depth = 5
Log Loss : 1.3026842715148657
for n_estimators = 200 and max depth = 10
Log Loss : 1.7049197843899575
for n_estimators = 500 and max depth = 2
Log Loss : 1.6520713435262695
for n_estimators = 500 and max depth = 3
Log Loss : 1.4745257278202597
for n_estimators = 500 and max depth = 5
Log Loss : 1.3199526198149358
for n_estimators = 500 and max depth = 10
Log Loss : 1.694972026834648
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6332163279836642
for n_estimators = 1000 and max depth = 3
Log Loss : 1.466804765821512
for n_estimators = 1000 and max depth = 5
Log Loss : 1.3316796052151383
for n_estimators = 1000 and max depth = 10
Log Loss : 1.670054736470742
For values of best alpha = 100 The train log loss is: 0.059596876418074145
For values of best alpha = 100 The cross validation log loss is: 1.2516985817635862
For values of best alpha = 100 The test log loss is: 1.3403002913968824
```

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

In [109]:

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

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

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

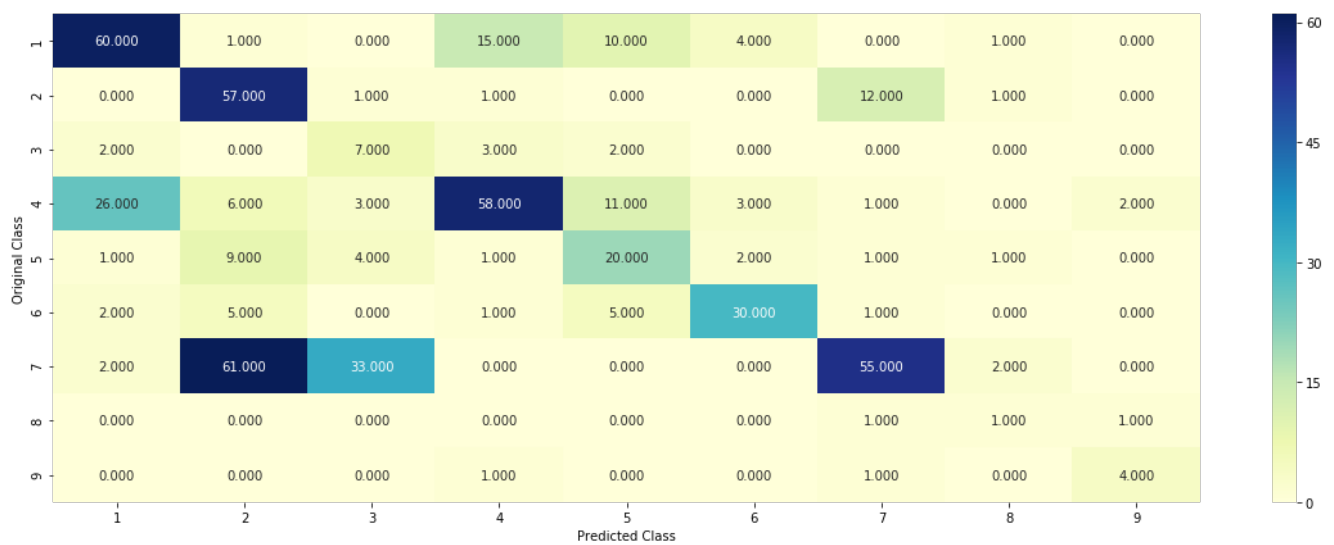
# -----
# video link: https://www.appliedaiaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
# -----

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

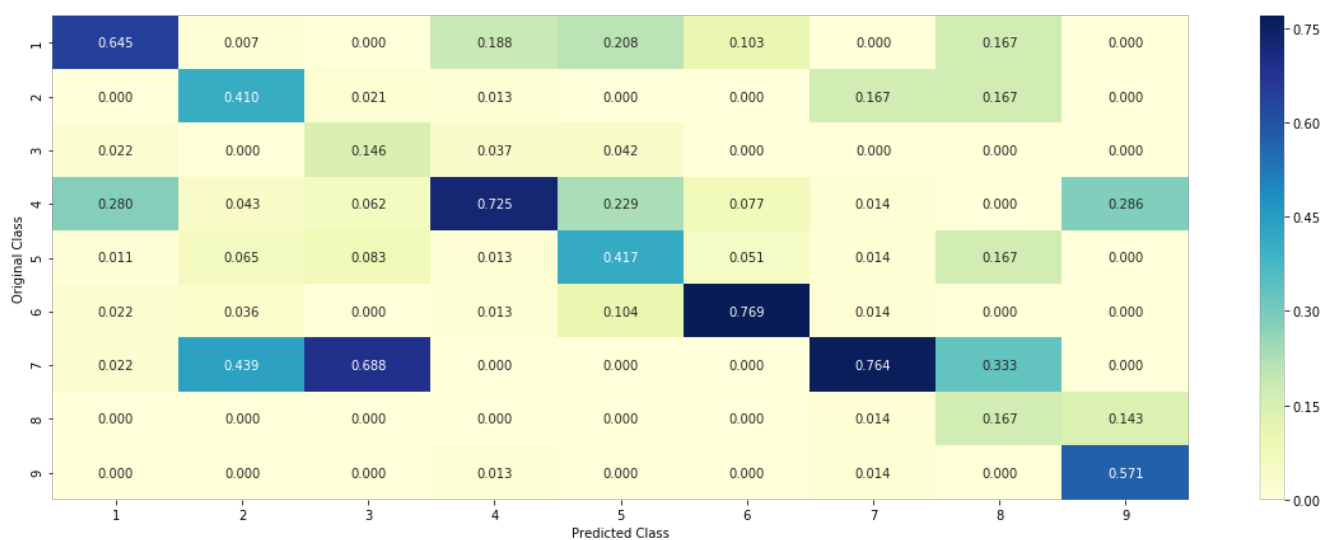
Log loss : 1.2516985817635862

Number of mis-classified points : 0.45112781954887216

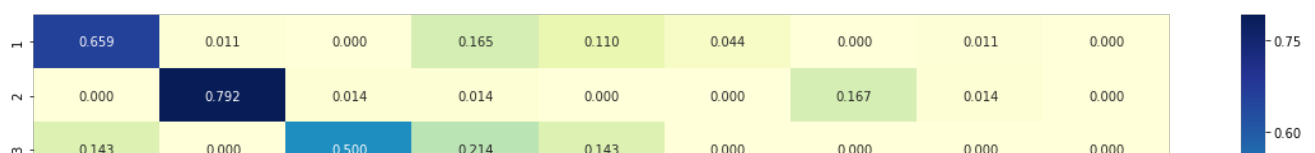
----- Confusion matrix -----

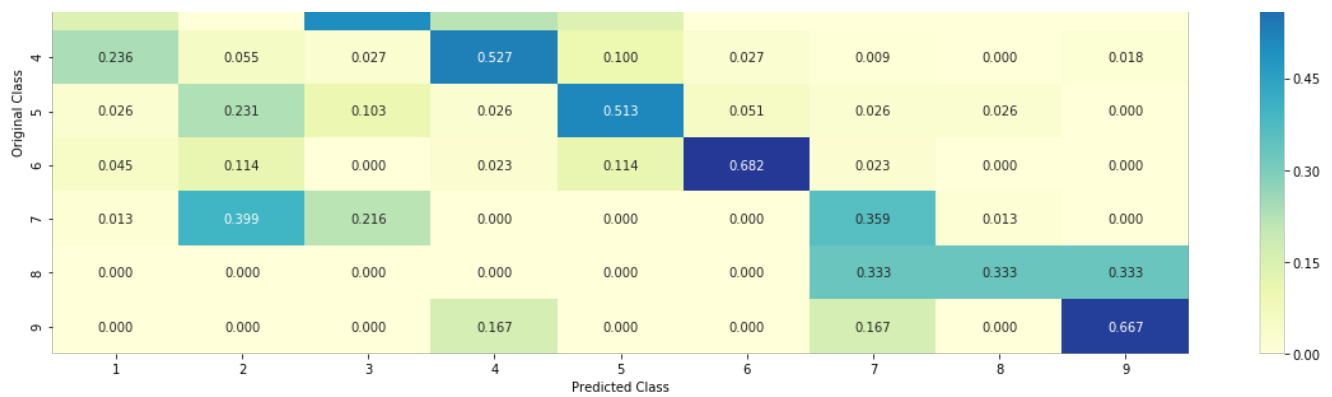


----- Precision matrix (Column Sum=1) -----



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





4.5.5. Feature Importance

4.5.5.1. Correctly Classified point

In [110]:

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

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

Predicted Class : 4

Predicted Class Probabilities: [[0.1895 0.0357 0.0924 0.5398 0.0362 0.0409 0.018 0.0292 0.0183]]

Actual Class : 1

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

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

4.5.5.2. Incorrectly Classified point

In [111]:

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

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0293 0.1425 0.2763 0.0311 0.0347 0.0517 0.3872 0.0276 0.0196]]
Actual Class : 7
-----
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Variation is important feature
Text is important feature
Text is important feature
Gene is important feature
Text is important feature
Gene is important feature
Text is important feature
Text is important feature
Variation is important feature
Gene is important feature
Gene is important feature
Text is important feature
Gene is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Text is important feature
Text is important feature
Text is important feature
Gene is important feature
Gene is important feature
```

4.7 Stack the models

4.7.1 testing with hyper parameter tuning

In [112]:

```
# read more about SGDClassifier() at http://scikit-
learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_i
ter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0
=0.0, power t=0.5,
```

```

# class_weight=None, warm_start=False, average=False, n_iter=None)

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

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

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

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

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

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

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

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

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

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

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

sig_clf1.fit(train_x_onehotCoding, train_y)
print("Logistic Regression : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf1.predict_proba(cv_x_onehotCoding))))
sig_clf2.fit(train_x_onehotCoding, train_y)
print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf2.predict_proba(cv_x_onehotCoding))))

```

```

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

```

Logistic Regression : Log Loss: 0.97
Support vector machines : Log Loss: 1.94
Naive Bayes : Log Loss: 1.14

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: 1.524
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.082
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 0.979
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.044

4.7.2 testing the model with the best hyper parameters

In [113]:

```

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

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

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

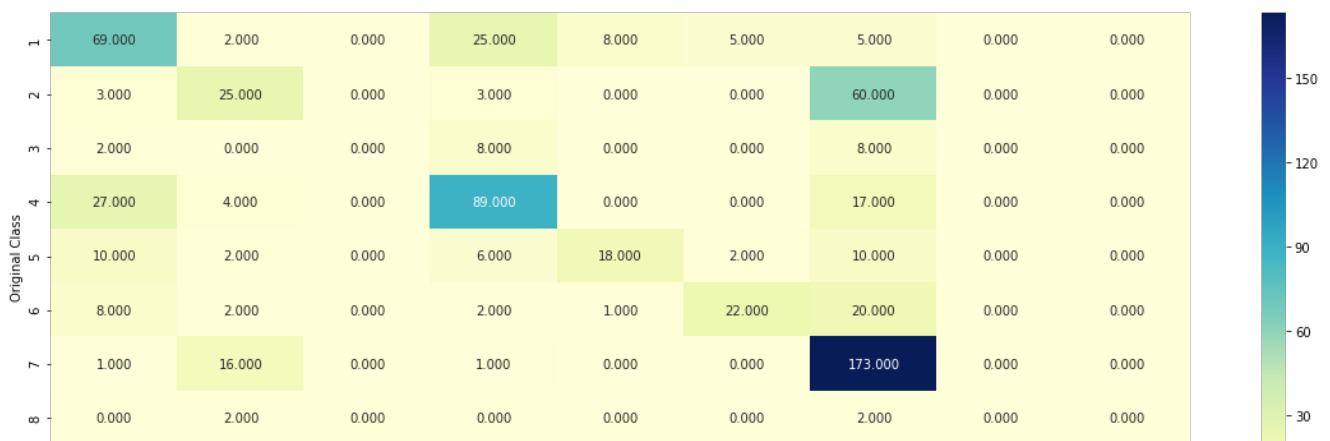
log_error = log_loss(test_y, scf.predict_proba(test_x_onehotCoding_LR))
print("Log loss (test) on the stacking classifier :",log_error)

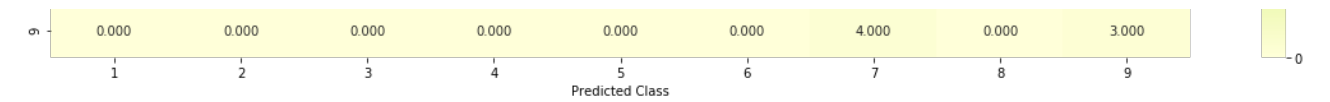
print("Number of missclassified point :", np.count_nonzero((scf.predict(test_x_onehotCoding_LR)-
test_y))/test_y.shape[0])
plot_confusion_matrix(test_y=test_y, predict_y=scf.predict(test_x_onehotCoding_LR))

```

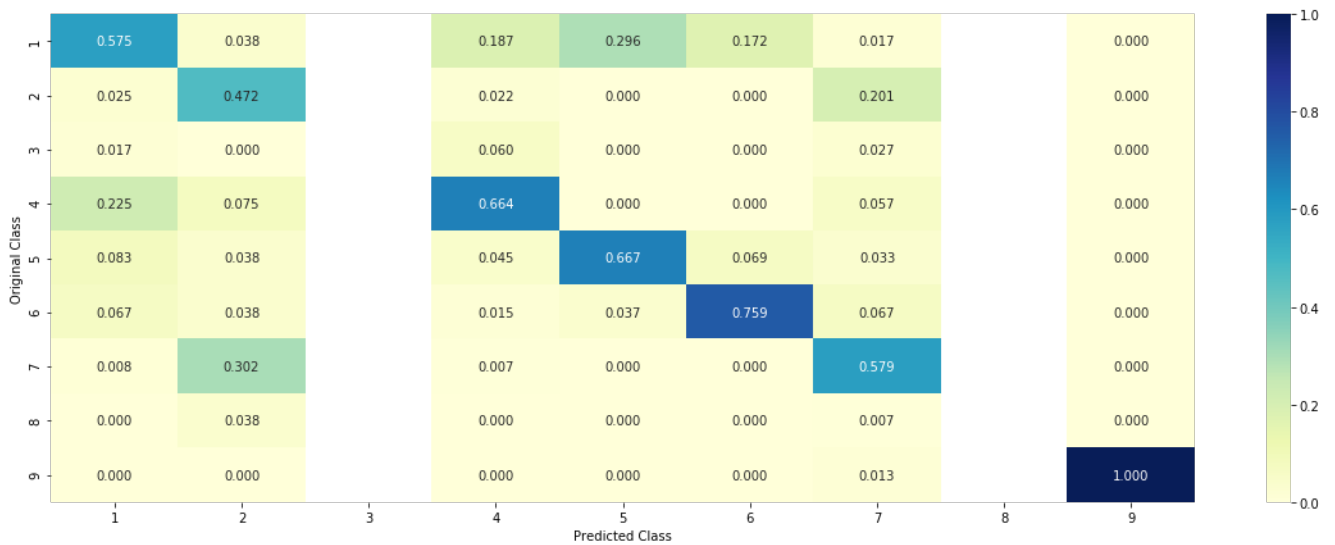
Log loss (train) on the stacking classifier : 1.0573489766384603
Log loss (CV) on the stacking classifier : 1.0763899785391848
Log loss (test) on the stacking classifier : 1.2068895117125609
Number of missclassified point : 0.4

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

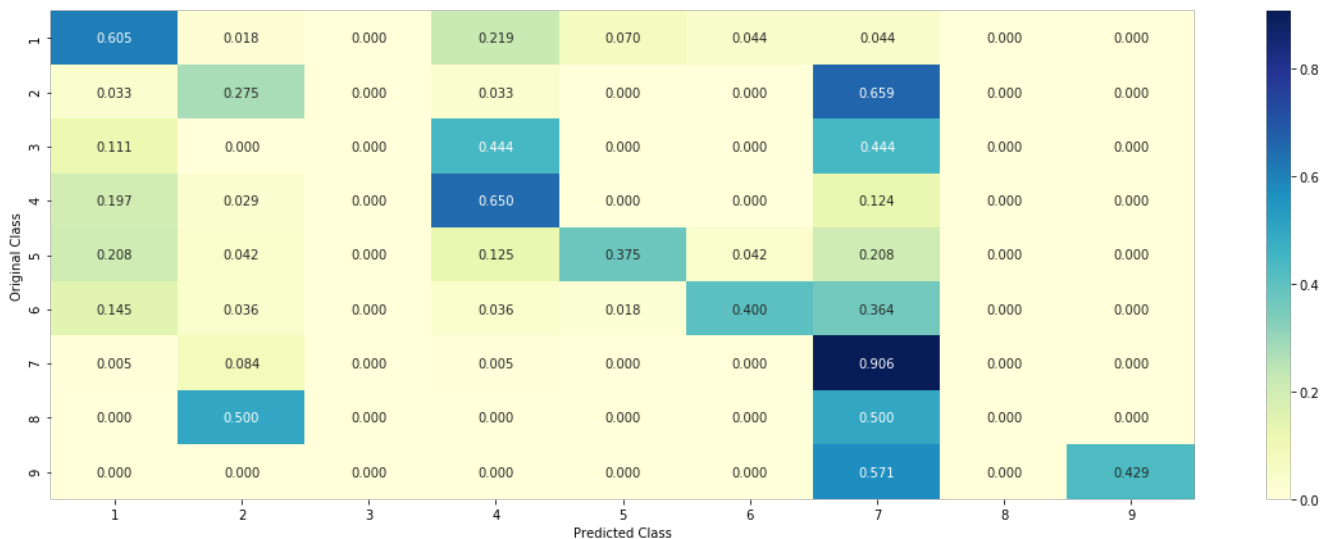




Precision matrix (Column Sum=1)



Recall matrix (Row sum=1)



4.7.3 Maximum Voting classifier

In [114]:

```
#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))
```

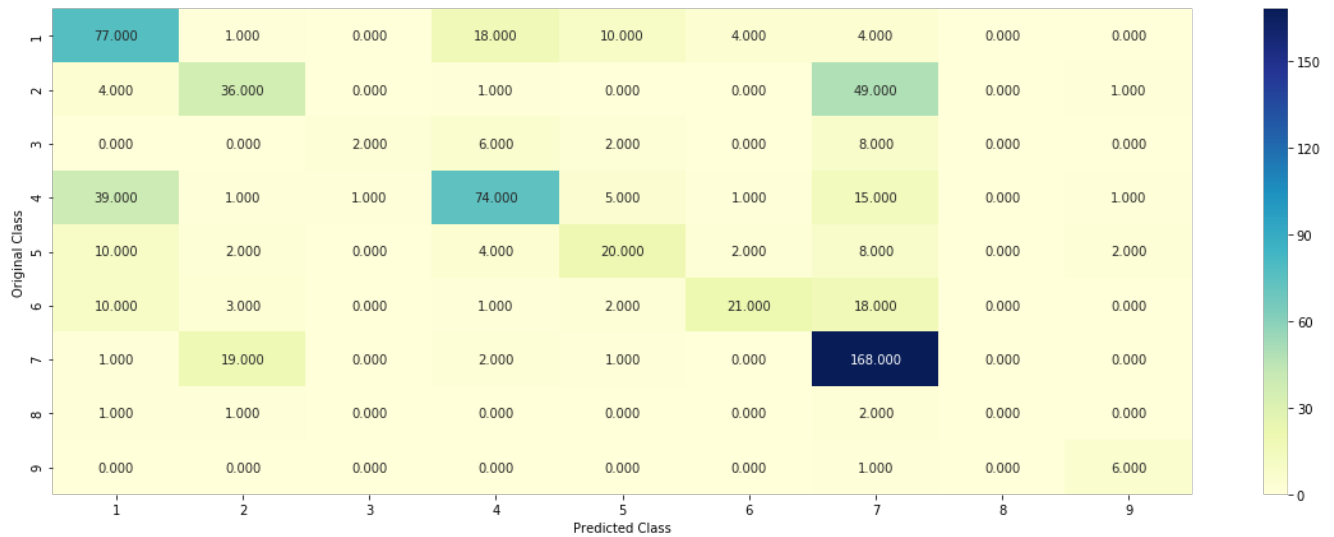
Log loss (train) on the VotingClassifier : 1.0888914138823633

Log loss (CV) on the VotingClassifier : 1.1557909492206915

Log loss (test) on the VotingClassifier : 1.2147180697456577

Number of missclassified point : 0.3924812030075188

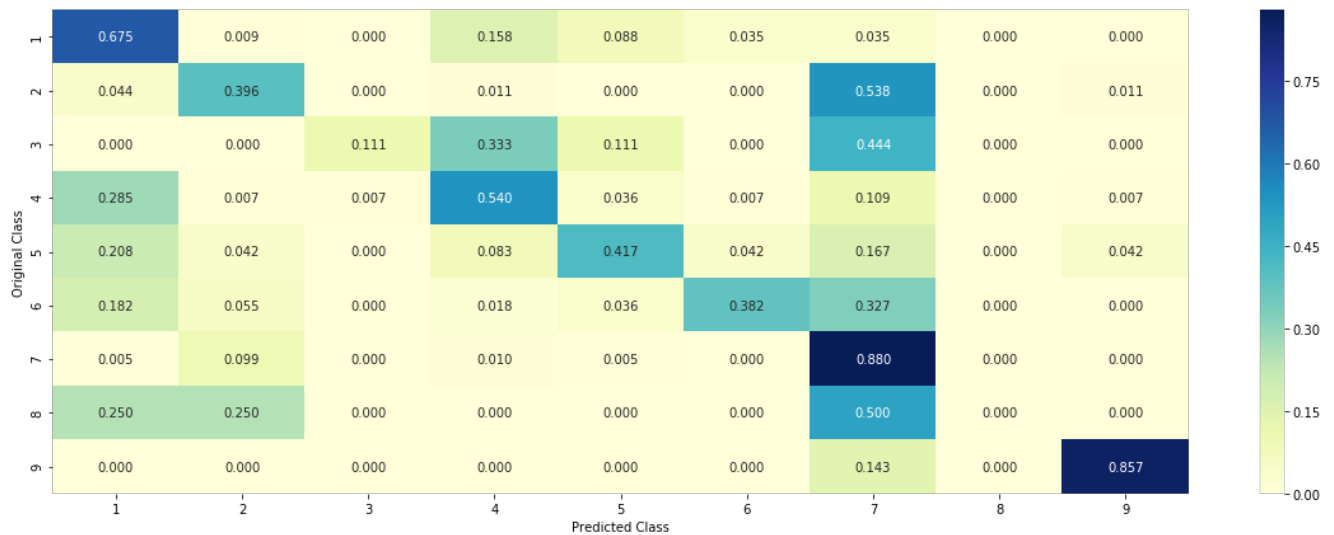
----- Confusion matrix -----



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



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



5. Results

In [283]:

```
from prettytable import PrettyTable

table = PrettyTable()
table.field_names = ["Model","Train loss","CV loss","Test loss","Error"]
table.add_row(["Naive Bayes",1.07,1.13,1.21,37.59])
table.add_row(["KNN",0.68,1.03,1.04,33.64])
table.add_row(["Logistic Regression-Class Balance",0.71,0.96,0.94,33.08])
table.add_row(["Logistic Regression-Without Class Balance",0.71,0.95,0.94,32.70])
table.add_row(["LinearSVM",0.92,1.05,1.13,33.08])
table.add_row(["Random Forest",0.57,1.16,1.18,37.96])
table.add_row(["Stacking",1.05,1.07,1.20,40.00])
table.add_row(["Maximum Voting classifier",1.08,1.15,1.21,39.24])

print(table.get_string(title="Results"))
```

Results					
Model	Train loss	CV loss	Test loss	Error	
Naive Bayes	1.07	1.13	1.21	37.59	
KNN	0.68	1.03	1.04	33.64	
Logistic Regression-Class Balance	0.71	0.96	0.94	33.08	
Logistic Regression-Without Class Balance	0.71	0.95	0.94	32.7	
LinearSVM	0.92	1.05	1.13	33.08	
Random Forest	0.57	1.16	1.18	37.96	
Stacking	1.05	1.07	1.2	40.0	
Maximum Voting classifier	1.08	1.15	1.21	39.24	

6. Conclusion

By observing the above results Logistic Regression model has given good results.