# 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=qxXRKVompI8

### 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 conatins 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 are 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 Learing Objectives and Constraints

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

Constraints:

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

## 2.3. Train, CV and Test Datasets

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

# 3. Exploratory Data Analysis

In [4]:

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

from mlxtend.classifier import StackingClassifier

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

## 3.1. Reading Data

### 3.1.1. Reading Gene and Variation Data

In [5]:

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

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

Out[5]:

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

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

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

### 3.1.2. Reading Text Data

In [6]:

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

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

Out[6]:

| | 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 [7]:

```python
# 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 [8]:

```
#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 : 465.46737670000005 seconds
```

In [9]:

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

Out[9]:

| | 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 [10]:

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

Out[10]:

| | 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 [11]:

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

In [12]:

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

Out[12]:

| | 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 (64:20:16)

In [13]:

```python
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 varaible '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
varaible '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 [14]:

```python
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 [15]:

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

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

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


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

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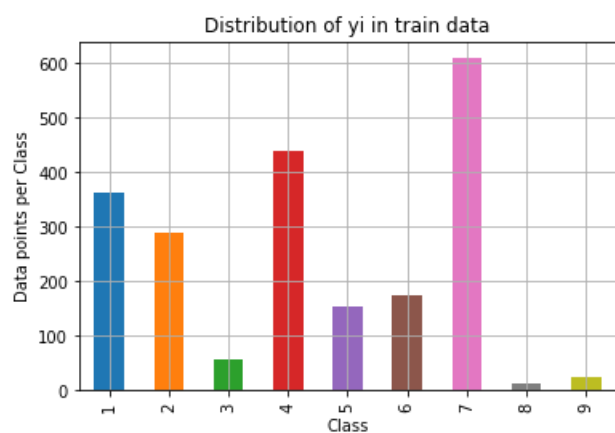
```python
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',test_class_distribution.values[i], '(', np.rou
nd((test_class_distribution.values[i]/test_df.shape[0]*100),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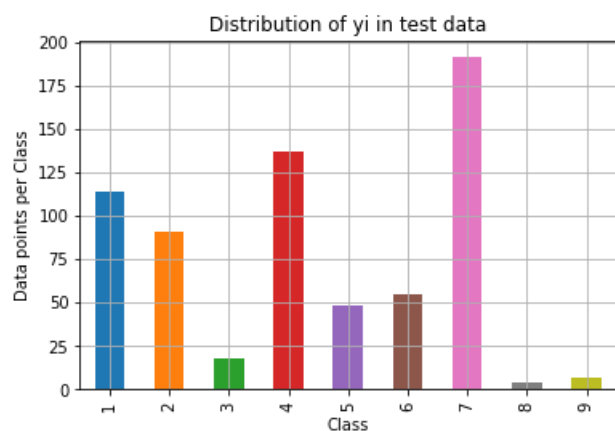


Distribution of yi in train data

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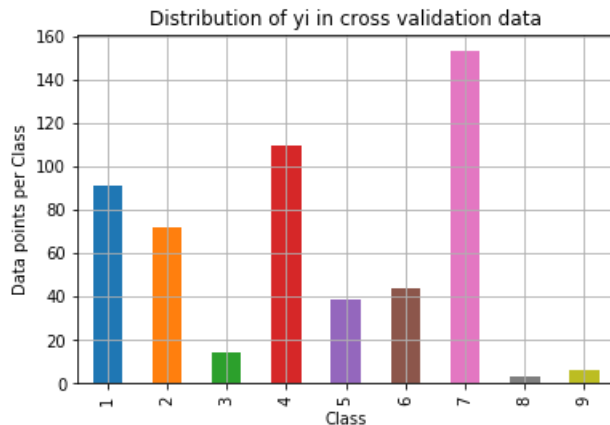


Distribution of yi in test data

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

```
Number of data points in class 9 : 7 ( 1.053 %)
Number of data points in class 8 : 4 ( 0.602 %)
--------------------------------------------------------------------------------
```



Distribution of yi in cross validation data

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

## 3.2 Prediction using a 'Random' Model

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

In [16]:

```python
# 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 corresonds to columns and axis=1 corresponds to rows in two
    diamensional 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 corresonds to columns and axis=1 corresponds to rows in two
    diamensional 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 (Columm 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 [17]:

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

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


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

predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)
```
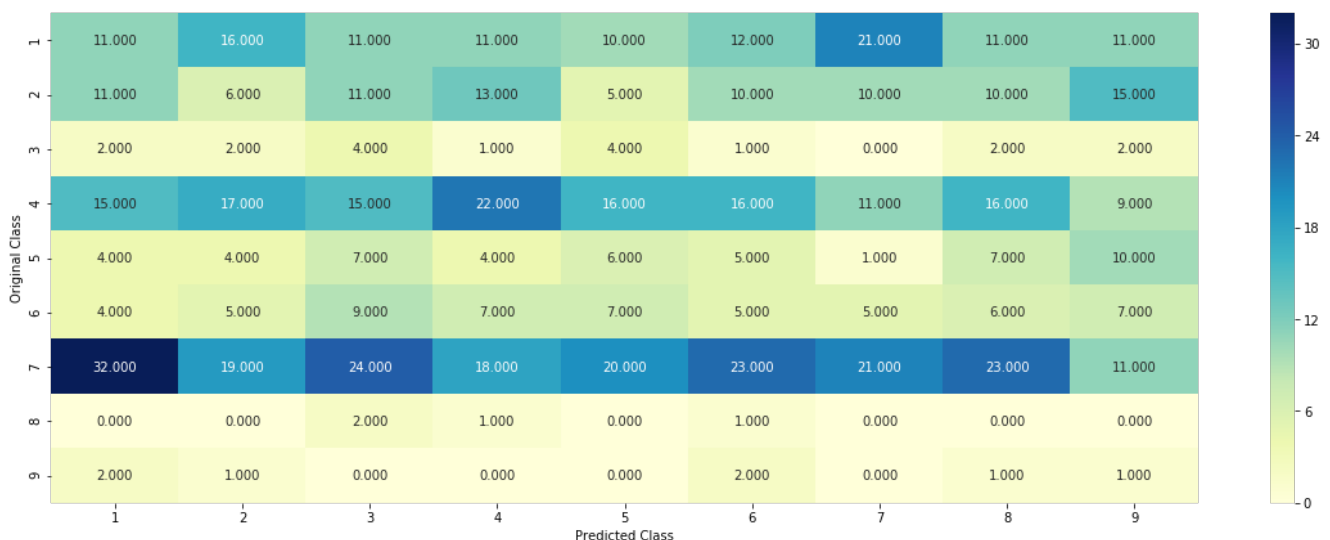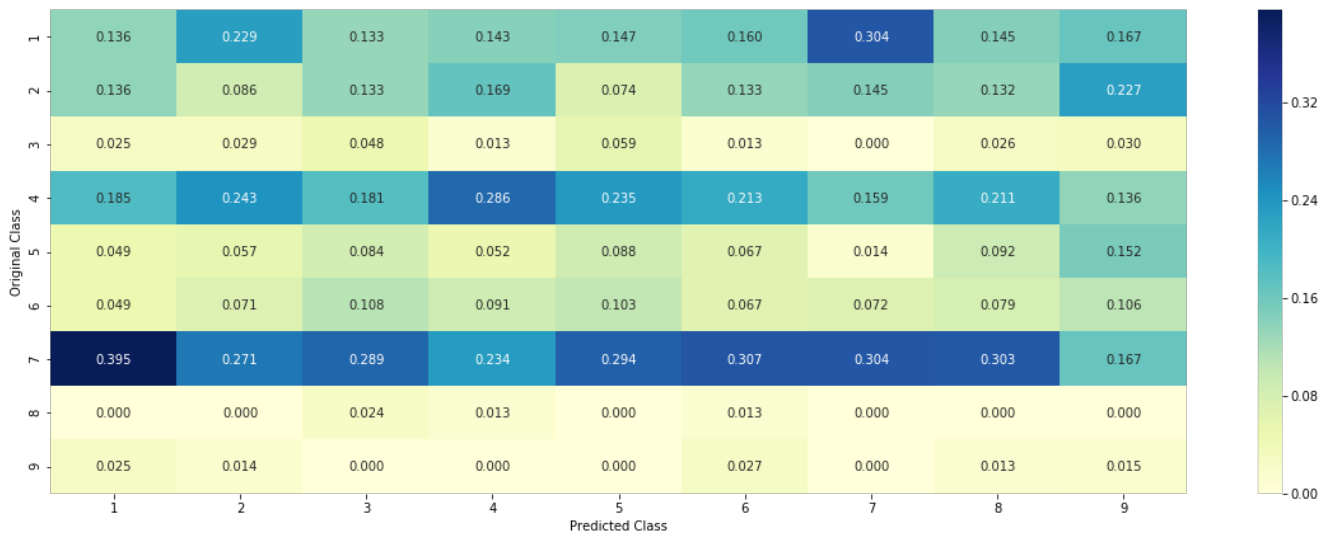
```
Log loss on Cross Validation Data using Random Model 2.4856839415650755
Log loss on Test Data using Random Model 2.407813444978706
-------------------- Confusion matrix --------------------
```

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



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



## 3.3 Univariate Analysis

In [18]:

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

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

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

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

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

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

            # cls_cnt.shape[0](numerator) will contain the number of time that particular feature occured
            # ccured 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, 0.068181818181818177,
    # 0.13636363636363635, 0.25, 0.19318181818181818, 0.03787878787878788, 0.03787878787878788,
    # 0.03787878787878788],
    #      'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224489795918366,
    # 0.27040816326530615, 0.061224489795918366, 0.066326530612244902, 0.051020408163265307, 0.051020408
    # 163265307, 0.056122448979591837],
    #      'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625, 0.068181818181818177,
    # 0.068181818181818177, 0.0625, 0.34659090909090912, 0.0625, 0.056818181818181816],
    #      'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.060606060606060608,
    # 0.078787878787878782, 0.1393939393939394, 0.34545454545454546, 0.060606060606060608,
    # 0.060606060606060608, 0.060606060606060608],
    #      'PTEN': [0.069182389937106917, 0.062893081761006289, 0.069182389937106917,
    # 0.46540880503144655, 0.075471698113207544, 0.062893081761006289, 0.069182389937106917, 0.062893081
    # 761006289, 0.062893081761006289],
    #      'KIT': [0.066225165562913912, 0.25165562913907286, 0.072847682119205295,
    # 0.072847682119205295, 0.066225165562913912, 0.066225165562913912, 0.27152317880794702,
    # 0.066225165562913912, 0.066225165562913912],
    #      'BRAF': [0.066666666666666666, 0.1799999999999999, 0.073333333333333334,
```

```
#            ... [0.000000000000000, 0.199999999999999, 0.0000000000000001,
0.073333333333333334, 0.093333333333333338, 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 da
ta
    gv_fea = []
    # for every feature values in the given data frame we will check if it is there in the train
data then we will add the feature to gv_fea
    # if not we will add [1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9] to gv_fea
    for index, row in df.iterrows():
        if row[feature] in dict(value_count).keys():
            gv_fea.append(gv_dict[row[feature]])
        else:
            gv_fea.append([1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9])
#            gv_fea.append([-1,-1,-1,-1,-1,-1,-1,-1,-1])
    return gv_fea
```

when we caculate the probability of a feature belongs to any particular class, we apply laplace smoothing
- (numerator + 10\*alpha) / (denominator + 90\*alpha)

### 3.2.1 Univariate Analysis on Gene Feature

**Q1.** Gene, What type of feature it is ?

**Ans.** Gene is a categorical variable

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

In [19]:

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

```
Number of Unique Genes : 235
BRCA1     170
TP53      109
EGFR       96
BRCA2      88
PTEN       84
BRAF       61
KIT        60
ALK        48
ERBB2      47
PDGFRA     39
Name: Gene, dtype: int64
```

In [20]:

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

Ans: There are 235 different categories of genes in the train data, and they are distibuted as fol
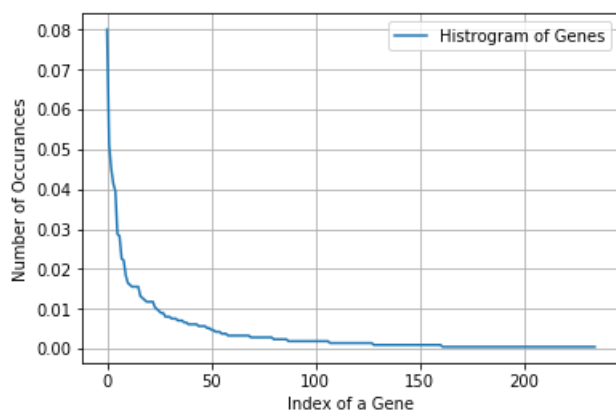lows

In [21]:

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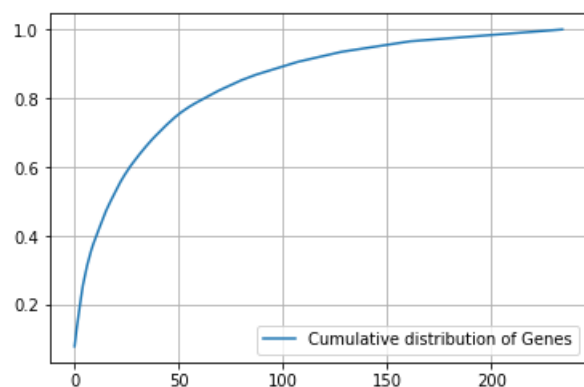
```
plt.grid()
plt.show()
```



In [22]:

```
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 [23]:

```
#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 [24]:

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

train_gene_feature_responseCoding is converted feature using respone coding method. The shape of g
ene feature: (2124, 9)

In [25]:

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

In [26]:

```
train_gene_feature_onehotCoding.shape
```

Out[26]:

(2124, 234)

In [27]:

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

Out[27]:

```
3121     KRAS
424      TP53
1961    NUP93
1133      MET
2502    BRCA1
Name: Gene, dtype: object
```

In [28]:

```
gene_vectorizer.get_feature_names()
```

Out[28]:

```
['abl1',
 'acvr1',
 'ago2',
 'akt1',
 'akt2',
 'akt3',
 'alk',
 'apc',
 'ar',
 'araf',
 'arid1a',
 'arid1b',
 'arid2',
 'arid5b',
 'asxl2',
 'atm',
 'atr',
 'axin1',
 'axl',
 'b2m',
 'bap1',
 'bard1',
 'bcl10',
 'bcl2l11',
 'bcor',
 'braf',
 'brca1',
 'brca2',
 'brd4',
 'brip1',
 'btk',
 'card11',
 'carm1'
```

```
    'calm1',
    'casp8',
    'cbl',
    'ccnd1',
    'ccnd2',
    'ccnd3',
    'ccne1',
    'cdh1',
    'cdk12',
    'cdk4',
    'cdk6',
    'cdk8',
    'cdkn1a',
    'cdkn1b',
    'cdkn2a',
    'cdkn2b',
    'cebpa',
    'chek2',
    'cic',
    'crebbp',
    'ctcf',
    'ctla4',
    'ctnnb1',
    'ddr2',
    'dicer1',
    'dnmt3a',
    'dnmt3b',
    'dusp4',
    'egfr',
    'eif1ax',
    'elf3',
    'ep300',
    'epas1',
    'epcam',
    'erbb2',
    'erbb3',
    'erbb4',
    'ercc2',
    'ercc3',
    'ercc4',
    'erg',
    'errfi1',
    'esr1',
    'etv1',
    'etv6',
    'ewsr1',
    'ezh2',
    'fanca',
    'fancc',
    'fat1',
    'fbxw7',
    'fgfr1',
    'fgfr2',
    'fgfr3',
    'fgfr4',
    'flt3',
    'foxa1',
    'foxl2',
    'foxo1',
    'foxp1',
    'gata3',
    'gli1',
    'gna11',
    'gnaq',
    'gnas',
    'h3f3a',
    'hist1h1c',
    'hla',
    'hnf1a',
    'hras',
    'idh1',
    'idh2',
    'igf1r',
    'ikbke',
    'ikzf1',
    'jak1',
    'jak2',
    'kdm5a',
```

'kdm5a',
'kdm5c',
'kdm6a',
'kdr',
'keap1',
'kit',
'klf4',
'kmt2a',
'kmt2c',
'kmt2d',
'knstrn',
'kras',
'lats1',
'lats2',
'map2k1',
'map2k2',
'map2k4',
'map3k1',
'mapk1',
'mdm4',
'med12',
'mef2b',
'met',
'mga',
'mlh1',
'mpl',
'msh2',
'msh6',
'mtor',
'myc',
'mycn',
'myd88',
'myod1',
'ncor1',
'nf1',
'nf2',
'nfe2l2',
'nfkbia',
'nkx2',
'notch1',
'notch2',
'nras',
'nsd1',
'ntrk1',
'ntrk2',
'ntrk3',
'nup93',
'pax8',
'pdgfra',
'pdgfrb',
'pik3ca',
'pik3cb',
'pik3cd',
'pik3r1',
'pik3r2',
'pik3r3',
'pim1',
'pms1',
'pms2',
'pole',
'ppm1d',
'ppp2r1a',
'prdm1',
'ptch1',
'pten',
'ptpn11',
'ptprd',
'ptprt',
'rab35',
'rac1',
'rad21',
'rad51c',
'rad51d',
'rad54l',
'raf1',
'rara',
'rasa1',
'rb1',

```
 'rb1',
 'rbm10',
 'ret',
 'rhoa',
 'rictor',
 'rit1',
 'ros1',
 'runx1',
 'rybp',
 'sdhb',
 'sdhc',
 'setd2',
 'sf3b1',
 'shoc2',
 'shq1',
 'smad2',
 'smad3',
 'smad4',
 'smarca4',
 'smarcb1',
 'smo',
 'sos1',
 'sox9',
 'spop',
 'src',
 'srsf2',
 'stag2',
 'stat3',
 'stk11',
 'tcf7l2',
 'tert',
 'tet1',
 'tet2',
 'tgfbr1',
 'tgfbr2',
 'tmprss2',
 'tp53',
 'tp53bp1',
 'tsc1',
 'tsc2',
 'u2af1',
 'vegfa',
 'vhl',
 'whsc1',
 'whsc1l1',
 'xpo1',
 'xrcc2',
 'yap1']
```

In [29]:

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

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of g
ene feature: (2124, 234)
```

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

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

In [30]:

```python
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_i
ter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0
```

```
=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

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

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


cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_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.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_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_lo
ss(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, p
redict_y, labels=clf.classes_, eps=1e-15))
```
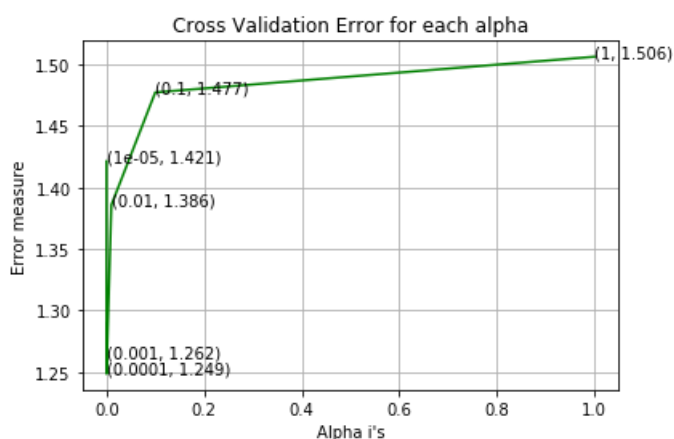
```
For values of alpha =  1e-05 The log loss is: 1.4211787073314368
For values of alpha =  0.0001 The log loss is: 1.248513528391795
For values of alpha =  0.001 The log loss is: 1.2624989157380022
For values of alpha =  0.01 The log loss is: 1.3856857071267357
For values of alpha =  0.1 The log loss is: 1.4771317929059231
For values of alpha =  1 The log loss is: 1.5060939535275977
```



```
For values of best alpha =  0.0001 The train log loss is: 1.0400679900693728
```

```
For values of best alpha =  0.0001 The cross validation log loss is: 1.248513528391795
For values of best alpha =  0.0001 The test log loss is: 1.2080285958178787
```

**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 [31]:

```
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.s
hape[0])*100)
```

```
Q6. How many data points in Test and CV datasets are covered by the  235  genes in train dataset?
Ans
1. In test data 643 out of 665 : 96.69172932330827
2. In cross validation data 511 out of  532 : 96.05263157894737
```

### 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 [32]:

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

```
Number of Unique Variations : 1932
Truncating_Mutations        57
Deletion                    51
Amplification               39
Fusions                     26
Overexpression               4
G12V                         3
Promoter_Hypermethylation    2
V321M                        2
F384L                        2
C618R                        2
Name: Variation, dtype: int64
```

In [33]:

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

```
Ans: There are 1932 different categories of variations in the train data, and they are distibuted
as follows
```

In [34]:

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

```
plt.ylabel( 
plt.legend()
plt.grid()
plt.show()
```



In [35]:

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

```
[0.02683616 0.05084746 0.06920904 ... 0.99905838 0.99952919 1.        ]
```



**Q9.** How to featurize this Variation feature ?

**Ans.** There are two ways we can featurize this variable check out this video:
https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/

1. One hot Encoding
2. Response coding

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

In [36]:

```
# 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 [37]:

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

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

In [38]:

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

In [39]:

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

train_variation_feature_onehotEncoded is converted feature using the onne-hot encoding method. The
shape of Variation feature: (2124, 1957)

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

Let's build a model just like the earlier!

In [40]:

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

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

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

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


cv_log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(train_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")
```
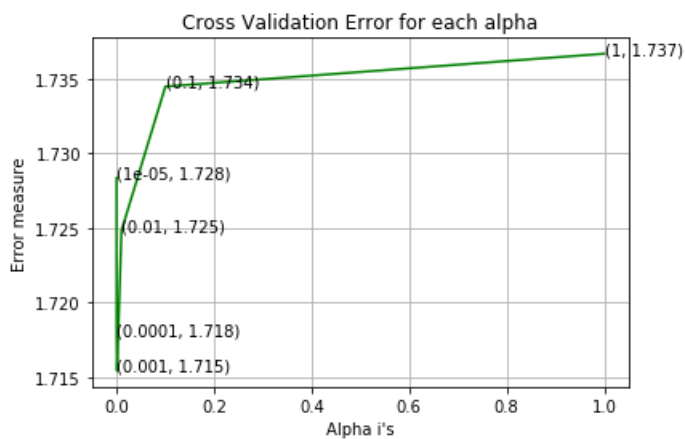
```
plt.show()


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: 1.7283455629543865
For values of alpha =  0.0001 The log loss is: 1.7178444980791865
For values of alpha =  0.001 The log loss is: 1.7154349746204964
For values of alpha =  0.01 The log loss is: 1.7248281031240964
For values of alpha =  0.1 The log loss is: 1.7344713769230264
For values of alpha =  1 The log loss is: 1.7366685264948203
```



```
For values of best alpha =  0.001 The train log loss is: 1.064100766557401
For values of best alpha =  0.001 The cross validation log loss is: 1.7154349746204964
For values of best alpha =  0.001 The test log loss is: 1.7119893151829093
```

**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 [41]:

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

```
Q12. How many data points are covered by total  1932  genes in test and cross validation data
sets?
Ans
1. In test data 68 out of 665 : 10.225563909774436
2. In cross validation data 56 out of  532 : 10.526315789473683
```

### 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 predicitng y_i?
5. Is the text feature stable across train, test and CV datasets?

In [42]:

```python
# 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 [43]:

```python
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 [44]:

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

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

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


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

Total number of unique words in train data : 52960

In [45]:

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

# dict_list[i] is build on i'th  class text data
# total_dict is buid on whole training text data
total_dict = extract_dictionary_paddle(train_df.head(1))
```

```
total_dict = extract_dictionary_paddle(train_df.head(1))


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

In [46]:

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

In [47]:

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

In [48]:

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

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

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

In [49]:

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

In [50]:

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

```
Counter({3: 5292, 4: 3678, 6: 3198, 5: 2692, 8: 2172, 7: 1769, 9: 1603, 10: 1505, 12: 1201, 11: 119
3, 15: 1095, 13: 904, 16: 835, 14: 755, 18: 711, 24: 634, 17: 610, 20: 569, 22: 493, 19: 478, 21:
453, 30: 399, 29: 362, 27: 351, 26: 347, 25: 346, 23: 341, 45: 324, 28: 321, 33: 296, 32: 284, 49:
261, 31: 256, 34: 248, 36: 241, 40: 237, 38: 227, 44: 223, 48: 222, 35: 218, 39: 210, 42: 204, 37:
187, 50: 185, 46: 164, 43: 157, 55: 156, 57: 151, 56: 151, 41: 150, 60: 145, 52: 145, 47: 143, 54:
135, 51: 133, 53: 132, 61: 124, 65: 116, 58: 112, 72: 111, 66: 111, 64: 111, 59: 107, 70: 106, 62:
105, 63: 96, 73: 95, 71: 93, 67: 93, 88: 92, 98: 91, 69: 90, 77: 89, 84: 88, 90: 87, 78: 86, 76:
81, 68: 81, 85: 79, 80: 76, 79: 76, 93: 73, 82: 73, 92: 70, 96: 69, 75: 69, 81: 67, 74: 67, 95: 6
6, 86: 65, 83: 64, 120: 62, 89: 60, 87: 60, 105: 59, 94: 57, 107: 55, 100: 55, 132: 53, 99: 52, 1
19: 51, 102: 51, 113: 50, 104: 49, 115: 48, 110: 47, 101: 47, 126: 46, 147: 45, 108: 45, 118: 44,
111: 44, 137: 43, 91: 43, 122: 42, 112: 42, 109: 42, 106: 41, 97: 41, 134: 40, 143: 39, 135: 39,
103: 39, 144: 38, 136: 38, 114: 38, 130: 37, 127: 37, 116: 37, 139: 36, 138: 36, 121: 36, 124: 35,
123: 34, 182: 33, 170: 33, 158: 33, 152: 33, 153: 32, 148: 32, 128: 32, 125: 32, 176: 31, 146: 30,
140: 30, 129: 30, 117: 30, 164: 29, 150: 29, 141: 29, 131: 29, 220: 28, 145: 28, 205: 27, 190: 27,
180: 27, 169: 27, 156: 27, 192: 26, 177: 26, 174: 26, 166: 26, 165: 26, 154: 26, 149: 26, 218: 25,
172: 25, 159: 25, 221: 24, 194: 24, 185: 24, 184: 24, 162: 24, 157: 24, 215: 23, 212: 23, 204: 23,
```

173: 23, 161: 23, 155: 23, 142: 23, 133: 23, 232: 22, 231: 22, 216: 22, 211: 22, 207: 22, 188: 22,
163: 22, 160: 22, 206: 21, 196: 21, 186: 21, 171: 21, 168: 21, 284: 20, 257: 20, 226: 20, 217: 20,
191: 20, 189: 20, 234: 19, 230: 19, 209: 19, 200: 19, 151: 19, 318: 18, 297: 18, 294: 18, 273: 18,
223: 18, 202: 18, 199: 18, 187: 18, 183: 18, 179: 18, 307: 17, 262: 17, 236: 17, 229: 17, 225: 17,
197: 17, 181: 17, 272: 16, 271: 16, 268: 16, 245: 16, 243: 16, 210: 16, 208: 16, 203: 16, 198: 16,
178: 16, 167: 16, 295: 15, 280: 15, 269: 15, 253: 15, 242: 15, 240: 15, 224: 15, 219: 15, 214: 15,
193: 15, 175: 15, 347: 14, 331: 14, 299: 14, 288: 14, 285: 14, 260: 14, 251: 14, 250: 14, 248: 14,
244: 14, 228: 14, 222: 14, 213: 14, 320: 13, 312: 13, 289: 13, 282: 13, 281: 13, 278: 13, 270: 13,
259: 13, 252: 13, 247: 13, 246: 13, 241: 13, 239: 13, 238: 13, 237: 13, 233: 13, 448: 12, 346: 12,
341: 12, 340: 12, 317: 12, 305: 12, 292: 12, 287: 12, 274: 12, 264: 12, 261: 12, 255: 12, 235: 12,
201: 12, 548: 11, 508: 11, 452: 11, 434: 11, 421: 11, 417: 11, 406: 11, 405: 11, 403: 11, 375: 11,
361: 11, 360: 11, 359: 11, 353: 11, 345: 11, 343: 11, 332: 11, 330: 11, 323: 11, 309: 11, 298: 11,
296: 11, 291: 11, 267: 11, 266: 11, 227: 11, 565: 10, 523: 10, 511: 10, 428: 10, 425: 10, 415: 10,
392: 10, 384: 10, 370: 10, 356: 10, 336: 10, 334: 10, 333: 10, 321: 10, 316: 10, 265: 10, 263: 10,
256: 10, 249: 10, 496: 9, 433: 9, 424: 9, 407: 9, 404: 9, 400: 9, 395: 9, 394: 9, 364: 9, 357: 9,
351: 9, 344: 9, 335: 9, 329: 9, 328: 9, 326: 9, 319: 9, 311: 9, 306: 9, 300: 9, 293: 9, 283: 9,
277: 9, 275: 9, 650: 8, 589: 8, 581: 8, 458: 8, 431: 8, 420: 8, 418: 8, 362: 8, 358: 8, 325: 8,
314: 8, 308: 8, 304: 8, 303: 8, 301: 8, 258: 8, 254: 8, 948: 7, 633: 7, 628: 7, 598: 7, 567: 7,
543: 7, 502: 7, 494: 7, 479: 7, 471: 7, 464: 7, 463: 7, 457: 7, 453: 7, 446: 7, 442: 7, 437: 7,
432: 7, 411: 7, 391: 7, 388: 7, 381: 7, 374: 7, 373: 7, 369: 7, 368: 7, 366: 7, 363: 7, 355: 7,
352: 7, 350: 7, 348: 7, 315: 7, 286: 7, 195: 7, 1374: 6, 924: 6, 915: 6, 877: 6, 840: 6, 826: 6,
761: 6, 728: 6, 679: 6, 674: 6, 664: 6, 638: 6, 637: 6, 608: 6, 594: 6, 542: 6, 540: 6, 535: 6,
500: 6, 490: 6, 475: 6, 470: 6, 468: 6, 465: 6, 456: 6, 444: 6, 439: 6, 436: 6, 426: 6, 423: 6,
412: 6, 410: 6, 409: 6, 401: 6, 398: 6, 396: 6, 390: 6, 385: 6, 383: 6, 380: 6, 378: 6, 377: 6,
372: 6, 349: 6, 342: 6, 338: 6, 302: 6, 290: 6, 1506: 5, 1284: 5, 1253: 5, 1197: 5, 1179: 5, 1037:
5, 1034: 5, 992: 5, 943: 5, 914: 5, 911: 5, 885: 5, 853: 5, 831: 5, 830: 5, 823: 5, 816: 5, 796:
5, 771: 5, 755: 5, 748: 5, 717: 5, 714: 5, 713: 5, 698: 5, 697: 5, 690: 5, 681: 5, 677: 5, 667: 5,
665: 5, 652: 5, 651: 5, 646: 5, 620: 5, 613: 5, 607: 5, 604: 5, 586: 5, 579: 5, 578: 5, 569: 5,
564: 5, 558: 5, 553: 5, 550: 5, 547: 5, 546: 5, 545: 5, 544: 5, 531: 5, 530: 5, 526: 5, 525: 5,
518: 5, 503: 5, 491: 5, 484: 5, 483: 5, 482: 5, 481: 5, 476: 5, 472: 5, 466: 5, 461: 5, 455: 5,
451: 5, 449: 5, 447: 5, 443: 5, 438: 5, 429: 5, 416: 5, 414: 5, 399: 5, 387: 5, 376: 5, 367: 5,
365: 5, 339: 5, 337: 5, 324: 5, 322: 5, 313: 5, 279: 5, 2016: 4, 1686: 4, 1641: 4, 1554: 4, 1446: 4
, 1385: 4, 1370: 4, 1337: 4, 1203: 4, 1170: 4, 1128: 4, 1114: 4, 1082: 4, 1053: 4, 1043: 4, 1008: 4
, 973: 4, 966: 4, 954: 4, 939: 4, 922: 4, 896: 4, 895: 4, 884: 4, 865: 4, 836: 4, 827: 4, 825: 4,
818: 4, 811: 4, 802: 4, 793: 4, 790: 4, 782: 4, 776: 4, 775: 4, 770: 4, 754: 4, 747: 4, 743: 4,
741: 4, 738: 4, 715: 4, 712: 4, 705: 4, 702: 4, 694: 4, 688: 4, 675: 4, 672: 4, 645: 4, 643: 4,
641: 4, 639: 4, 625: 4, 624: 4, 622: 4, 610: 4, 606: 4, 605: 4, 603: 4, 600: 4, 597: 4, 595: 4,
592: 4, 591: 4, 590: 4, 587: 4, 582: 4, 580: 4, 575: 4, 573: 4, 571: 4, 568: 4, 562: 4, 559: 4,
555: 4, 551: 4, 537: 4, 536: 4, 532: 4, 529: 4, 528: 4, 524: 4, 521: 4, 517: 4, 514: 4, 512: 4,
510: 4, 509: 4, 501: 4, 495: 4, 493: 4, 469: 4, 467: 4, 454: 4, 440: 4, 430: 4, 427: 4, 419: 4,
413: 4, 397: 4, 393: 4, 389: 4, 386: 4, 379: 4, 371: 4, 354: 4, 327: 4, 310: 4, 276: 4, 4154: 3,
3731: 3, 3660: 3, 3296: 3, 3269: 3, 3100: 3, 3058: 3, 2412: 3, 2346: 3, 2291: 3, 2179: 3, 2123: 3,
2087: 3, 2069: 3, 2039: 3, 2005: 3, 1951: 3, 1934: 3, 1865: 3, 1860: 3, 1824: 3, 1800: 3, 1787: 3,
1782: 3, 1780: 3, 1759: 3, 1743: 3, 1728: 3, 1714: 3, 1652: 3, 1628: 3, 1614: 3, 1613: 3, 1604: 3,
1591: 3, 1527: 3, 1511: 3, 1505: 3, 1447: 3, 1414: 3, 1406: 3, 1401: 3, 1372: 3, 1363: 3, 1354: 3,
1343: 3, 1334: 3, 1327: 3, 1300: 3, 1296: 3, 1278: 3, 1277: 3, 1275: 3, 1271: 3, 1262: 3, 1258: 3,
1255: 3, 1248: 3, 1208: 3, 1199: 3, 1188: 3, 1187: 3, 1185: 3, 1177: 3, 1172: 3, 1169: 3, 1155: 3,
1152: 3, 1139: 3, 1134: 3, 1118: 3, 1116: 3, 1109: 3, 1102: 3, 1097: 3, 1090: 3, 1088: 3, 1072: 3,
1054: 3, 1047: 3, 1044: 3, 1042: 3, 1038: 3, 1024: 3, 1013: 3, 1007: 3, 1006: 3, 999: 3, 990: 3, 97
1: 3, 968: 3, 961: 3, 958: 3, 956: 3, 953: 3, 942: 3, 940: 3, 935: 3, 933: 3, 916: 3, 909: 3, 902:
3, 886: 3, 876: 3, 870: 3, 869: 3, 861: 3, 855: 3, 854: 3, 852: 3, 851: 3, 849: 3, 848: 3, 845: 3,
835: 3, 820: 3, 808: 3, 805: 3, 803: 3, 798: 3, 792: 3, 788: 3, 784: 3, 774: 3, 767: 3, 764: 3,
763: 3, 760: 3, 758: 3, 757: 3, 752: 3, 749: 3, 736: 3, 734: 3, 733: 3, 732: 3, 731: 3, 730: 3,
727: 3, 726: 3, 720: 3, 718: 3, 704: 3, 700: 3, 689: 3, 685: 3, 682: 3, 680: 3, 678: 3, 671: 3,
662: 3, 661: 3, 656: 3, 655: 3, 649: 3, 640: 3, 629: 3, 626: 3, 623: 3, 621: 3, 614: 3, 609: 3,
602: 3, 574: 3, 566: 3, 561: 3, 557: 3, 556: 3, 552: 3, 533: 3, 522: 3, 513: 3, 507: 3, 506: 3,
504: 3, 499: 3, 497: 3, 492: 3, 488: 3, 478: 3, 474: 3, 473: 3, 462: 3, 459: 3, 450: 3, 445: 3,
435: 3, 402: 3, 382: 3, 9703: 2, 7080: 2, 6418: 2, 5923: 2, 5439: 2, 5340: 2, 5014: 2, 4891: 2, 473
7: 2, 4533: 2, 4461: 2, 4411: 2, 4248: 2, 4213: 2, 4198: 2, 4175: 2, 4095: 2, 4028: 2, 3900: 2, 386
5: 2, 3831: 2, 3829: 2, 3804: 2, 3781: 2, 3734: 2, 3638: 2, 3608: 2, 3606: 2, 3517: 2, 3443: 2, 335
3: 2, 3314: 2, 3268: 2, 3242: 2, 3207: 2, 3187: 2, 3175: 2, 3160: 2, 3046: 2, 3036: 2, 3026: 2, 297
8: 2, 2955: 2, 2951: 2, 2903: 2, 2897: 2, 2848: 2, 2727: 2, 2705: 2, 2668: 2, 2631: 2, 2625: 2, 262
1: 2, 2586: 2, 2581: 2, 2569: 2, 2568: 2, 2516: 2, 2505: 2, 2500: 2, 2476: 2, 2474: 2, 2457: 2, 241
8: 2, 2414: 2, 2409: 2, 2382: 2, 2378: 2, 2343: 2, 2340: 2, 2317: 2, 2245: 2, 2223: 2, 2217: 2, 220
1: 2, 2196: 2, 2168: 2, 2167: 2, 2150: 2, 2149: 2, 2131: 2, 2121: 2, 2104: 2, 2101: 2, 2089: 2, 208
2: 2, 2077: 2, 2053: 2, 2036: 2, 2026: 2, 2017: 2, 2013: 2, 2007: 2, 2002: 2, 1963: 2, 1960: 2, 194
9: 2, 1946: 2, 1942: 2, 1930: 2, 1915: 2, 1900: 2, 1877: 2, 1868: 2, 1854: 2, 1849: 2, 1839: 2, 182
1: 2, 1818: 2, 1789: 2, 1779: 2, 1772: 2, 1768: 2, 1761: 2, 1752: 2, 1750: 2, 1740: 2, 1729: 2, 172
4: 2, 1720: 2, 1709: 2, 1702: 2, 1691: 2, 1690: 2, 1689: 2, 1676: 2, 1672: 2, 1671: 2, 1658: 2, 165
4: 2, 1624: 2, 1619: 2, 1609: 2, 1606: 2, 1598: 2, 1596: 2, 1593: 2, 1584: 2, 1561: 2, 1558: 2, 155
5: 2, 1546: 2, 1544: 2, 1543: 2, 1542: 2, 1540: 2, 1536: 2, 1535: 2, 1517: 2, 1516: 2, 1508: 2, 149
8: 2, 1485: 2, 1484: 2, 1482: 2, 1476: 2, 1472: 2, 1443: 2, 1441: 2, 1440: 2, 1429: 2, 1428: 2, 142
7: 2, 1423: 2, 1408: 2, 1403: 2, 1399: 2, 1398: 2, 1388: 2, 1382: 2, 1362: 2, 1357: 2, 1350: 2, 134
8: 2, 1347: 2, 1344: 2, 1338: 2, 1332: 2, 1331: 2, 1328: 2, 1322: 2, 1320: 2, 1318: 2, 1314: 2, 130
8: 2, 1306: 2, 1303: 2, 1298: 2, 1292: 2, 1289: 2, 1282: 2, 1274: 2, 1268: 2, 1266: 2, 1259: 2, 125
0: 2, 1246: 2, 1244: 2, 1243: 2, 1240: 2, 1239: 2, 1238: 2, 1237: 2, 1234: 2, 1225: 2, 1219: 2, 121
6: 2, 1212: 2, 1210: 2, 1209: 2, 1202: 2, 1196: 2, 1194: 2, 1192: 2, 1186: 2, 1176: 2, 1174: 2, 116

6: 2, 1161: 2, 1154: 2, 1153: 2, 1151: 2, 1145: 2, 1138: 2, 1136: 2, 1131: 2, 1129: 2, 1126: 2, 112
5: 2, 1124: 2, 1120: 2, 1117: 2, 1115: 2, 1111: 2, 1098: 2, 1093: 2, 1091: 2, 1087: 2, 1077: 2, 107
0: 2, 1041: 2, 1040: 2, 1036: 2, 1030: 2, 1029: 2, 1017: 2, 1009: 2, 1002: 2, 997: 2, 995: 2, 994:
2, 993: 2, 988: 2, 987: 2, 984: 2, 980: 2, 970: 2, 960: 2, 959: 2, 952: 2, 946: 2, 937: 2, 930: 2,
928: 2, 921: 2, 920: 2, 912: 2, 904: 2, 900: 2, 898: 2, 893: 2, 892: 2, 891: 2, 890: 2, 883: 2,
878: 2, 875: 2, 873: 2, 871: 2, 866: 2, 863: 2, 862: 2, 846: 2, 844: 2, 843: 2, 841: 2, 837: 2,
833: 2, 832: 2, 822: 2, 813: 2, 812: 2, 809: 2, 806: 2, 804: 2, 800: 2, 799: 2, 797: 2, 791: 2,
781: 2, 780: 2, 773: 2, 772: 2, 768: 2, 766: 2, 765: 2, 759: 2, 753: 2, 750: 2, 744: 2, 742: 2,
737: 2, 729: 2, 725: 2, 724: 2, 721: 2, 716: 2, 708: 2, 707: 2, 706: 2, 703: 2, 701: 2, 699: 2,
696: 2, 695: 2, 693: 2, 687: 2, 686: 2, 683: 2, 676: 2, 673: 2, 670: 2, 668: 2, 663: 2, 660: 2,
659: 2, 658: 2, 654: 2, 647: 2, 644: 2, 634: 2, 632: 2, 631: 2, 627: 2, 618: 2, 616: 2, 615: 2,
612: 2, 611: 2, 596: 2, 593: 2, 588: 2, 585: 2, 584: 2, 577: 2, 576: 2, 572: 2, 563: 2, 560: 2,
539: 2, 538: 2, 527: 2, 520: 2, 519: 2, 516: 2, 515: 2, 489: 2, 487: 2, 486: 2, 477: 2, 460: 2,
441: 2, 408: 2, 155000: 1, 117788: 1, 80538: 1, 68590: 1, 67717: 1, 64828: 1, 64631: 1, 64241: 1,
61638: 1, 57111: 1, 53167: 1, 49045: 1, 48405: 1, 47256: 1, 46419: 1, 44764: 1, 42719: 1, 42544: 1
, 42086: 1, 41978: 1, 40880: 1, 40551: 1, 38916: 1, 38609: 1, 38329: 1, 37646: 1, 37548: 1, 37256:
1, 36526: 1, 36173: 1, 34763: 1, 34156: 1, 33780: 1, 33467: 1, 31866: 1, 31572: 1, 29222: 1, 27955
: 1, 27684: 1, 27162: 1, 26143: 1, 25981: 1, 25932: 1, 25846: 1, 25545: 1, 24848: 1, 24584: 1, 245
60: 1, 24464: 1, 24335: 1, 24194: 1, 23932: 1, 23190: 1, 22909: 1, 22745: 1, 22542: 1, 21930: 1, 2
1660: 1, 21535: 1, 20802: 1, 20732: 1, 20575: 1, 20570: 1, 20380: 1, 20262: 1, 20181: 1, 19841: 1,
19668: 1, 19478: 1, 19192: 1, 19103: 1, 19038: 1, 19034: 1, 18791: 1, 18685: 1, 18586: 1, 18550: 1
, 18470: 1, 18436: 1, 18252: 1, 18190: 1, 18039: 1, 17974: 1, 17737: 1, 17728: 1, 17537: 1, 17494:
1, 17480: 1, 17443: 1, 17429: 1, 17345: 1, 17252: 1, 17082: 1, 17073: 1, 17002: 1, 16842: 1, 16819
: 1, 16757: 1, 16409: 1, 15841: 1, 15815: 1, 15688: 1, 15679: 1, 15621: 1, 15590: 1, 15513: 1, 154
74: 1, 15275: 1, 15172: 1, 15127: 1, 15000: 1, 14909: 1, 14852: 1, 14791: 1, 14601: 1, 14582: 1, 1
4517: 1, 14510: 1, 14448: 1, 14288: 1, 14267: 1, 14261: 1, 14029: 1, 13988: 1, 13781: 1, 13709: 1,
13566: 1, 13554: 1, 13541: 1, 13504: 1, 13448: 1, 13191: 1, 13088: 1, 13064: 1, 13040: 1, 12988: 1
, 12887: 1, 12863: 1, 12735: 1, 12734: 1, 12731: 1, 12707: 1, 12682: 1, 12673: 1, 12669: 1, 12668:
1, 12627: 1, 12531: 1, 12500: 1, 12493: 1, 12484: 1, 12467: 1, 12455: 1, 12449: 1, 12432: 1, 12374
: 1, 12326: 1, 12245: 1, 12226: 1, 12161: 1, 12135: 1, 12120: 1, 12101: 1, 12087: 1, 12080: 1, 118
70: 1, 11857: 1, 11848: 1, 11804: 1, 11762: 1, 11738: 1, 11732: 1, 11606: 1, 11597: 1, 11579: 1, 1
1566: 1, 11563: 1, 11509: 1, 11478: 1, 11235: 1, 11117: 1, 10986: 1, 10933: 1, 10931: 1, 10901: 1,
10884: 1, 10729: 1, 10668: 1, 10632: 1, 10551: 1, 10522: 1, 10499: 1, 10476: 1, 10473: 1, 10462: 1
, 10394: 1, 10366: 1, 10300: 1, 10224: 1, 10187: 1, 10072: 1, 10044: 1, 10025: 1, 9960: 1, 9954: 1
, 9912: 1, 9905: 1, 9849: 1, 9843: 1, 9830: 1, 9827: 1, 9812: 1, 9722: 1, 9655: 1, 9484: 1, 9436: 1
, 9426: 1, 9389: 1, 9387: 1, 9368: 1, 9328: 1, 9302: 1, 9291: 1, 9272: 1, 9173: 1, 9161: 1, 9155: 1
, 9135: 1, 9125: 1, 9106: 1, 9097: 1, 9083: 1, 9078: 1, 9034: 1, 9028: 1, 9026: 1, 9010: 1, 9000: 1
, 8956: 1, 8954: 1, 8921: 1, 8896: 1, 8861: 1, 8805: 1, 8665: 1, 8618: 1, 8512: 1, 8499: 1, 8491: 1
, 8485: 1, 8482: 1, 8467: 1, 8445: 1, 8411: 1, 8395: 1, 8375: 1, 8311: 1, 8303: 1, 8280: 1, 8238: 1
, 8207: 1, 8197: 1, 8193: 1, 8182: 1, 8174: 1, 8135: 1, 8129: 1, 8105: 1, 8070: 1, 8036: 1, 8029: 1
, 7984: 1, 7947: 1, 7936: 1, 7924: 1, 7920: 1, 7901: 1, 7834: 1, 7825: 1, 7788: 1, 7783: 1, 7769: 1
, 7742: 1, 7724: 1, 7716: 1, 7707: 1, 7705: 1, 7701: 1, 7696: 1, 7665: 1, 7599: 1, 7598: 1, 7584: 1
, 7526: 1, 7515: 1, 7482: 1, 7471: 1, 7451: 1, 7445: 1, 7409: 1, 7404: 1, 7375: 1, 7367: 1, 7353: 1
, 7330: 1, 7309: 1, 7293: 1, 7290: 1, 7285: 1, 7269: 1, 7241: 1, 7235: 1, 7231: 1, 7222: 1, 7215: 1
, 7139: 1, 7127: 1, 7115: 1, 7072: 1, 7065: 1, 7058: 1, 7052: 1, 7048: 1, 7045: 1, 7031: 1, 7021: 1
, 7005: 1, 6959: 1, 6953: 1, 6945: 1, 6926: 1, 6922: 1, 6914: 1, 6887: 1, 6884: 1, 6876: 1, 6867: 1
, 6836: 1, 6830: 1, 6810: 1, 6809: 1, 6803: 1, 6797: 1, 6796: 1, 6790: 1, 6741: 1, 6683: 1, 6670: 1
, 6622: 1, 6613: 1, 6606: 1, 6600: 1, 6598: 1, 6590: 1, 6588: 1, 6557: 1, 6551: 1, 6545: 1, 6542: 1
, 6537: 1, 6529: 1, 6524: 1, 6452: 1, 6437: 1, 6432: 1, 6408: 1, 6391: 1, 6390: 1, 6360: 1, 6351: 1
, 6348: 1, 6332: 1, 6326: 1, 6320: 1, 6272: 1, 6268: 1, 6218: 1, 6201: 1, 6182: 1, 6168: 1, 6162: 1
, 6145: 1, 6142: 1, 6136: 1, 6084: 1, 6083: 1, 6060: 1, 6058: 1, 6045: 1, 6044: 1, 6039: 1, 6031: 1
, 6022: 1, 6015: 1, 6013: 1, 5989: 1, 5969: 1, 5961: 1, 5953: 1, 5949: 1, 5926: 1, 5914: 1, 5860: 1
, 5854: 1, 5848: 1, 5843: 1, 5822: 1, 5815: 1, 5806: 1, 5790: 1, 5780: 1, 5773: 1, 5772: 1, 5769: 1
, 5765: 1, 5726: 1, 5719: 1, 5692: 1, 5663: 1, 5662: 1, 5634: 1, 5622: 1, 5619: 1, 5616: 1, 5612: 1
, 5596: 1, 5578: 1, 5569: 1, 5562: 1, 5555: 1, 5511: 1, 5509: 1, 5476: 1, 5473: 1, 5459: 1, 5418: 1
, 5399: 1, 5394: 1, 5366: 1, 5363: 1, 5354: 1, 5343: 1, 5330: 1, 5318: 1, 5301: 1, 5299: 1, 5291: 1
, 5284: 1, 5282: 1, 5279: 1, 5260: 1, 5259: 1, 5255: 1, 5254: 1, 5240: 1, 5237: 1, 5224: 1, 5214: 1
, 5212: 1, 5202: 1, 5189: 1, 5147: 1, 5126: 1, 5116: 1, 5113: 1, 5089: 1, 5060: 1, 5047: 1, 5033: 1
, 5024: 1, 5011: 1, 5009: 1, 4999: 1, 4994: 1, 4990: 1, 4986: 1, 4970: 1, 4940: 1, 4934: 1, 4929: 1
, 4923: 1, 4914: 1, 4898: 1, 4896: 1, 4890: 1, 4869: 1, 4864: 1, 4862: 1, 4849: 1, 4844: 1, 4829: 1
, 4818: 1, 4816: 1, 4815: 1, 4814: 1, 4803: 1, 4800: 1, 4788: 1, 4786: 1, 4776: 1, 4769: 1, 4768: 1
, 4749: 1, 4734: 1, 4727: 1, 4723: 1, 4713: 1, 4710: 1, 4703: 1, 4662: 1, 4658: 1, 4651: 1, 4623: 1
, 4618: 1, 4600: 1, 4589: 1, 4562: 1, 4559: 1, 4554: 1, 4545: 1, 4531: 1, 4528: 1, 4523: 1, 4515: 1
, 4513: 1, 4494: 1, 4493: 1, 4480: 1, 4477: 1, 4468: 1, 4459: 1, 4457: 1, 4455: 1, 4442: 1, 4437: 1
, 4434: 1, 4410: 1, 4402: 1, 4396: 1, 4393: 1, 4390: 1, 4381: 1, 4380: 1, 4368: 1, 4364: 1, 4363: 1
, 4349: 1, 4346: 1, 4337: 1, 4328: 1, 4325: 1, 4316: 1, 4315: 1, 4307: 1, 4306: 1, 4304: 1, 4295: 1
, 4294: 1, 4293: 1, 4291: 1, 4280: 1, 4259: 1, 4249: 1, 4246: 1, 4245: 1, 4231: 1, 4226: 1, 4214: 1
, 4187: 1, 4181: 1, 4179: 1, 4161: 1, 4146: 1, 4145: 1, 4119: 1, 4113: 1, 4105: 1, 4100: 1, 4097: 1
, 4096: 1, 4086: 1, 4084: 1, 4071: 1, 4070: 1, 4058: 1, 4057: 1, 4054: 1, 4050: 1, 4048: 1, 4047: 1
, 4040: 1, 4037: 1, 4036: 1, 4025: 1, 4014: 1, 4011: 1, 4009: 1, 4008: 1, 4002: 1, 3997: 1, 3983: 1
, 3972: 1, 3966: 1, 3956: 1, 3950: 1, 3940: 1, 3939: 1, 3921: 1, 3919: 1, 3909: 1, 3888: 1, 3882: 1
, 3881: 1, 3864: 1, 3861: 1, 3858: 1, 3816: 1, 3815: 1, 3814: 1, 3808: 1, 3792: 1, 3787: 1, 3778: 1
, 3770: 1, 3765: 1, 3764: 1, 3761: 1, 3755: 1, 3749: 1, 3739: 1, 3735: 1, 3732: 1, 3716: 1, 3711: 1
, 3708: 1, 3702: 1, 3699: 1, 3695: 1, 3689: 1, 3687: 1, 3686: 1, 3677: 1, 3676: 1, 3656: 1, 3649: 1
, 3635: 1, 3629: 1, 3626: 1, 3619: 1, 3613: 1, 3611: 1, 3607: 1, 3604: 1, 3598: 1, 3595: 1, 3577: 1
, 3575: 1, 3568: 1, 3564: 1, 3563: 1, 3562: 1, 3553: 1, 3548: 1, 3547: 1, 3543: 1, 3541: 1, 3529: 1
, 3525: 1, 3513: 1, 3510: 1, 3506: 1, 3505: 1, 3503: 1, 3498: 1, 3497: 1, 3493: 1, 3492: 1, 3487: 1

, 3486: 1, 3474: 1, 3463: 1, 3462: 1, 3459: 1, 3456: 1, 3452: 1, 3449: 1, 3448: 1, 3439: 1, 3438: 1
, 3436: 1, 3434: 1, 3433: 1, 3430: 1, 3428: 1, 3426: 1, 3416: 1, 3413: 1, 3404: 1, 3403: 1, 3402: 1
, 3401: 1, 3397: 1, 3389: 1, 3387: 1, 3379: 1, 3377: 1, 3375: 1, 3365: 1, 3364: 1, 3358: 1, 3356: 1
, 3349: 1, 3336: 1, 3329: 1, 3325: 1, 3305: 1, 3303: 1, 3302: 1, 3301: 1, 3292: 1, 3286: 1, 3284: 1
, 3281: 1, 3280: 1, 3271: 1, 3264: 1, 3261: 1, 3256: 1, 3255: 1, 3251: 1, 3249: 1, 3241: 1, 3237: 1
, 3235: 1, 3234: 1, 3233: 1, 3232: 1, 3229: 1, 3228: 1, 3208: 1, 3206: 1, 3204: 1, 3191: 1, 3184: 1
, 3176: 1, 3174: 1, 3169: 1, 3159: 1, 3158: 1, 3146: 1, 3143: 1, 3139: 1, 3134: 1, 3130: 1, 3128: 1
, 3121: 1, 3120: 1, 3114: 1, 3101: 1, 3099: 1, 3094: 1, 3077: 1, 3076: 1, 3073: 1, 3071: 1, 3068: 1
, 3067: 1, 3062: 1, 3059: 1, 3056: 1, 3051: 1, 3044: 1, 3031: 1, 3030: 1, 3018: 1, 3015: 1, 3013: 1
, 3012: 1, 3011: 1, 3007: 1, 3006: 1, 2997: 1, 2993: 1, 2983: 1, 2977: 1, 2976: 1, 2968: 1, 2964: 1
, 2957: 1, 2944: 1, 2943: 1, 2942: 1, 2940: 1, 2937: 1, 2935: 1, 2934: 1, 2932: 1, 2910: 1, 2906: 1
, 2904: 1, 2894: 1, 2893: 1, 2892: 1, 2885: 1, 2877: 1, 2875: 1, 2871: 1, 2867: 1, 2866: 1, 2863: 1
, 2861: 1, 2860: 1, 2854: 1, 2853: 1, 2852: 1, 2849: 1, 2839: 1, 2837: 1, 2829: 1, 2828: 1, 2819: 1
, 2813: 1, 2810: 1, 2806: 1, 2802: 1, 2798: 1, 2793: 1, 2787: 1, 2786: 1, 2772: 1, 2759: 1, 2758: 1
, 2756: 1, 2747: 1, 2746: 1, 2744: 1, 2738: 1, 2737: 1, 2736: 1, 2733: 1, 2719: 1, 2716: 1, 2712: 1
, 2708: 1, 2706: 1, 2704: 1, 2699: 1, 2688: 1, 2686: 1, 2680: 1, 2679: 1, 2674: 1, 2672: 1, 2671: 1
, 2666: 1, 2661: 1, 2658: 1, 2656: 1, 2654: 1, 2652: 1, 2647: 1, 2641: 1, 2639: 1, 2636: 1, 2628: 1
, 2627: 1, 2626: 1, 2622: 1, 2614: 1, 2611: 1, 2607: 1, 2604: 1, 2601: 1, 2597: 1, 2594: 1, 2592: 1
, 2583: 1, 2580: 1, 2578: 1, 2576: 1, 2570: 1, 2567: 1, 2565: 1, 2557: 1, 2553: 1, 2551: 1, 2543: 1
, 2535: 1, 2534: 1, 2531: 1, 2527: 1, 2526: 1, 2523: 1, 2521: 1, 2519: 1, 2510: 1, 2508: 1, 2507: 1
, 2504: 1, 2499: 1, 2497: 1, 2494: 1, 2493: 1, 2492: 1, 2491: 1, 2489: 1, 2475: 1, 2471: 1, 2468: 1
, 2467: 1, 2464: 1, 2462: 1, 2460: 1, 2459: 1, 2456: 1, 2450: 1, 2447: 1, 2442: 1, 2432: 1, 2429: 1
, 2413: 1, 2411: 1, 2407: 1, 2405: 1, 2402: 1, 2401: 1, 2398: 1, 2396: 1, 2395: 1, 2391: 1, 2386: 1
, 2384: 1, 2379: 1, 2376: 1, 2374: 1, 2373: 1, 2371: 1, 2369: 1, 2367: 1, 2363: 1, 2357: 1, 2356: 1
, 2354: 1, 2353: 1, 2341: 1, 2332: 1, 2327: 1, 2323: 1, 2322: 1, 2321: 1, 2313: 1, 2310: 1, 2305: 1
, 2304: 1, 2303: 1, 2300: 1, 2298: 1, 2292: 1, 2288: 1, 2286: 1, 2283: 1, 2281: 1, 2277: 1, 2276: 1
, 2275: 1, 2273: 1, 2270: 1, 2268: 1, 2264: 1, 2259: 1, 2258: 1, 2252: 1, 2246: 1, 2241: 1, 2240: 1
, 2239: 1, 2234: 1, 2230: 1, 2229: 1, 2227: 1, 2222: 1, 2221: 1, 2220: 1, 2219: 1, 2215: 1, 2209: 1
, 2207: 1, 2202: 1, 2192: 1, 2191: 1, 2189: 1, 2188: 1, 2182: 1, 2178: 1, 2176: 1, 2175: 1, 2172: 1
, 2170: 1, 2164: 1, 2162: 1, 2158: 1, 2155: 1, 2154: 1, 2153: 1, 2146: 1, 2145: 1, 2143: 1, 2142: 1
, 2140: 1, 2137: 1, 2135: 1, 2132: 1, 2122: 1, 2120: 1, 2117: 1, 2115: 1, 2114: 1, 2113: 1, 2112: 1
, 2109: 1, 2108: 1, 2106: 1, 2092: 1, 2090: 1, 2080: 1, 2074: 1, 2067: 1, 2065: 1, 2063: 1, 2055: 1
, 2054: 1, 2050: 1, 2048: 1, 2047: 1, 2045: 1, 2042: 1, 2035: 1, 2034: 1, 2032: 1, 2031: 1, 2022: 1
, 2020: 1, 2018: 1, 2010: 1, 2009: 1, 2008: 1, 2001: 1, 2000: 1, 1998: 1, 1997: 1, 1996: 1, 1983: 1
, 1979: 1, 1976: 1, 1974: 1, 1969: 1, 1968: 1, 1966: 1, 1965: 1, 1961: 1, 1956: 1, 1953: 1, 1945: 1
, 1944: 1, 1941: 1, 1931: 1, 1929: 1, 1928: 1, 1920: 1, 1918: 1, 1917: 1, 1916: 1, 1911: 1, 1910: 1
, 1909: 1, 1906: 1, 1905: 1, 1902: 1, 1901: 1, 1895: 1, 1890: 1, 1886: 1, 1885: 1, 1884: 1, 1880: 1
, 1879: 1, 1876: 1, 1875: 1, 1873: 1, 1871: 1, 1869: 1, 1867: 1, 1863: 1, 1862: 1, 1858: 1, 1857: 1
, 1856: 1, 1852: 1, 1850: 1, 1847: 1, 1845: 1, 1844: 1, 1836: 1, 1830: 1, 1828: 1, 1827: 1, 1826: 1
, 1825: 1, 1823: 1, 1819: 1, 1817: 1, 1816: 1, 1812: 1, 1811: 1, 1804: 1, 1803: 1, 1799: 1, 1794: 1
, 1793: 1, 1792: 1, 1791: 1, 1790: 1, 1783: 1, 1778: 1, 1777: 1, 1776: 1, 1775: 1, 1771: 1, 1766: 1
, 1765: 1, 1763: 1, 1760: 1, 1757: 1, 1755: 1, 1754: 1, 1751: 1, 1748: 1, 1741: 1, 1738: 1, 1734: 1
, 1731: 1, 1726: 1, 1723: 1, 1721: 1, 1719: 1, 1718: 1, 1713: 1, 1712: 1, 1708: 1, 1707: 1, 1706: 1
, 1703: 1, 1701: 1, 1699: 1, 1693: 1, 1688: 1, 1687: 1, 1677: 1, 1675: 1, 1673: 1, 1669: 1, 1668: 1
, 1667: 1, 1664: 1, 1662: 1, 1660: 1, 1659: 1, 1656: 1, 1653: 1, 1651: 1, 1648: 1, 1645: 1, 1640: 1
, 1639: 1, 1634: 1, 1633: 1, 1632: 1, 1631: 1, 1620: 1, 1618: 1, 1616: 1, 1612: 1, 1611: 1, 1610: 1
, 1608: 1, 1607: 1, 1605: 1, 1603: 1, 1602: 1, 1595: 1, 1592: 1, 1589: 1, 1588: 1, 1586: 1, 1581: 1
, 1580: 1, 1578: 1, 1577: 1, 1576: 1, 1574: 1, 1570: 1, 1568: 1, 1565: 1, 1560: 1, 1559: 1, 1552: 1
, 1550: 1, 1549: 1, 1548: 1, 1547: 1, 1541: 1, 1539: 1, 1537: 1, 1534: 1, 1532: 1, 1531: 1, 1530: 1
, 1528: 1, 1526: 1, 1524: 1, 1523: 1, 1520: 1, 1519: 1, 1514: 1, 1513: 1, 1512: 1, 1502: 1, 1500: 1
, 1499: 1, 1495: 1, 1491: 1, 1490: 1, 1489: 1, 1488: 1, 1487: 1, 1486: 1, 1483: 1, 1479: 1, 1478: 1
, 1473: 1, 1468: 1, 1461: 1, 1460: 1, 1456: 1, 1449: 1, 1445: 1, 1438: 1, 1437: 1, 1432: 1, 1431: 1
, 1425: 1, 1424: 1, 1422: 1, 1421: 1, 1420: 1, 1419: 1, 1416: 1, 1415: 1, 1412: 1, 1411: 1, 1407: 1
, 1405: 1, 1402: 1, 1396: 1, 1392: 1, 1391: 1, 1390: 1, 1389: 1, 1384: 1, 1383: 1, 1381: 1, 1377: 1
, 1373: 1, 1371: 1, 1368: 1, 1366: 1, 1365: 1, 1364: 1, 1361: 1, 1359: 1, 1358: 1, 1356: 1, 1355: 1
, 1352: 1, 1351: 1, 1346: 1, 1342: 1, 1341: 1, 1335: 1, 1333: 1, 1330: 1, 1329: 1, 1326: 1, 1325: 1
, 1324: 1, 1323: 1, 1319: 1, 1312: 1, 1310: 1, 1304: 1, 1302: 1, 1294: 1, 1293: 1, 1291: 1, 1288: 1
, 1285: 1, 1283: 1, 1281: 1, 1273: 1, 1272: 1, 1270: 1, 1269: 1, 1267: 1, 1265: 1, 1264: 1, 1260: 1
, 1254: 1, 1251: 1, 1247: 1, 1245: 1, 1241: 1, 1235: 1, 1232: 1, 1230: 1, 1229: 1, 1224: 1, 1222: 1
, 1220: 1, 1218: 1, 1217: 1, 1215: 1, 1214: 1, 1211: 1, 1206: 1, 1204: 1, 1201: 1, 1200: 1, 1193: 1
, 1191: 1, 1178: 1, 1175: 1, 1173: 1, 1171: 1, 1168: 1, 1167: 1, 1165: 1, 1162: 1, 1160: 1, 1159: 1
, 1158: 1, 1157: 1, 1156: 1, 1149: 1, 1147: 1, 1146: 1, 1143: 1, 1142: 1, 1141: 1, 1140: 1, 1127: 1
, 1123: 1, 1122: 1, 1119: 1, 1113: 1, 1110: 1, 1108: 1, 1107: 1, 1106: 1, 1105: 1, 1103: 1, 1095: 1
, 1094: 1, 1092: 1, 1089: 1, 1086: 1, 1085: 1, 1084: 1, 1083: 1, 1080: 1, 1079: 1, 1078: 1, 1073: 1
, 1069: 1, 1067: 1, 1066: 1, 1063: 1, 1062: 1, 1061: 1, 1060: 1, 1059: 1, 1058: 1, 1051: 1, 1050: 1
, 1049: 1, 1048: 1, 1045: 1, 1039: 1, 1035: 1, 1033: 1, 1028: 1, 1025: 1, 1022: 1, 1020: 1, 1019: 1
, 1018: 1, 1004: 1, 1003: 1, 1001: 1, 1000: 1, 998: 1, 996: 1, 991: 1, 985: 1, 983: 1, 982: 1, 979:
1, 978: 1, 977: 1, 976: 1, 975: 1, 972: 1, 967: 1, 964: 1, 962: 1, 951: 1, 950: 1, 945: 1, 944: 1,
941: 1, 936: 1, 934: 1, 931: 1, 929: 1, 927: 1, 926: 1, 925: 1, 923: 1, 918: 1, 917: 1, 913: 1,
910: 1, 907: 1, 905: 1, 903: 1, 901: 1, 899: 1, 889: 1, 888: 1, 887: 1, 882: 1, 880: 1, 879: 1,
874: 1, 868: 1, 867: 1, 864: 1, 859: 1, 858: 1, 857: 1, 850: 1, 847: 1, 839: 1, 838: 1, 834: 1,
829: 1, 828: 1, 821: 1, 819: 1, 817: 1, 815: 1, 814: 1, 810: 1, 807: 1, 801: 1, 795: 1, 794: 1,
789: 1, 787: 1, 786: 1, 785: 1, 783: 1, 779: 1, 778: 1, 777: 1, 769: 1, 756: 1, 751: 1, 746: 1,
745: 1, 735: 1, 723: 1, 719: 1, 710: 1, 709: 1, 692: 1, 691: 1, 684: 1, 669: 1, 666: 1, 653: 1,
642: 1, 636: 1, 635: 1, 630: 1, 617: 1, 601: 1, 599: 1, 583: 1, 570: 1, 554: 1, 541: 1, 534: 1,
505: 1, 498: 1, 485: 1, 480: 1})

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

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

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

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


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

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

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


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

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

```
For values of alpha =  1e-05 The log loss is: 1.3700028384794973
For values of alpha =  0.0001 The log loss is: 1.3522429849572446
For values of alpha =  0.001 The log loss is: 1.269035824190863
For values of alpha =  0.01 The log loss is: 1.2920334469908283
For values of alpha =  0.1 The log loss is: 1.4819885933164294
For values of alpha =  1 The log loss is: 1.6717696998108351
```

```
For values of best alpha =  0.001 The train log loss is: 0.7465950275703206
For values of best alpha =  0.001 The cross validation log loss is: 1.269035824190863
For values of best alpha =  0.001 The test log loss is: 1.1708911818712113
```

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

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

In [52]:

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

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

In [53]:

```python
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")
```

```
97.341 % of word of test data appeared in train data
96.611 % of word of Cross Validation appeared in train data
```

# 4. Machine Learning Models

In [54]:

```python
#Data preparation for ML models.

#Misc. functionns for ML models

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

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

In [55]:

```python
def report_log_loss(train_x, train_y, test_x, test_y,  clf):
    clf.fit(train_x, train_y)
```

```
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x, train_y)
    sig_clf_probs = sig_clf.predict_proba(test_x)
    return log_loss(test_y, sig_clf_probs, eps=1e-15)
```

In [56]:

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

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

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

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

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

## Stacking the three types of features

In [57]:

```
# merging gene, variance and text features

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

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

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

test_x_onehotCoding = hstack((test_gene_var_onehotCoding, test_text_feature_onehotCoding)).tocsr()
```

```
test_y = np.array(list(test_df['Class']))

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


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

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

In [58]:

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

```
One hot encoding features :
(number of data points * number of features) in train data =  (2124, 55151)
(number of data points * number of features) in test data =   (665, 55151)
(number of data points * number of features) in cross validation data = (532, 55151)
```

In [59]:

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

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

## 4.1. Base Line Model

### 4.1.1. Naive Bayes

#### 4.1.1.1. Hyper parameter tuning

In [60]:

```
# 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/
# ----------------------
```

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

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

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


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


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

```
for alpha = 1e-05
Log Loss : 1.3162269190737566
for alpha = 0.0001
Log Loss : 1.3139102307631054
for alpha = 0.001
Log Loss : 1.3156114626971793
for alpha = 0.1
Log Loss : 1.2981730002798886
for alpha = 1
Log Loss : 1.3101887444204294
for alpha = 10
Log Loss : 1.4079155398575487
for alpha = 100
Log Loss : 1.3837868289591662
for alpha = 1000
Log Loss : 1.3376611452835814
```

Cross Validation Error for each alpha

```
For values of best alpha =  0.1 The train log loss is: 0.8669233610691601
For values of best alpha =  0.1 The cross validation log loss is: 1.2981730002798886
For values of best alpha =  0.1 The test log loss is: 1.339314694768685
```

**4.1.1.2. Testing the model with best hyper paramters**

In [61]:

```python
# 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 : 1.2981730002798886
Number of missclassified point : 0.4191729323308271
-------------------- Confusion matrix --------------------
```

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



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



### 4.1.1.3. Feature Importance, Correctly classified point

In [62]:

```
test_point_index = 1
no_feature = 100
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
```

```
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding[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 : 2
Predicted Class Probabilities: [[0.0916 0.4596 0.0181 0.1269 0.0423 0.0345 0.2195 0.0043 0.0032]]
Actual Class : 7
--------------------------------------------------
17 Text feature [identified] present in test data point [True]
19 Text feature [including] present in test data point [True]
22 Text feature [molecular] present in test data point [True]
25 Text feature [confirmed] present in test data point [True]
26 Text feature [sequencing] present in test data point [True]
28 Text feature [clinical] present in test data point [True]
30 Text feature [another] present in test data point [True]
31 Text feature [case] present in test data point [True]
32 Text feature [recently] present in test data point [True]
33 Text feature [different] present in test data point [True]
34 Text feature [using] present in test data point [True]
35 Text feature [well] present in test data point [True]
37 Text feature [patient] present in test data point [True]
40 Text feature [harbor] present in test data point [True]
41 Text feature [12] present in test data point [True]
42 Text feature [mutations] present in test data point [True]
43 Text feature [found] present in test data point [True]
44 Text feature [may] present in test data point [True]
45 Text feature [also] present in test data point [True]
46 Text feature [revealed] present in test data point [True]
47 Text feature [15] present in test data point [True]
49 Text feature [potential] present in test data point [True]
50 Text feature [kinase] present in test data point [True]
51 Text feature [cases] present in test data point [True]
54 Text feature [number] present in test data point [True]
55 Text feature [additional] present in test data point [True]
56 Text feature [complete] present in test data point [True]
57 Text feature [specific] present in test data point [True]
59 Text feature [common] present in test data point [True]
60 Text feature [need] present in test data point [True]
61 Text feature [go] present in test data point [True]
62 Text feature [10] present in test data point [True]
63 Text feature [reported] present in test data point [True]
64 Text feature [pcr] present in test data point [True]
65 Text feature [previously] present in test data point [True]
67 Text feature [described] present in test data point [True]
69 Text feature [however] present in test data point [True]
70 Text feature [observed] present in test data point [True]
71 Text feature [respectively] present in test data point [True]
72 Text feature [studies] present in test data point [True]
73 Text feature [similarly] present in test data point [True]
74 Text feature [harboring] present in test data point [True]
75 Text feature [distinct] present in test data point [True]
76 Text feature [one] present in test data point [True]
77 Text feature [similar] present in test data point [True]
78 Text feature [single] present in test data point [True]
80 Text feature [addition] present in test data point [True]
81 Text feature [achieved] present in test data point [True]
82 Text feature [analysis] present in test data point [True]
83 Text feature [three] present in test data point [True]
84 Text feature [performed] present in test data point [True]
85 Text feature [informed] present in test data point [True]
87 Text feature [40] present in test data point [True]
89 Text feature [due] present in test data point [True]
91 Text feature [although] present in test data point [True]
92 Text feature [various] present in test data point [True]
93 Text feature [several] present in test data point [True]
94 Text feature [33] present in test data point [True]
95 Text feature [inhibitor] present in test data point [True]
96 Text feature [two] present in test data point [True]
97 Text feature [tumor] present in test data point [True]
98 Text feature [16] present in test data point [True]
```

```
99 Text feature [known] present in test data point [True]
Out of the top  100  features  63 are present in query point
```

### 4.1.1.4. Feature Importance, Incorrectly classified point

```
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: [[0.0812 0.2141 0.0161 0.113  0.0376 0.0307 0.5005 0.0038 0.0028]]
Actual Class : 7
--------------------------------------------------
16 Text feature [kinase] present in test data point [True]
17 Text feature [downstream] present in test data point [True]
19 Text feature [presence] present in test data point [True]
22 Text feature [activation] present in test data point [True]
23 Text feature [shown] present in test data point [True]
24 Text feature [inhibitor] present in test data point [True]
25 Text feature [well] present in test data point [True]
26 Text feature [contrast] present in test data point [True]
27 Text feature [recently] present in test data point [True]
28 Text feature [cells] present in test data point [True]
29 Text feature [obtained] present in test data point [True]
30 Text feature [expressing] present in test data point [True]
31 Text feature [previously] present in test data point [True]
32 Text feature [growth] present in test data point [True]
34 Text feature [also] present in test data point [True]
35 Text feature [suggest] present in test data point [True]
36 Text feature [cell] present in test data point [True]
37 Text feature [however] present in test data point [True]
38 Text feature [independent] present in test data point [True]
39 Text feature [factor] present in test data point [True]
40 Text feature [compared] present in test data point [True]
41 Text feature [mutations] present in test data point [True]
42 Text feature [found] present in test data point [True]
43 Text feature [higher] present in test data point [True]
44 Text feature [showed] present in test data point [True]
45 Text feature [10] present in test data point [True]
46 Text feature [similar] present in test data point [True]
47 Text feature [treated] present in test data point [True]
48 Text feature [addition] present in test data point [True]
49 Text feature [activated] present in test data point [True]
51 Text feature [may] present in test data point [True]
53 Text feature [constitutive] present in test data point [True]
54 Text feature [studies] present in test data point [True]
55 Text feature [interestingly] present in test data point [True]
56 Text feature [12] present in test data point [True]
57 Text feature [followed] present in test data point [True]
59 Text feature [total] present in test data point [True]
64 Text feature [observed] present in test data point [True]
65 Text feature [mechanism] present in test data point [True]
66 Text feature [using] present in test data point [True]
68 Text feature [approximately] present in test data point [True]
69 Text feature [phosphorylation] present in test data point [True]
70 Text feature [including] present in test data point [True]
71 Text feature [without] present in test data point [True]
73 Text feature [identified] present in test data point [True]
74 Text feature [respectively] present in test data point [True]
75 Text feature [although] present in test data point [True]
76 Text feature [suggests] present in test data point [True]
77 Text feature [reported] present in test data point [True]
79 Text feature [new] present in test data point [True]
80 Text feature [whereas] present in test data point [True]
```

80 Text feature [whereas] present in test data point [True]
82 Text feature [measured] present in test data point [True]
83 Text feature [report] present in test data point [True]
84 Text feature [thus] present in test data point [True]
85 Text feature [due] present in test data point [True]
88 Text feature [approved] present in test data point [True]
89 Text feature [mutation] present in test data point [True]
90 Text feature [enhanced] present in test data point [True]
92 Text feature [three] present in test data point [True]
93 Text feature [two] present in test data point [True]
95 Text feature [proliferation] present in test data point [True]
96 Text feature [fig] present in test data point [True]
97 Text feature [revealed] present in test data point [True]
98 Text feature [figure] present in test data point [True]
99 Text feature [show] present in test data point [True]
Out of the top  100  features  65 are present in query point

## 4.2. K Nearest Neighbour Classification

### 4.2.1. Hyper parameter tuning

In [64]:

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

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


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


alpha = [5, 11, 15, 21, 31, 41, 51, 99]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = KNeighborsClassifier(n_neighbors=i)
    clf.fit(train_x_responseCoding, train_y)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_x_responseCoding, train_y)
    sig_clf_probs = sig_clf.predict_proba(cv_x_responseCoding)
    cv_log_error_array.append(log_loss(cv_y, sig_clf_probs, labels=clf.classes_, eps=1e-15))
    # to avoid rounding error while multiplying probabilites we use log-probability 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.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()


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

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

```
for alpha = 5
Log Loss : 1.1106865245096922
for alpha = 11
Log Loss : 1.0975875836482032
for alpha = 15
Log Loss : 1.0921265613501825
for alpha = 21
Log Loss : 1.0956683795240185
for alpha = 31
Log Loss : 1.0959451958194335
for alpha = 41
Log Loss : 1.1086096315570855
for alpha = 51
Log Loss : 1.119749289492384
for alpha = 99
Log Loss : 1.164288735084412
```



```
For values of best alpha =  15 The train log loss is: 0.6855555333780055
For values of best alpha =  15 The cross validation log loss is: 1.0921265613501825
For values of best alpha =  15 The test log loss is: 1.072279017870719
```

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

In [65]:

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

# methods of
```

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

Log loss : 1.0921265613501825
Number of mis-classified points : 0.37969924812030076
------------------- Confusion matrix -------------------



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



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

### 4.2.3.Sample Query point -1

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

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

```
Predicted Class : 2
Actual Class : 7
The  15  nearest neighbours of the test points belongs to classes [2 2 7 2 7 7 7 2 7 7 2 2 7 2 3]
Fequency of nearest points : Counter({2: 7, 7: 7, 3: 1})
```

### 4.2.4. Sample Query Point-2

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

test_point_index = 100

predicted_cls = sig_clf.predict(test_x_responseCoding[test_point_index].reshape(1,-1))
print("Predicted Class :", predicted_cls[0])
print("Actual Class :", test_y[test_point_index])
neighbors = clf.kneighbors(test_x_responseCoding[test_point_index].reshape(1, -1), alpha[best_alpha
])
print("the k value for knn is",alpha[best_alpha],"and the nearest neighbours of the test points be
longs to classes",train_y[neighbors[1][0]])
print("Fequency of nearest points :",Counter(train_y[neighbors[1][0]]))
```
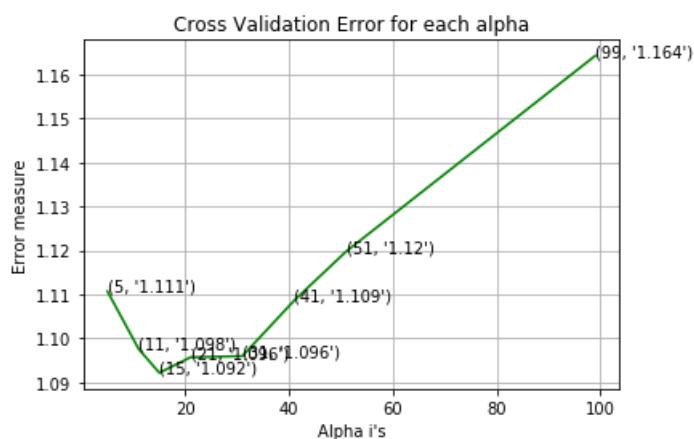
```
Predicted Class : 7
Actual Class : 7
the k value for knn is 15 and the nearest neighbours of the test points belongs to classes [7 2 7
7 7 7 5 2 7 2 6 6 6 5 6]
Fequency of nearest points : Counter({7: 6, 6: 4, 2: 3, 5: 2})
```

## 4.3. Logistic Regression

### 4.3.1. With Class balancing

#### 4.3.1.1. Hyper paramter tuning

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

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

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


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

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

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


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

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

```
for alpha = 1e-06
Log Loss : 1.376204958131797
for alpha = 1e-05
Log Loss : 1.3607096345898757
for alpha = 0.0001
Log Loss : 1.2996855092894215
for alpha = 0.001
Log Loss : 1.1364170518260186
for alpha = 0.01
Log Loss : 1.1920182966197623
for alpha = 0.1
Log Loss : 1.5286844610158679
for alpha = 1
Log Loss : 1.7302794707119695
for alpha = 10
Log Loss : 1.7514657596874037
for alpha = 100
Log Loss : 1.7536308618432714
```



```
For values of best alpha =  0.001 The train log loss is: 0.618166137635762
For values of best alpha =  0.001 The cross validation log loss is: 1.1364170518260186
For values of best alpha =  0.001 The test log loss is: 1.1100215541545624
```

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

In [69]:

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

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

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

```
Log loss : 1.1364170518260186
Number of mis-classified points : 0.3815789473684211
-------------------- Confusion matrix --------------------
```

| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 52.000 | 0.000 | 1.000 | 19.000 | 6.000 | 4.000 | 9.000 | 0.000 | 0.000 |
| 2 | 3.000 | 27.000 | 1.000 | 0.000 | 0.000 | 0.000 | 41.000 | 0.000 | 0.000 |
| 3 | 1.000 | 0.000 | 5.000 | 4.000 | 2.000 | 0.000 | 2.000 | 0.000 | 0.000 |
| 4 | 19.000 | 0.000 | 0.000 | 76.000 | 3.000 | 0.000 | 12.000 | 0.000 | 0.000 |
| 5 | 8.000 | 0.000 | 1.000 | 5.000 | 13.000 | 4.000 | 8.000 | 0.000 | 0.000 |
| 6 | 8.000 | 0.000 | 1.000 | 3.000 | 3.000 | 20.000 | 9.000 | 0.000 | 0.000 |
| 7 | 1.000 | 11.000 | 5.000 | 2.000 | 1.000 | 0.000 | 131.000 | 2.000 | 0.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 3.000 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.000 | 5.000 |

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

| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.565 | 0.000 | 0.071 | 0.174 | 0.214 | 0.143 | 0.042 | 0.000 | 0.000 |
| 2 | 0.033 | 0.711 | 0.071 | 0.000 | 0.000 | 0.000 | 0.190 | 0.000 | 0.000 |
| 3 | 0.011 | 0.000 | 0.357 | 0.037 | 0.071 | 0.000 | 0.009 | 0.000 | 0.000 |
| 4 | 0.207 | 0.000 | 0.000 | 0.697 | 0.107 | 0.000 | 0.056 | 0.000 | 0.000 |
| 5 | 0.087 | 0.000 | 0.071 | 0.046 | 0.464 | 0.143 | 0.037 | 0.000 | 0.000 |
| 6 | 0.087 | 0.000 | 0.071 | 0.028 | 0.107 | 0.714 | 0.042 | 0.000 | 0.000 |
| 7 | 0.011 | 0.289 | 0.357 | 0.018 | 0.036 | 0.000 | 0.606 | 1.000 | 0.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.014 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.005 | 0.000 | 1.000 |

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

| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.571 | 0.000 | 0.011 | 0.209 | 0.066 | 0.044 | 0.099 | 0.000 | 0.000 |
| 2 | 0.042 | 0.375 | 0.014 | 0.000 | 0.000 | 0.000 | 0.569 | 0.000 | 0.000 |
| 3 | 0.071 | 0.000 | 0.357 | 0.286 | 0.143 | 0.000 | 0.143 | 0.000 | 0.000 |
| 4 | 0.173 | 0.000 | 0.000 | 0.691 | 0.027 | 0.000 | 0.109 | 0.000 | 0.000 |
| 5 | 0.205 | 0.000 | 0.026 | 0.128 | 0.333 | 0.103 | 0.205 | 0.000 | 0.000 |
| 6 | 0.182 | 0.000 | 0.023 | 0.068 | 0.068 | 0.455 | 0.205 | 0.000 | 0.000 |
| 7 | 0.007 | 0.072 | 0.033 | 0.013 | 0.007 | 0.000 | 0.856 | 0.013 | 0.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.167 | 0.000 | 0.833 |

### 4.3.1.3. Feature Importance

In [70]:

```python
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 importent features are present in our query point")
    print("-"*50)
    print("The features that are most importent of the ",predicted_cls[0]," class:")
    print (tabulate(tabulte_list, headers=["Index",'Feature name', 'Present or Not']))
```

### 4.3.1.3.1. Correctly Classified point

In [71]:

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

```
Predicted Class : 2
Predicted Class Probabilities: [[0.0246 0.6476 0.01   0.0341 0.0329 0.017  0.2223 0.0055 0.0061]]
Actual Class : 7
--------------------------------------------------
202 Text feature [t315i] present in test data point [True]
227 Text feature [narrower] present in test data point [True]
Out of the top  500  features  2 are present in query point
```

### 4.3.1.3.2. Incorrectly Classified point

In [72]:

```
test_point_index = 100
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding[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: [[0.0923 0.1855 0.0134 0.0902 0.0339 0.0221 0.5488 0.0061 0.0077]]
Actual Class : 7
--------------------------------------------------
24 Text feature [constitutive] present in test data point [True]
44 Text feature [activated] present in test data point [True]
46 Text feature [constitutivelv] present in test data point [True]
```

```
73 Text feature [erk1] present in test data point [True]
200 Text feature [starved] present in test data point [True]
202 Text feature [phospho] present in test data point [True]
208 Text feature [ligand] present in test data point [True]
228 Text feature [downstream] present in test data point [True]
284 Text feature [extracellular] present in test data point [True]
377 Text feature [technology] present in test data point [True]
Out of the top  500  features  10 are present in query point
```

## 4.3.2. Without Class balancing

### 4.3.2.1. Hyper paramter tuning

In [73]:

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

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

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



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

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

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

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

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

```
for alpha = 1e-06
Log Loss : 1.3469646794324086
for alpha = 1e-05
Log Loss : 1.358646218121377
for alpha = 0.0001
Log Loss : 1.3291719274321534
for alpha = 0.001
Log Loss : 1.208675471831997
for alpha = 0.01
Log Loss : 1.2397204283254306
for alpha = 0.1
Log Loss : 1.3720210518253768
for alpha = 1
Log Loss : 1.620487349363535
```



For values of best alpha =  0.001 The train log loss is: 0.6086512192063078
For values of best alpha =  0.001 The cross validation log loss is: 1.208675471831997
For values of best alpha =  0.001 The test log loss is: 1.1258587901632542

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

In [74]:

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

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

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

```
#------------------------------

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

Log loss : 1.208675471831997
Number of mis-classified points : 0.37781954887218044
-------------------- Confusion matrix --------------------



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



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

Predicted Class

### 4.3.2.3. Feature Importance, Correctly Classified point

In [75]:

```python
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(train_x_onehotCoding,train_y)
test_point_index = 1
no_feature = 500
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:",
np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]),4))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.coef_)[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0],
test_df['TEXT'].iloc[test_point_index],test_df['Gene'].iloc[test_point_index],test_df['Variation']
.iloc[test_point_index], no_feature)
```
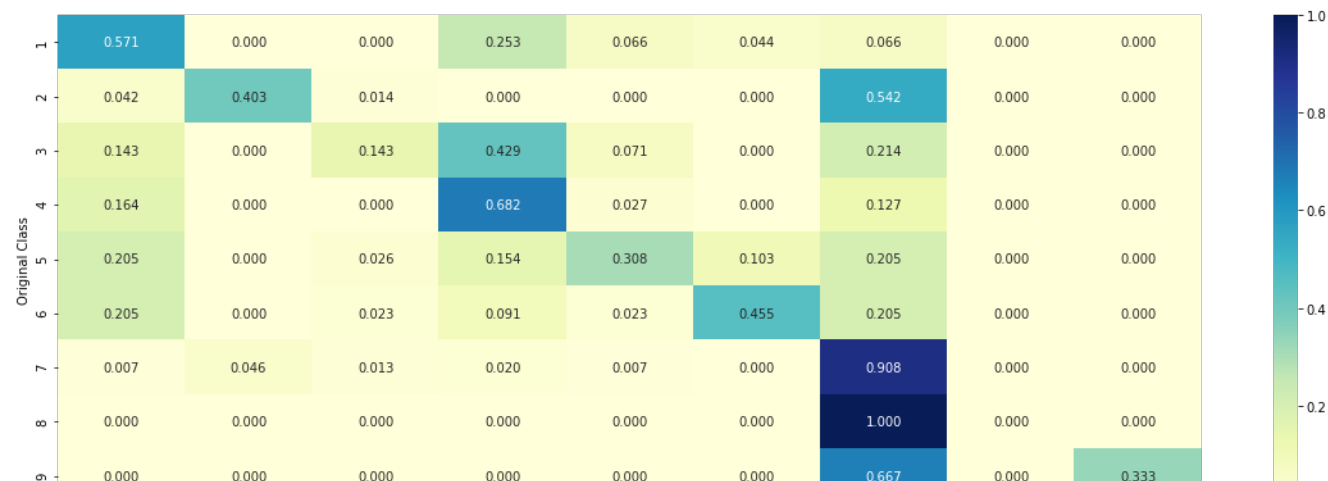
```
Predicted Class : 2
Predicted Class Probabilities: [[0.0304 0.656  0.0082 0.0432 0.0321 0.0173 0.2077 0.0039 0.0013]]
Actual Class : 7
--------------------------------------------------
224 Text feature [t315i] present in test data point [True]
253 Text feature [narrower] present in test data point [True]
Out of the top  500  features  2 are present in query point
```

### 4.3.2.4. Feature Importance, Inorrectly Classified point

In [76]:

```python
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: [[0.0941 0.1867 0.012  0.097  0.0347 0.0226 0.5463 0.0046 0.002 ]]
Actual Class : 7
--------------------------------------------------
77 Text feature [constitutive] present in test data point [True]
107 Text feature [constitutively] present in test data point [True]
127 Text feature [activated] present in test data point [True]
161 Text feature [erk1] present in test data point [True]
285 Text feature [phospho] present in test data point [True]
332 Text feature [extracellular] present in test data point [True]
350 Text feature [starved] present in test data point [True]
354 Text feature [technology] present in test data point [True]
394 Text feature [downstream] present in test data point [True]
Out of the top  500  features  9 are present in query point
```

## 4.4. Linear Support Vector Machines

### 4.4.1. Hyper paramter tuning

In [77]:

```python
# 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, t
ol=0.001,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', ra
ndom_state=None)

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



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

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

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


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

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

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

```
for C = 1e-05
Log Loss : 1.3651970530999085
for C = 0.0001
Log Loss : 1.374949223854886
for C = 0.001
Log Loss : 1.2674235291279028
for C = 0.01
Log Loss : 1.1806234085166019
for C = 0.1
Log Loss : 1.4287727750362658
for C = 1
Log Loss : 1.7517051728388966
for C = 10
Log Loss : 1.7539922678053463
for C = 100
Log Loss : 1.7540195871673772
```



```
For values of best alpha =  0.01 The train log loss is: 0.7465566172226077
For values of best alpha =  0.01 The cross validation log loss is: 1.1806234085166019
For values of best alpha =  0.01 The test log loss is: 1.1699143602769757
```

### 4.4.2. Testing model with best hyper parameters

In [78]:

```
# 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, t
ol=0.001,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', ra
ndom_state=None)

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


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

```
Log loss : 1.1806234085166019
Number of mis-classified points : 0.38345864661654133
-------------------- Confusion matrix --------------------
```

**Confusion Matrix (counts)** — Original Class (rows) vs Predicted Class (columns)

| Original \ Predicted | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 46.000 | 1.000 | 1.000 | 24.000 | 9.000 | 4.000 | 6.000 | 0.000 | 0.000 |
| 2 | 3.000 | 31.000 | 1.000 | 0.000 | 0.000 | 0.000 | 37.000 | 0.000 | 0.000 |
| 3 | 1.000 | 0.000 | 6.000 | 2.000 | 3.000 | 0.000 | 2.000 | 0.000 | 0.000 |
| 4 | 16.000 | 0.000 | 3.000 | 70.000 | 7.000 | 3.000 | 11.000 | 0.000 | 0.000 |
| 5 | 6.000 | 0.000 | 1.000 | 4.000 | 19.000 | 3.000 | 6.000 | 0.000 | 0.000 |
| 6 | 5.000 | 0.000 | 2.000 | 2.000 | 5.000 | 22.000 | 8.000 | 0.000 | 0.000 |
| 7 | 1.000 | 12.000 | 5.000 | 3.000 | 2.000 | 0.000 | 130.000 | 0.000 | 0.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 3.000 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 2.000 | 0.000 | 4.000 |

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

| Original \ Predicted | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.590 | 0.023 | 0.053 | 0.229 | 0.200 | 0.125 | 0.029 |  | 0.000 |
| 2 | 0.038 | 0.705 | 0.053 | 0.000 | 0.000 | 0.000 | 0.180 |  | 0.000 |
| 3 | 0.013 | 0.000 | 0.316 | 0.019 | 0.067 | 0.000 | 0.010 |  | 0.000 |
| 4 | 0.205 | 0.000 | 0.158 | 0.667 | 0.156 | 0.094 | 0.054 |  | 0.000 |
| 5 | 0.077 | 0.000 | 0.053 | 0.038 | 0.422 | 0.094 | 0.029 |  | 0.000 |
| 6 | 0.064 | 0.000 | 0.105 | 0.019 | 0.111 | 0.688 | 0.039 |  | 0.000 |
| 7 | 0.013 | 0.273 | 0.263 | 0.029 | 0.044 | 0.000 | 0.634 |  | 0.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.015 |  | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.010 |  | 1.000 |

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

| Original \ Predicted | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.505 | 0.011 | 0.011 | 0.264 | 0.099 | 0.044 | 0.066 | 0.000 | 0.000 |
| 2 | 0.042 | 0.431 | 0.014 | 0.000 | 0.000 | 0.000 | 0.514 | 0.000 | 0.000 |
| 3 | 0.071 | 0.000 | 0.429 | 0.143 | 0.214 | 0.000 | 0.143 | 0.000 | 0.000 |
| 4 | 0.145 | 0.000 | 0.027 | 0.636 | 0.064 | 0.027 | 0.100 | 0.000 | 0.000 |
| 5 | 0.154 | 0.000 | 0.026 | 0.103 | 0.487 | 0.077 | 0.154 | 0.000 | 0.000 |
| 6 | 0.114 | 0.000 | 0.045 | 0.045 | 0.114 | 0.500 | 0.182 | 0.000 | 0.000 |
| 7 | 0.007 | 0.078 | 0.033 | 0.020 | 0.013 | 0.000 | 0.850 | 0.000 | 0.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.333 | 0.000 | 0.667 |

### 4.3.3. Feature Importance

#### 4.3.3.1. For Correctly classified point

```
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 : 2
Predicted Class Probabilities: [[0.0596 0.5598 0.0122 0.0879 0.0427 0.0231 0.2057 0.0053 0.0036]]
Actual Class : 7
--------------------------------------------------
49 Text feature [t315i] present in test data point [True]
109 Text feature [narrower] present in test data point [True]
436 Text feature [cml] present in test data point [True]
Out of the top  500  features  3 are present in query point
```

### 4.3.3.2. For Incorrectly classified point

```
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: [[0.1304 0.1548 0.017  0.1183 0.0462 0.0326 0.4916 0.0052 0.0038]]
Actual Class : 7
--------------------------------------------------
34 Text feature [constitutive] present in test data point [True]
52 Text feature [constitutively] present in test data point [True]
70 Text feature [starved] present in test data point [True]
102 Text feature [activated] present in test data point [True]
142 Text feature [erk1] present in test data point [True]
232 Text feature [ligand] present in test data point [True]
233 Text feature [phospho] present in test data point [True]
254 Text feature [technology] present in test data point [True]
332 Text feature [extracellular] present in test data point [True]
355 Text feature [fgf1] present in test data point [True]
396 Text feature [serum] present in test data point [True]
398 Text feature [expressing] present in test data point [True]
418 Text feature [activation] present in test data point [True]
424 Text feature [remain] present in test data point [True]
452 Text feature [downstream] present in test data point [True]
459 Text feature [calvaria] present in test data point [True]
467 Text feature [ornitz] present in test data point [True]
472 Text feature [mitogen] present in test data point [True]
Out of the top  500  features  18 are present in query point
```

## 4.5 Random Forest Classifier

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

In [81]:

```python
# -------------------------------
# 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 : 1.254927457081247
for n_estimators = 100 and max depth =  10
Log Loss : 1.1920712824451987
for n_estimators = 200 and max depth =  5
Log Loss : 1.2344628349548514
for n_estimators = 200 and max depth =  10
Log Loss : 1.1772584109862803
for n_estimators = 500 and max depth =  5
Log Loss : 1.2260897606901613
for n_estimators = 500 and max depth =  10
Log Loss : 1.1678492035652386
for n_estimators = 1000 and max depth =  5
Log Loss : 1.2206060983283764
for n_estimators = 1000 and max depth =  10
Log Loss : 1.167424064224008
for n_estimators = 2000 and max depth =  5
Log Loss : 1.2164484837029024
for n_estimators = 2000 and max depth =  10
Log Loss : 1.1649179335443909
For values of best estimator =  2000 The train log loss is: 0.6985174287433517
For values of best estimator =  2000 The cross validation log loss is: 1.1649179335443909
For values of best estimator =  2000 The test log loss is: 1.1577784570208531
```

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

In [82]:

```
# --------------------------------
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_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/
# --------------------------------

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 : 1.1649179335443909
Number of mis-classified points : 0.38721804511278196
------------------- Confusion matrix -------------------
```

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



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



### 4.5.3. Feature Importance

#### 4.5.3.1. Correctly Classified point
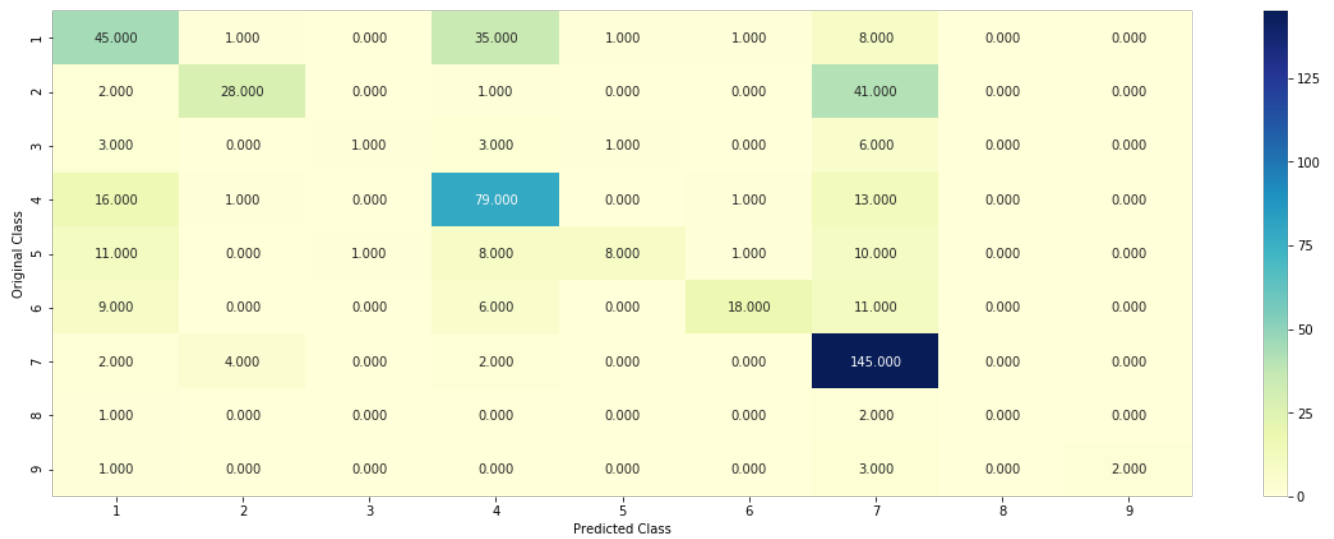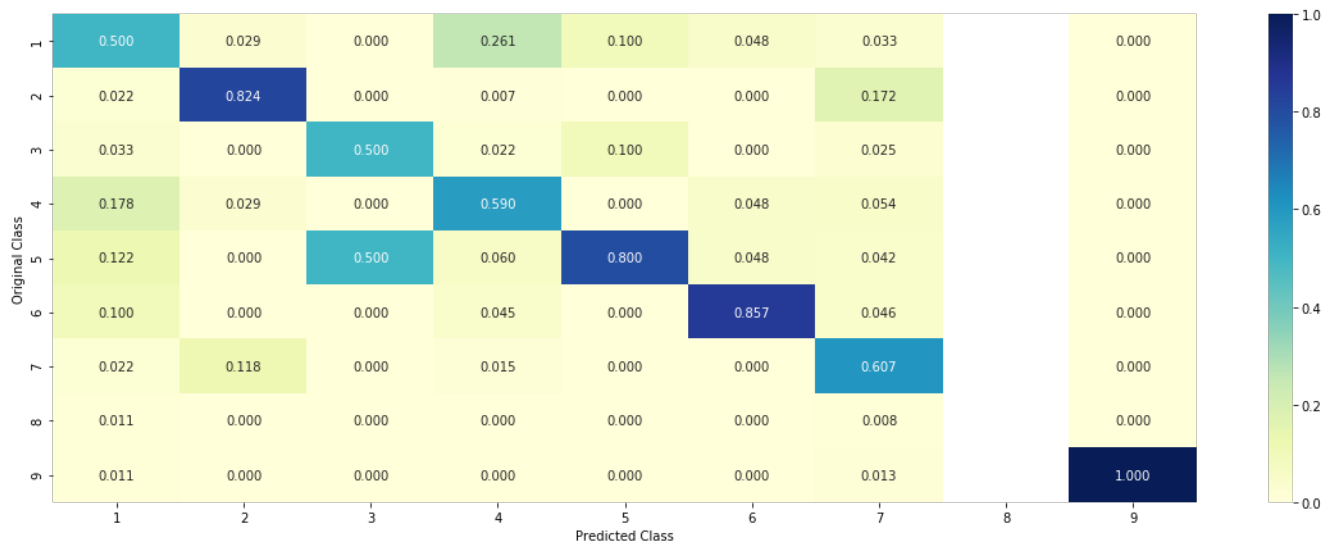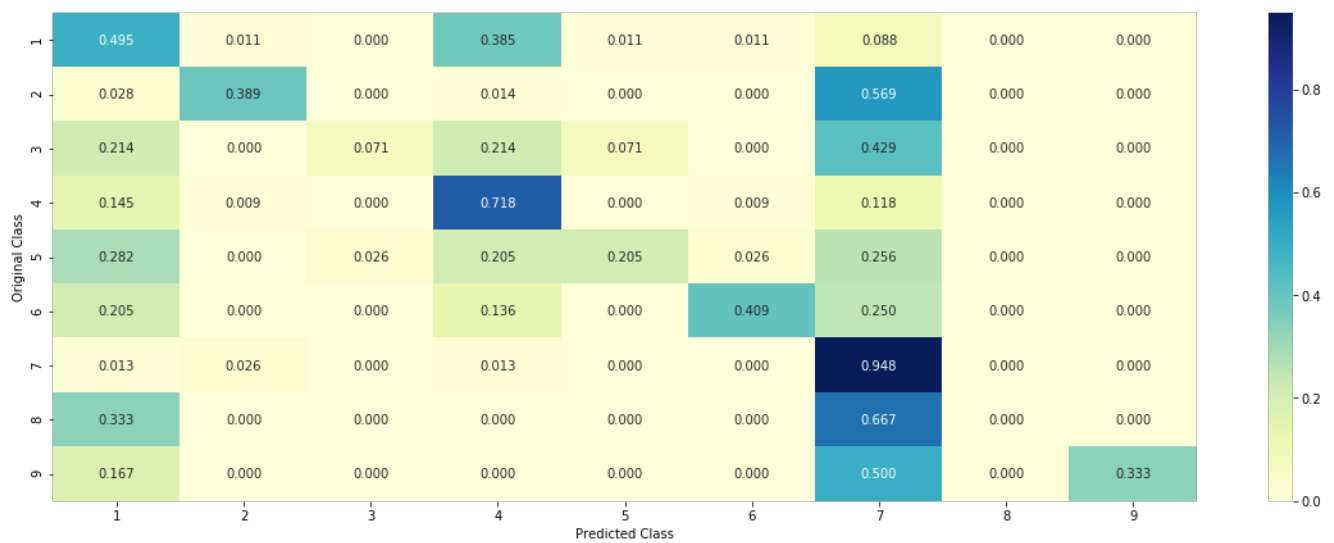
In [83]:

```python
# 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 : 2
Predicted Class Probabilities: [[0.0278 0.4485 0.0133 0.0285 0.0322 0.0239 0.4183 0.0039 0.0035]]
Actual Class : 7
--------------------------------------------------
0 Text feature [kinase] present in test data point [True]
1 Text feature [tyrosine] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [phosphorylation] present in test data point [True]
5 Text feature [constitutive] present in test data point [True]
7 Text feature [activation] present in test data point [True]
8 Text feature [inhibitors] present in test data point [True]
10 Text feature [inhibitor] present in test data point [True]
11 Text feature [treatment] present in test data point [True]
14 Text feature [function] present in test data point [True]
15 Text feature [signaling] present in test data point [True]
17 Text feature [erk] present in test data point [True]
18 Text feature [extracellular] present in test data point [True]
20 Text feature [therapy] present in test data point [True]
21 Text feature [oncogenic] present in test data point [True]
22 Text feature [kinases] present in test data point [True]
23 Text feature [growth] present in test data point [True]
31 Text feature [loss] present in test data point [True]
33 Text feature [protein] present in test data point [True]
35 Text feature [variants] present in test data point [True]
36 Text feature [drug] present in test data point [True]
37 Text feature [receptor] present in test data point [True]
40 Text feature [akt] present in test data point [True]
42 Text feature [months] present in test data point [True]
43 Text feature [efficacy] present in test data point [True]
46 Text feature [inhibition] present in test data point [True]
47 Text feature [patients] present in test data point [True]
48 Text feature [ic50] present in test data point [True]
49 Text feature [ba] present in test data point [True]
50 Text feature [phospho] present in test data point [True]
51 Text feature [cells] present in test data point [True]
52 Text feature [resistance] present in test data point [True]
53 Text feature [treated] present in test data point [True]
57 Text feature [expressing] present in test data point [True]
60 Text feature [inhibited] present in test data point [True]
62 Text feature [cell] present in test data point [True]
66 Text feature [clinical] present in test data point [True]
67 Text feature [f3] present in test data point [True]
71 Text feature [trial] present in test data point [True]
72 Text feature [tki] present in test data point [True]
76 Text feature [tkis] present in test data point [True]
86 Text feature [atp] present in test data point [True]
87 Text feature [type] present in test data point [True]
91 Text feature [imatinib] present in test data point [True]
93 Text feature [lines] present in test data point [True]
94 Text feature [factor] present in test data point [True]
97 Text feature [kit] present in test data point [True]
99 Text feature [assays] present in test data point [True]
Out of the top  100  features  48 are present in query point
```

### 4.5.3.2. Inorrectly Classified point

```
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("Actuall 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: [[0.0916 0.1288 0.0205 0.1183 0.0522 0.0346 0.5423 0.0059 0.0059]]
Actuall Class : 7
--------------------------------------------------
0 Text feature [kinase] present in test data point [True]
1 Text feature [tyrosine] present in test data point [True]
3 Text feature [phosphorylation] present in test data point [True]
5 Text feature [constitutive] present in test data point [True]
6 Text feature [missense] present in test data point [True]
7 Text feature [activation] present in test data point [True]
10 Text feature [inhibitor] present in test data point [True]
11 Text feature [treatment] present in test data point [True]
13 Text feature [activated] present in test data point [True]
14 Text feature [function] present in test data point [True]
15 Text feature [signaling] present in test data point [True]
17 Text feature [erk] present in test data point [True]
18 Text feature [extracellular] present in test data point [True]
22 Text feature [kinases] present in test data point [True]
23 Text feature [growth] present in test data point [True]
28 Text feature [downstream] present in test data point [True]
29 Text feature [functional] present in test data point [True]
33 Text feature [protein] present in test data point [True]
37 Text feature [receptor] present in test data point [True]
39 Text feature [constitutively] present in test data point [True]
44 Text feature [erk1] present in test data point [True]
47 Text feature [patients] present in test data point [True]
50 Text feature [phospho] present in test data point [True]
51 Text feature [cells] present in test data point [True]
53 Text feature [treated] present in test data point [True]
54 Text feature [proliferation] present in test data point [True]
57 Text feature [expressing] present in test data point [True]
62 Text feature [cell] present in test data point [True]
66 Text feature [clinical] present in test data point [True]
80 Text feature [dna] present in test data point [True]
83 Text feature [mitogen] present in test data point [True]
87 Text feature [type] present in test data point [True]
93 Text feature [lines] present in test data point [True]
94 Text feature [factor] present in test data point [True]
97 Text feature [kit] present in test data point [True]
98 Text feature [starved] 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)

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

```python
# 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 = [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 : 2.080563875466587
for n_estimators = 10 and max depth =  3
Log Loss : 1.7594690722298234
for n_estimators = 10 and max depth =  5
Log Loss : 1.4850610276142582
for n_estimators = 10 and max depth =  10
Log Loss : 2.029501709688584
for n_estimators = 50 and max depth =  2
Log Loss : 1.7304304466273224
for n_estimators = 50 and max depth =  3
Log Loss : 1.4877338318390345
for n_estimators = 50 and max depth =  5
Log Loss : 1.4385162944322603
for n_estimators = 50 and max depth =  10
Log Loss : 1.7266516162918693
for n_estimators = 100 and max depth =  2
Log Loss : 1.5811384686559244
for n_estimators = 100 and max depth =  3
Log Loss : 1.4716479163155611
for n_estimators = 100 and max depth =  5
Log Loss : 1.4487378644754785
for n_estimators = 100 and max depth =  10
Log Loss : 1.7651083172945004
for n_estimators = 200 and max depth =  2
Log Loss : 1.6170710892033717
for n_estimators = 200 and max depth =  3
Log Loss : 1.476796700629157
for n_estimators = 200 and max depth =  5
Log Loss : 1.4825577826825083
for n_estimators = 200 and max depth =  10
Log Loss : 1.7386487090688711
for n_estimators = 500 and max depth =  2
Log Loss : 1.695083360025382
for n_estimators = 500 and max depth =  3
Log Loss : 1.5537666460223531
for n_estimators = 500 and max depth =  5
Log Loss : 1.496842676964451
for n_estimators = 500 and max depth =  10
Log Loss : 1.7765181482008163
for n_estimators = 1000 and max depth =  2
Log Loss : 1.653865794432428
for n_estimators = 1000 and max depth =  3
Log Loss : 1.553591789233817
for n_estimators = 1000 and max depth =  5
Log Loss : 1.4827588794781739
for n_estimators = 1000 and max depth =  10
Log Loss : 1.8003937930287932
For values of best alpha =  50 The train log loss is: 0.0666884290578537
For values of best alpha =  50 The cross validation log loss is: 1.4385162944322603
For values of best alpha =  50 The test log loss is: 1.3467763044012273
```

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

In [86]:

```
# --------------------------------
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_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/
# --------------------------------

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 : 1.4385162944322603
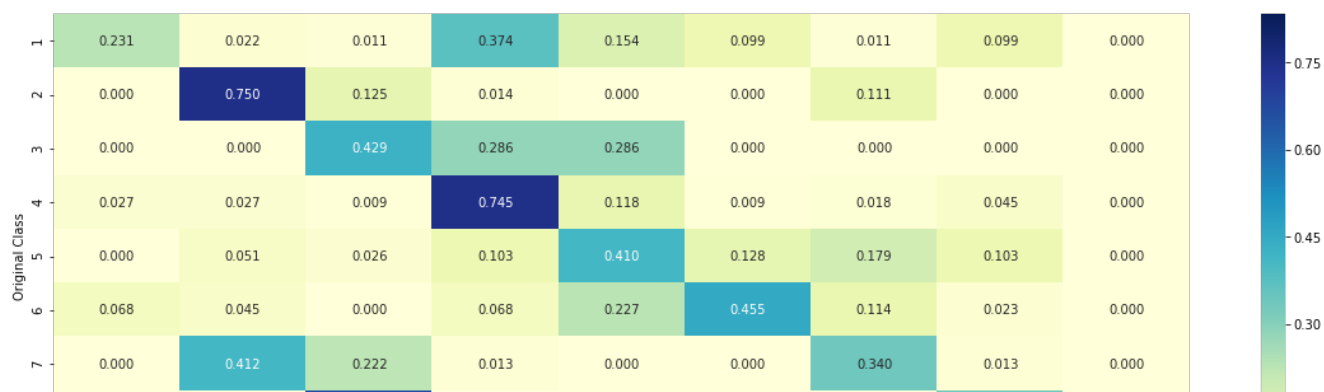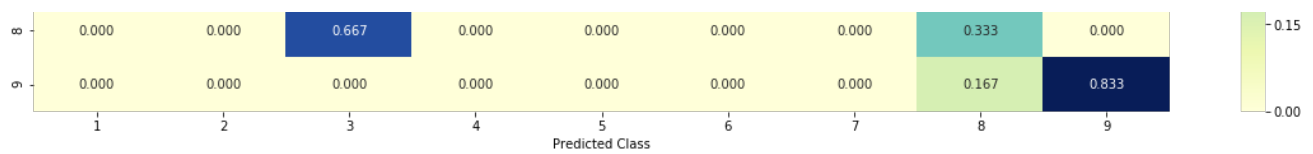Number of mis-classified points : 0.5169172932330827
-------------------- Confusion matrix --------------------



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



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

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
|---|---|---|---|---|---|---|---|---|---|---|
| 8 | 0.000 | 0.000 | 0.667 | 0.000 | 0.000 | 0.000 | 0.000 | 0.333 | 0.000 | 0.15 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.167 | 0.833 | 0.00 |

Predicted Class

### 4.5.5. Feature Importance

#### 4.5.5.1. Correctly Classified point

In [87]:

```python
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 : 2
Predicted Class Probabilities: [[0.0118 0.5681 0.1048 0.0161 0.0311 0.0204 0.1982 0.043  0.0066]]
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
Text is important feature
Text is important feature
Gene 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
Gene is important feature
Gene is important feature
Text is important feature
Variation is important feature
Text is important feature
Text is important feature
Gene is important feature
Gene is important feature
```

#### 4.5.5.2. Incorrectly Classified point

```python
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])
print (clf.feature_importances_)
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
for i in indices:
    if i<9:
        print("Gene is important feature")
    elif i<18:
        print("Variation is important feature")
    else:
        print("Text is important feature")
```

```
Predicted Class : 3
Predicted Class Probabilities: [[0.024  0.1942 0.2927 0.0404 0.0586 0.0341 0.2467 0.0945 0.0149]]
Actual Class : 7
[0.01556987 0.01135167 0.00191613 0.02320069 0.00548511 0.00991622
 0.07943479 0.00122621 0.00352945 0.14017882 0.09741066 0.0171485
 0.14357602 0.06066384 0.05721822 0.13663835 0.00277327 0.00556372
 0.02266877 0.02542624 0.00229206 0.05189376 0.01091213 0.02355446
 0.04527902 0.001995   0.00317702]
--------------------------------------------------
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
Text is important feature
Text is important feature
Gene 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
Gene is important feature
Gene is important feature
Text is important feature
Variation is important feature
Text is important feature
Text is important feature
Gene is important feature
Gene is important feature
```

## 4.7 Stack the models

### 4.7.1 testing with hyper parameter tuning

```python
# 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
```

```python
# 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 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, t
ol=0.001,
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', ra
ndom_state=None)

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


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

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

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

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


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

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


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

sig_clf1.fit(train_x_onehotCoding, train_y)
print("Logistic Regression :  Log Loss: %0.2f" % (log_loss(cv_y, sig_clf1.predict_proba(cv_x_onehot
Coding))))
sig_clf2.fit(train_x_onehotCoding, train_y)
print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y,
sig_clf2.predict_proba(cv_x_onehotCoding))))
sig_clf3.fit(train_x_onehotCoding, train_y)
print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predict_proba(cv_x_onehotCoding)))
)
```

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

```
Logistic Regression :  Log Loss: 1.12
Support vector machines : Log Loss: 1.75
Naive Bayes : Log Loss: 1.32
--------------------------------------------------
Stacking Classifer : for the value of alpha: 0.000100 Log Loss: 2.179
Stacking Classifer : for the value of alpha: 0.001000 Log Loss: 2.043
Stacking Classifer : for the value of alpha: 0.010000 Log Loss: 1.541
Stacking Classifer : for the value of alpha: 0.100000 Log Loss: 1.157
Stacking Classifer : for the value of alpha: 1.000000 Log Loss: 1.247
Stacking Classifer : for the value of alpha: 10.000000 Log Loss: 1.488
```

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

In [90]:

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

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

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

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

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

```
Log loss (train) on the stacking classifier : 0.6623333825200146
Log loss (CV) on the stacking classifier : 1.1571015584573332
Log loss (test) on the stacking classifier : 1.1752133744472877
Number of missclassified point : 0.38646616541353385
-------------------- Confusion matrix --------------------
```

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



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



### 4.7.3 Maximum Voting classifier

In [91]:

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

```
Log loss (train) on the VotingClassifier : 0.9068218744593947
Log loss (CV) on the VotingClassifier : 1.2300558139735123
Log loss (test) on the VotingClassifier : 1.235705422647174
Number of missclassified point : 0.37593984962406013
-------------------- Confusion matrix --------------------
```

First matrix (counts):

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 68.000 | 2.000 | 0.000 | 23.000 | 9.000 | 2.000 | 9.000 | 0.000 | 1.000 |
| 2 | 2.000 | 27.000 | 0.000 | 3.000 | 1.000 | 2.000 | 56.000 | 0.000 | 0.000 |
| 3 | 1.000 | 0.000 | 6.000 | 5.000 | 1.000 | 0.000 | 5.000 | 0.000 | 0.000 |
| 4 | 21.000 | 2.000 | 0.000 | 102.000 | 2.000 | 0.000 | 10.000 | 0.000 | 0.000 |
| 5 | 12.000 | 5.000 | 1.000 | 6.000 | 16.000 | 3.000 | 5.000 | 0.000 | 0.000 |
| 6 | 6.000 | 2.000 | 0.000 | 4.000 | 4.000 | 25.000 | 14.000 | 0.000 | 0.000 |
| 7 | 1.000 | 15.000 | 2.000 | 4.000 | 2.000 | 0.000 | 166.000 | 0.000 | 1.000 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 2.000 | 0.000 | 2.000 |
| 9 | 2.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 5.000 |

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

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.602 | 0.038 | 0.000 | 0.156 | 0.257 | 0.062 | 0.034 | | 0.111 |
| 2 | 0.018 | 0.509 | 0.000 | 0.020 | 0.029 | 0.062 | 0.210 | | 0.000 |
| 3 | 0.009 | 0.000 | 0.667 | 0.034 | 0.029 | 0.000 | 0.019 | | 0.000 |
| 4 | 0.186 | 0.038 | 0.000 | 0.694 | 0.057 | 0.000 | 0.037 | | 0.000 |
| 5 | 0.106 | 0.094 | 0.111 | 0.041 | 0.457 | 0.094 | 0.019 | | 0.000 |
| 6 | 0.053 | 0.038 | 0.000 | 0.027 | 0.114 | 0.781 | 0.052 | | 0.000 |
| 7 | 0.009 | 0.283 | 0.222 | 0.027 | 0.057 | 0.000 | 0.622 | | 0.111 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.007 | | 0.222 |
| 9 | 0.018 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | | 0.556 |

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

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.596 | 0.018 | 0.000 | 0.202 | 0.079 | 0.018 | 0.079 | 0.000 | 0.009 |
| 2 | 0.022 | 0.297 | 0.000 | 0.033 | 0.011 | 0.022 | 0.615 | 0.000 | 0.000 |
| 3 | 0.056 | 0.000 | 0.333 | 0.278 | 0.056 | 0.000 | 0.278 | 0.000 | 0.000 |
| 4 | 0.153 | 0.015 | 0.000 | 0.745 | 0.015 | 0.000 | 0.073 | 0.000 | 0.000 |
| 5 | 0.250 | 0.104 | 0.021 | 0.125 | 0.333 | 0.062 | 0.104 | 0.000 | 0.000 |
| 6 | 0.109 | 0.036 | 0.000 | 0.073 | 0.073 | 0.455 | 0.255 | 0.000 | 0.000 |
| 7 | 0.005 | 0.079 | 0.010 | 0.021 | 0.010 | 0.000 | 0.869 | 0.000 | 0.005 |
| 8 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.500 | 0.000 | 0.500 |
| 9 | 0.286 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.714 |