# Personalized cancer diagnosis

**Here I am using TfidfVectorizer and I am Taking top 1000 features over here...and i am applying Feature Engineering, So that our Log loss error decreases**

# 1. Business Problem

## 1.1. Description

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

Data: Memorial Sloan Kettering Cancer Center (MSKCC)

Download training_variants.zip and training_text.zip from Kaggle.

***Context:***

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

***Problem statement :***

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

## 1.2. Source/Useful Links

Some articles and reference blogs about the problem statement

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

## 1.3. Real-world/Business objectives and constraints.

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

# 2. Machine Learning Problem Formulation

## 2.1. Data

### 2.1.1. Data Overview

- Source: https://www.kaggle.com/c/msk-redefining-cancer-treatment/data
- We have two data files: one 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 [88]:  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
```

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

from mlxtend.classifier import StackingClassifier

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

## 3.1. Reading Data

### 3.1.1. Reading Gene and Variation Data

```python
In [89]: 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[89]:

|   | 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 [90]:
```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[90]:

|   | 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 [91]:
```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 t
he data
                if not word in stop_words:
                    string += word + " "
```

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

In [92]:
```
#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_tim
e, "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 : 185.98104735499282 seconds
```

In [93]:
```
#merging both gene_variations and text data based on ID
result = pd.merge(data, data_text,on='ID', how='left')
result.head()
```

Out[93]:

| | 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 [94]: result[result.isnull().any(axis=1)]
```

Out[94]:

|      | 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 [95]: result.loc[result['TEXT'].isnull(),'TEXT'] = result['Gene'] +' '+result
         ['Variation']
```

```
In [96]: result[result['ID']==1109]
```

Out[96]:

|      | 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 [97]: 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 o
         f output varaible 'y_true' [stratify=y_true]
         X_train, test_df, y_train, y_test = train_test_split(result, y_true, st
```

```
    ratify=y_true, test_size=0.2)
    # split the train data into train and cross validation by maintaining s
    ame distribution of output varaible 'y_train' [stratify=y_train]
    train_df, cv_df, y_train, y_cv = train_test_split(X_train, y_train, str
    atify=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 [98]:
```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 [99]:
```python
# it returns a dict, keys as class labels and values as the number of d
ata 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/num
py.argsort.html
# -(train_class_distribution.values): the minus sign will give us in de
```

```python
creasing 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_distri
bution.values[i], '(', np.round((train_class_distribution.values[i]/tra
in_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/num
py.argsort.html
# -(train_class_distribution.values): the minus sign will give us in de
creasing 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_distrib
ution.values[i], '(', np.round((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/num
py.argsort.html
# -(train_class_distribution.values): the minus sign will give us in de
```
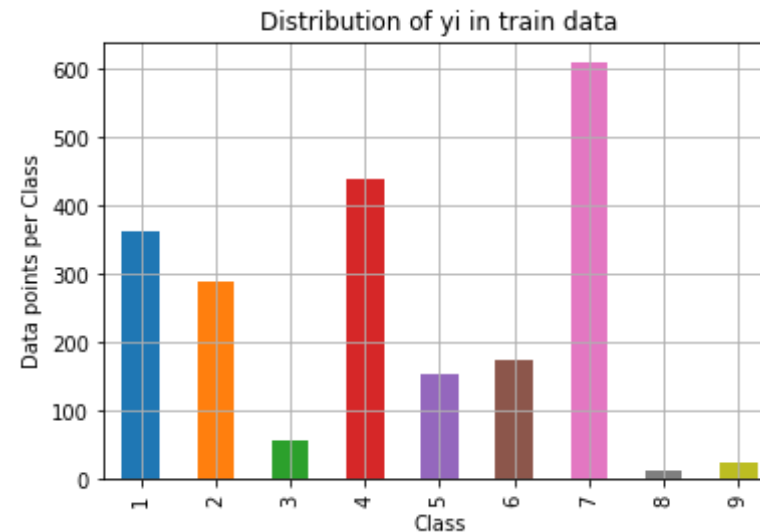
```
creasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',cv_class_distribut
ion.values[i], '(', np.round((cv_class_distribution.values[i]/cv_df.sha
pe[0]*100), 3), '%)')
```
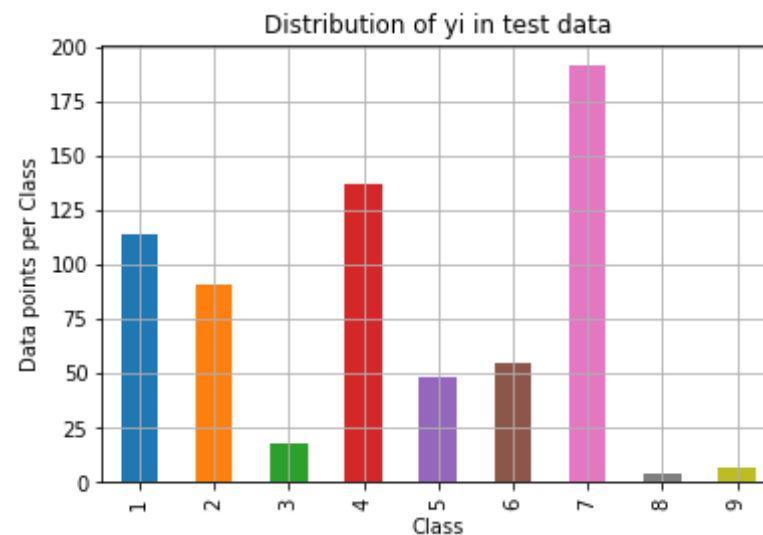
Distribution of yi in train data



```
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 %)
------------------------------------------------------------------------
---------
```
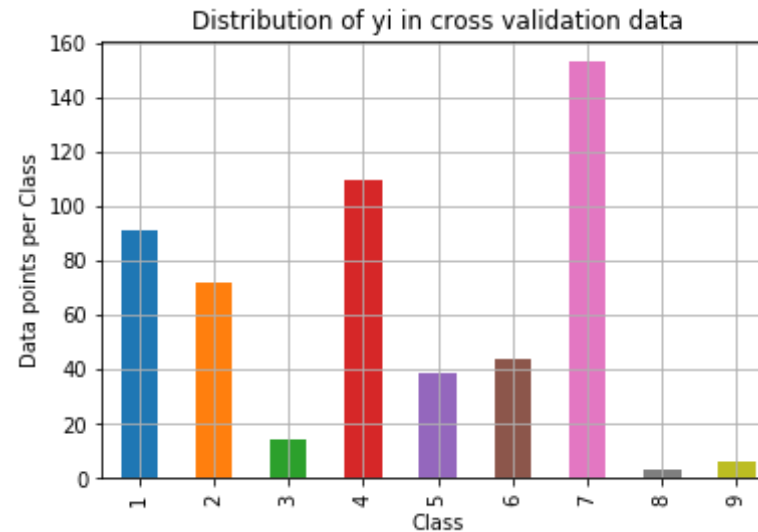
Distribution of yi in test data

```
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 [100]:  # 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)
```

```python
    # 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')
```

```python
    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=la
bels, 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=la
bels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()
```

In [101]:
```python
# 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_p
redicted_y, eps=1e-15))

predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)
```
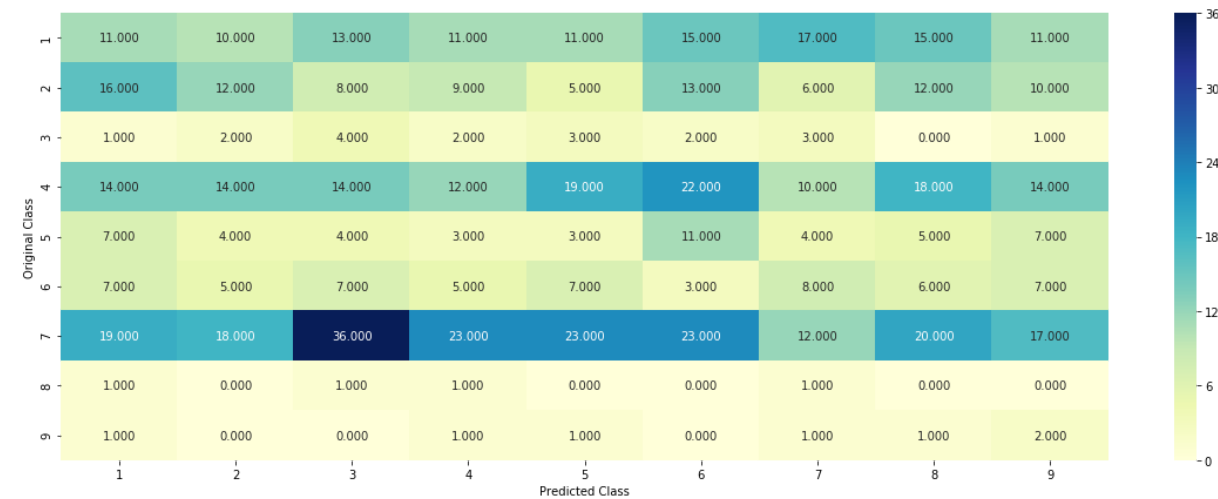
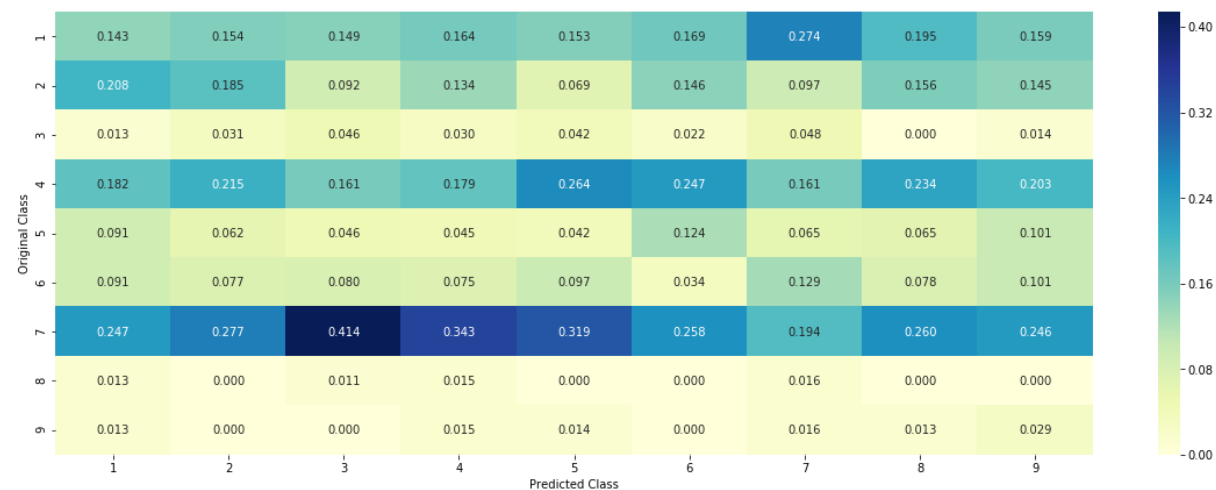Log loss on Cross Validation Data using Random Model 2.481523753995881
Log loss on Test Data using Random Model 2.5388681276861225
-------------------- Confusion matrix --------------------



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

| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.143 | 0.154 | 0.149 | 0.164 | 0.153 | 0.169 | 0.274 | 0.195 | 0.159 |
| 2 | 0.208 | 0.185 | 0.092 | 0.134 | 0.069 | 0.146 | 0.097 | 0.156 | 0.145 |
| 3 | 0.013 | 0.031 | 0.046 | 0.030 | 0.042 | 0.022 | 0.048 | 0.000 | 0.014 |
| 4 | 0.182 | 0.215 | 0.161 | 0.179 | 0.264 | 0.247 | 0.161 | 0.234 | 0.203 |
| 5 | 0.091 | 0.062 | 0.046 | 0.045 | 0.042 | 0.124 | 0.065 | 0.065 | 0.101 |
| 6 | 0.091 | 0.077 | 0.080 | 0.075 | 0.097 | 0.034 | 0.129 | 0.078 | 0.101 |
| 7 | 0.247 | 0.277 | 0.414 | 0.343 | 0.319 | 0.258 | 0.194 | 0.260 | 0.246 |
| 8 | 0.013 | 0.000 | 0.011 | 0.015 | 0.000 | 0.000 | 0.016 | 0.000 | 0.000 |
| 9 | 0.013 | 0.000 | 0.000 | 0.015 | 0.014 | 0.000 | 0.016 | 0.013 | 0.029 |

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



| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.096 | 0.088 | 0.114 | 0.096 | 0.096 | 0.132 | 0.149 | 0.132 | 0.096 |
| 2 | 0.176 | 0.132 | 0.088 | 0.099 | 0.055 | 0.143 | 0.066 | 0.132 | 0.110 |
| 3 | 0.056 | 0.111 | 0.222 | 0.111 | 0.167 | 0.111 | 0.167 | 0.000 | 0.056 |
| 4 | 0.102 | 0.102 | 0.102 | 0.088 | 0.139 | 0.161 | 0.073 | 0.131 | 0.102 |
| 5 | 0.146 | 0.083 | 0.083 | 0.062 | 0.062 | 0.229 | 0.083 | 0.104 | 0.146 |
| 6 | 0.127 | 0.091 | 0.127 | 0.091 | 0.127 | 0.055 | 0.145 | 0.109 | 0.127 |
| 7 | 0.099 | 0.094 | 0.188 | 0.120 | 0.120 | 0.120 | 0.063 | 0.105 | 0.089 |
| 8 | 0.250 | 0.000 | 0.250 | 0.250 | 0.000 | 0.000 | 0.250 | 0.000 | 0.000 |
| 9 | 0.143 | 0.000 | 0.000 | 0.143 | 0.143 | 0.000 | 0.143 | 0.143 | 0.286 |

## 3.3 Univariate Analysis

```
In [102]:  # code for response coding with Laplace smoothing.
           # alpha : used for laplace smoothing
           # feature: ['gene', 'variation']
```

```python
# df: ['train_df', 'test_df', 'cv_df']
# algorithm
# ----------
# Consider all unique values and the number of occurances of given feat
ure in train data dataframe
# build a vector (1*9) , the first element = (number of times it occure
d in class1 + 10*alpha / number of time it occurred in total data+90*al
pha)
# gv_dict is like a look up table, for every gene it store a (1*9) repr
esentation 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_f
ea'
# 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:
    #         {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
```

```python
    # 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 arr
ay for each gene/variation
    gv_dict = dict()

    # denominator will contain the number of time that particular featu
re occured in whole data
    for i, denominator in value_count.items():
        # vec will contain (p(yi==1/Gi) probability of gene/variation b
elongs 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['Ge
ne']=='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[f
eature]==i)]

            # cls_cnt.shape[0](numerator) will contain the number of ti
```

```python
me that particular feature occured 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 a
s 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, 0.068181
818181818177, 0.13636363636363635, 0.25, 0.19318181818181818, 0.0378787
8787878788, 0.03787878787878788, 0.03787878787878788],
    #     'TP53': [0.32142857142857145, 0.061224489795918366, 0.061224
489795918366, 0.27040816326530615, 0.061224489795918366, 0.066326530612
244902, 0.051020408163265307, 0.051020408163265307, 0.05612244897959183
7],
    #     'EGFR': [0.056818181818181816, 0.21590909090909091, 0.0625,
 0.068181818181818177, 0.068181818181818177, 0.0625, 0.3465909090909091
2, 0.0625, 0.056818181818181816],
    #     'BRCA2': [0.13333333333333333, 0.060606060606060608, 0.06060
6060606060608, 0.078787878787878782, 0.1393939393939394, 0.345454545454
54546, 0.060606060606060608, 0.060606060606060608, 0.06060606060606060
8],
    #     'PTEN': [0.069182389937106917, 0.062893081761006289, 0.06918
2389937106917, 0.46540880503144655, 0.075471698113207544, 0.06289308176
1006289, 0.069182389937106917, 0.062893081761006289, 0.0628930817610062
89],
    #     'KIT': [0.066225165562913912, 0.25165562913907286, 0.0728476
82119205295, 0.072847682119205295, 0.066225165562913912, 0.066225165562
913912, 0.2715231788079 4702, 0.066225165562913912, 0.06622516556291391
2],
    #     'BRAF': [0.066666666666666666, 0.17999999999999999, 0.073333
333333333334, 0.073333333333333334, 0.093333333333333338, 0.08000000000
0000002, 0.29999999999999999, 0.066666666666666666, 0.06666666666666666
6],
    #     ...
```

```
#       }
    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 e
ach feature value in the data
    gv_fea = []
    # for every feature values in the given data frame we will check if
 it is there in the train data then we will add the feature to gv_fea
    # if not we will add [1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9,1/9] to gv_fe
a
    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 [103]: 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 : 225
BRCA1     183
EGFR       99
TP53       96
PTEN       87
BRCA2      71
KIT        59
BRAF       57
ALK        48
ERBB2      46
FGFR2      37
Name: Gene, dtype: int64
```
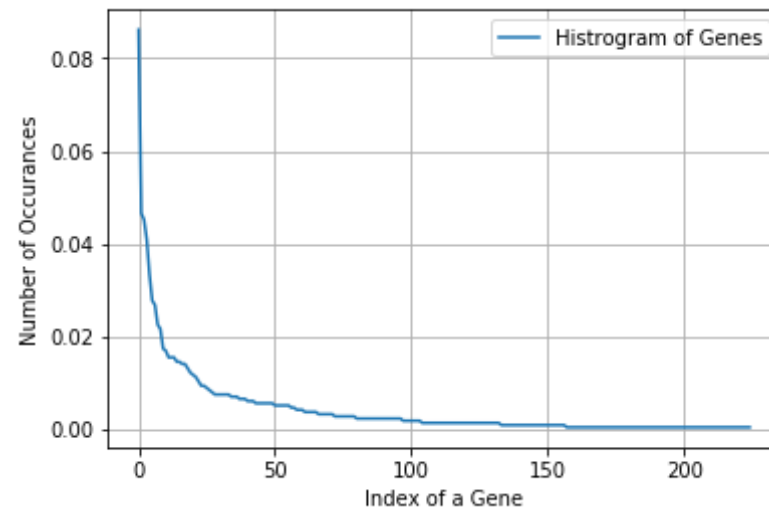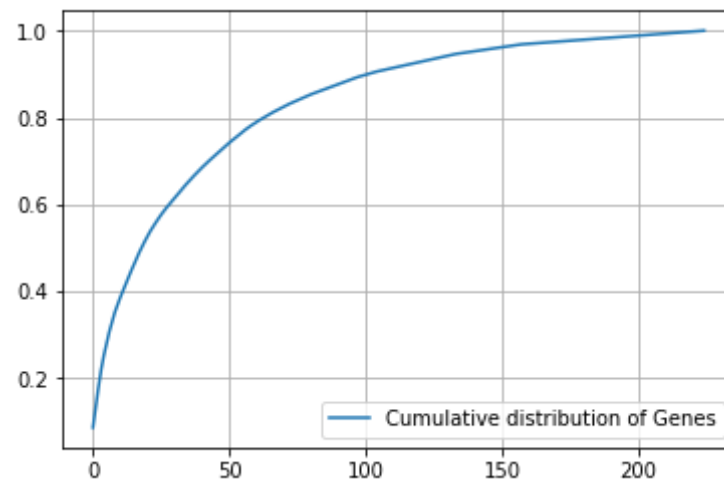
In [104]:
```
print("Ans: There are", unique_genes.shape[0] ,"different categories of
    genes in the train data, and they are distibuted as follows",)
```

```
Ans: There are 225 different categories of genes in the train data, and
they are distibuted as follows
```

In [105]:
```
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 [106]:
```python
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 [107]:
```
#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 [108]:
```
print("train_gene_feature_responseCoding is converted feature using respone coding method. The shape of gene feature:", train_gene_feature_responseCoding.shape)
```

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

In [109]:
```
# one-hot encoding of Gene feature.
gene_vectorizer = TfidfVectorizer()
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])
```

```
test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gen
e'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

In [110]: `train_df['Gene'].head()`

Out[110]: 
```
537      SMAD2
2808     BRCA2
3286       RET
2860     BRCA2
1692      PMS2
Name: Gene, dtype: object
```

In [111]: 
```python
print("train_gene_feature_onehotCoding is converted feature using one-h
ot encoding method. The shape of gene feature:", train_gene_feature_one
hotCoding.shape)
```

```
train_gene_feature_onehotCoding is converted feature using one-hot enco
ding method. The shape of gene feature: (2124, 224)
```

**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 [112]: 
```python
alpha = [10 ** x for x in range(-5, 1)] # hyperparam for SGD classifie
r.
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.clas
```

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

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_arra
y[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 vali
dation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps
=1e-15))
predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log l
oss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```
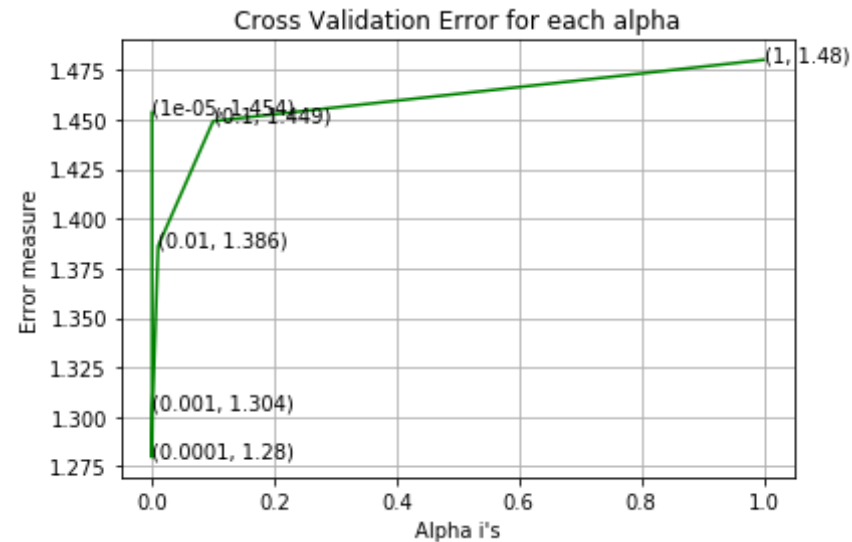
```
For values of alpha =  1e-05 The log loss is: 1.4535720391701341
For values of alpha =  0.0001 The log loss is: 1.2796360647519256
For values of alpha =  0.001 The log loss is: 1.3040318646899807
For values of alpha =  0.01 The log loss is: 1.386143387364092
```

```
For values of alpha =   0.1 The log loss is: 1.4491930791170367
For values of alpha =   1 The log loss is: 1.4800699070552643
```

Cross Validation Error for each alpha



```
For values of best alpha =   0.0001 The train log loss is: 1.01434767022
80119
For values of best alpha =   0.0001 The cross validation log loss is: 1.
2796360647519256
For values of best alpha =   0.0001 The test log loss is: 1.183349939962
4332
```

**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 [113]:  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'])))].shap
e[0]

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

```
Q6. How many data points in Test and CV datasets are covered by the  22
5  genes in train dataset?
Ans
1. In test data 640 out of 665 : 96.2406015037594
2. In cross validation data 507 out of  532 : 95.30075187969925
```

### 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 [114]:
```python
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 : 1940
Truncating_Mutations    57
Amplification           46
Deletion                42
Fusions                 16
Overexpression           4
Q61H                     3
T58I                     3
G12V                     3
T167A                    2
```
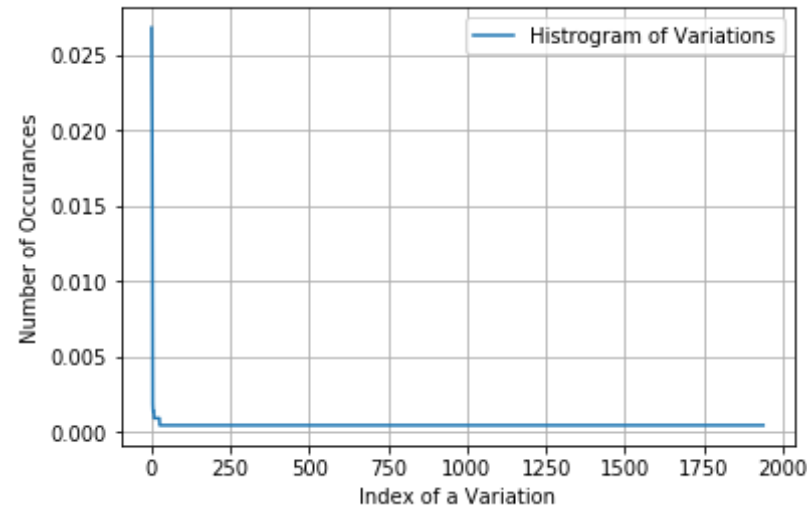
```
G12A                     2
Name: Variation, dtype: int64
```

In [115]: 
```python
print("Ans: There are", unique_variations.shape[0] ,"different categori
es of variations in the train data, and they are distibuted as follows"
,)
```

Ans: There are 1940 different categories of variations in the train dat
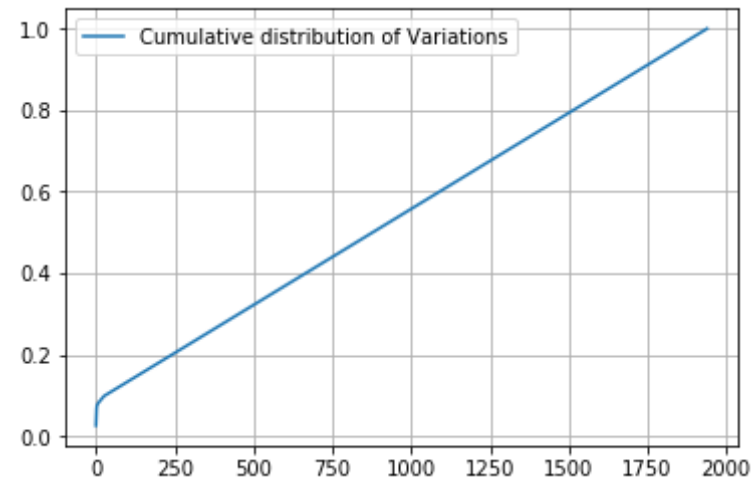a, and they are distibuted as follows

In [116]: 
```python
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.legend()
plt.grid()
plt.show()
```



In [117]: 
```python
c = np.cumsum(h)
print(c)
plt.plot(c,label='Cumulative distribution of Variations')
```

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

```
[0.02683616 0.04849341 0.06826742 ... 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 [118]:  # 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, "V
ariation", cv_df))
```

In [119]:
```
print("train_variation_feature_responseCoding is a converted feature us
ing the response coding method. The shape of Variation feature:", train
_variation_feature_responseCoding.shape)
```

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

In [120]:
```
# one-hot encoding of variation feature.
variation_vectorizer = TfidfVectorizer()
train_variation_feature_onehotCoding = variation_vectorizer.fit_transfo
rm(train_df['Variation'])
test_variation_feature_onehotCoding = variation_vectorizer.transform(te
st_df['Variation'])
cv_variation_feature_onehotCoding = variation_vectorizer.transform(cv_d
f['Variation'])
```

In [121]:
```
print("train_variation_feature_onehotEncoded is converted feature using
 the onne-hot encoding method. The shape of Variation feature:", train_
variation_feature_onehotCoding.shape)
```

train_variation_feature_onehotEncoded is converted feature using the on
ne-hot encoding method. The shape of Variation feature: (2124, 1968)

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

Let's build a model just like the earlier!

In [122]:
```
alpha = [10 ** x for x in range(-5, 1)]

cv_log_error_array=[]
```

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

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_arra
y[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 vali
dation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps
=1e-15))
predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log l
oss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
For values of alpha =  1e-05 The log loss is: 1.7028898672825914
For values of alpha =  0.0001 The log loss is: 1.6919908762776148
For values of alpha =  0.001 The log loss is: 1.69567590160075134
For values of alpha =  0.01 The log loss is: 1.6989173758756633
For values of alpha =  0.1 The log loss is: 1.7106118310238287
For values of alpha =  1 The log loss is: 1.7125669272790929
```



Cross Validation Error for each alpha

```
For values of best alpha =  0.0001 The train log loss is: 0.73009697523
11777
For values of best alpha =  0.0001 The cross validation log loss is: 1.
6919908762776148
For values of best alpha =  0.0001 The test log loss is: 1.682996162344
8697
```

**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 [123]:
```python
print("Q12. How many data points are covered by total ", unique_variati
ons.shape[0], " genes in test and cross validation data sets?")
test_coverage=test_df[test_df['Variation'].isin(list(set(train_df['Vari
ation'])))].shape[0]
cv_coverage=cv_df[cv_df['Variation'].isin(list(set(train_df['Variation'
])))].shape[0]
print('Ans\n1. In test data',test_coverage, 'out of',test_df.shape[0],
":",(test_coverage/test_df.shape[0])*100)
print('2. In cross validation data',cv_coverage, 'out of ',cv_df.shape[
0],":" ,(cv_coverage/cv_df.shape[0])*100)
```

```
Q12. How many data points are covered by total  1940  genes in test and
cross validation data sets?
Ans
1. In test data 79 out of 665 : 11.879699248120302
2. In cross validation data 54 out of  532 : 10.150375939849624
```

### 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 [124]:
```python
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 [125]: 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 )/(t
          otal_dict.get(word,0)+90)))
                      text_feature_responseCoding[row_index][i] = math.exp(sum_pr
          ob/len(row['TEXT'].split()))
                      row_index += 1
              return text_feature_responseCoding
```

```
In [126]: # building a CountVectorizer with all the words that occured minimum 3
           times in train data
          text_vectorizer = TfidfVectorizer(min_df=3,max_features=1000)
          train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_d
          f['TEXT'])
          # getting all the feature names (words)
          train_text_features= text_vectorizer.get_feature_names()

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

          # zip(list(text_features),text_fea_counts) will zip a word with its num
          ber of times it occured
          text_fea_dict = dict(zip(list(train_text_features),train_text_fea_count
          s))


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

          Total number of unique words in train data : 1000
```

```
In [127]: 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)


          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 [128]: #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 [129]: # 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_t
          ext_feature_responseCoding.sum(axis=1)).T

In [130]: # don't forget to normalize every feature
          train_text_feature_onehotCoding = normalize(train_text_feature_onehotCo
```

```
ding, axis=0)

# we use the same vectorizer that was trained on train data
test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEX
T'])
# don't forget to normalize every feature
test_text_feature_onehotCoding = normalize(test_text_feature_onehotCodi
ng, 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 [131]:
```
#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 [132]:
```
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.clas
ses_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_cv
, predict_y, labels=clf.classes_, eps=1e-15))

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_arra
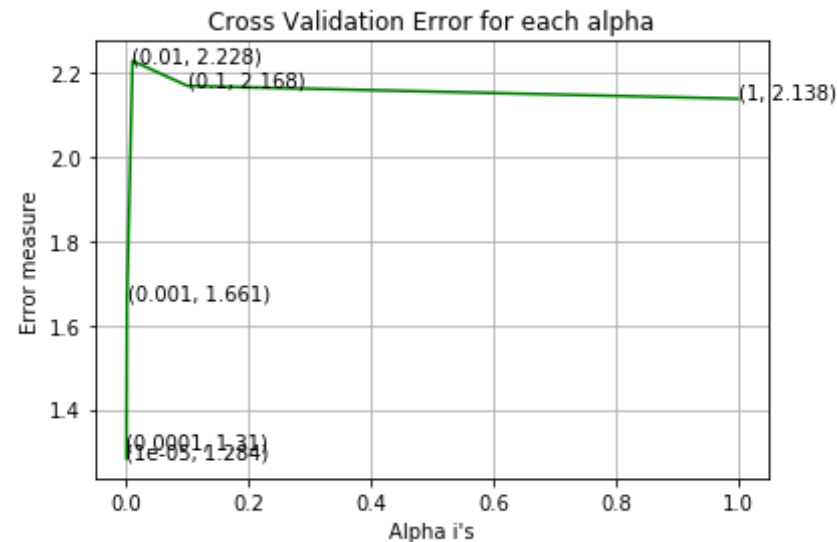```

```python
            y[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 vali
dation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps
=1e-15))
predict_y = sig_clf.predict_proba(test_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log l
oss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
For values of alpha =  1e-05 The log loss is: 1.2840223712982748
For values of alpha =  0.0001 The log loss is: 1.3100752177964912
For values of alpha =  0.001 The log loss is: 1.6610719668238374
For values of alpha =  0.01 The log loss is: 2.227883349080134
For values of alpha =  0.1 The log loss is: 2.1683742775323456
For values of alpha =  1 The log loss is: 2.137918196713479
```

Cross Validation Error for each alpha

For values of best alpha =  1e-05 The train log loss is: 0.7520981542039824
For values of best alpha =  1e-05 The cross validation log loss is: 1.2840223712982748
For values of best alpha =  1e-05 The test log loss is: 1.1223565604180978

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

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

```
In [136]:  def get_intersec_text(df):
               df_text_vec = TfidfVectorizer(min_df=3,max_features=1000)
               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_coun
           ts))
               len1 = len(set(df_text_features))
```

```
        len2 = len(set(train_text_features) & set(df_text_features))
        return len1,len2
```

In [137]:
```
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 appe
ared in train data")
```

```
94.2 % of word of test data appeared in train data
93.7 % of word of Cross Validation appeared in train data
```

## 4. Machine Learning Models

In [138]:
```
#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 probabilit
ies 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 [139]:
```
def report_log_loss(train_x, train_y, test_x, test_y,  clf):
    clf.fit(train_x, train_y)
```

```python
        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 [140]:
```python
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = TfidfVectorizer()
    var_count_vec = TfidfVectorizer()
    text_count_vec = TfidfVectorizer(min_df=3,max_features=1000)

    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 p
oint [{}]".format(word,yes_no))
        else:
            word = text_vec.get_feature_names()[v-(fea1_len+fea2_len)]
            yes_no = True if word in text.split() else False
            if yes_no:
                word_present += 1
                print(i, "Text feature [{}] present in test data point
```

```
   [{}]".format(word,yes_no))

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

# Feature Engineering

Here We are combining gene and variance and we are fitting TF-IDF Vectorizer to it... and then
we are going to transform train[text] , Test[text] and cv[text] using tfidf vectorizer of gene and
variance combined....

Bythis we can ensure that we will get some more information...

In [141]:
```
gene_variation_data=[]
for i in data["Gene"].values:
    gene_variation_data.append(i)
for i in data["Variation"].values:
    gene_variation_data.append(i)
cnt_vec=TfidfVectorizer()
Fit_vect=cnt_vec.fit(gene_variation_data)
train_text_vec=cnt_vec.transform(train_df["TEXT"])
cv_text_vec=cnt_vec.transform(cv_df["TEXT"])
test_text_vec=cnt_vec.transform(test_df["TEXT"])
train_text_vec=normalize(train_text_vec,axis=0)
cv_text_vec=normalize(cv_text_vec,axis=0)
test_text_vec=normalize(test_text_vec,axis=0)
```

## Stacking the three types of features

In [143]:
```
train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding,t
rain_variation_feature_onehotCoding))
test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding,tes
t_variation_feature_onehotCoding))
```

```python
cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding,cv_vari
ation_feature_onehotCoding))

train_gene_var_updt=hstack((train_gene_var_onehotCoding,train_text_vec
))
train_x_onehotCoding = hstack((train_gene_var_updt, train_text_feature_
onehotCoding)).tocsr()
train_y = np.array(list(train_df['Class']))


test_gene_var_updt=hstack((test_gene_var_onehotCoding,test_text_vec))
test_x_onehotCoding = hstack((test_gene_var_updt, test_text_feature_one
hotCoding)).tocsr()
test_y = np.array(list(test_df['Class']))


cv_gene_var_updt=hstack((cv_gene_var_onehotCoding,cv_text_vec))
cv_x_onehotCoding = hstack((cv_gene_var_updt, cv_text_feature_onehotCod
ing)).tocsr()
cv_y = np.array(list(cv_df['Class']))


train_gene_var_responseCoding = np.hstack((train_gene_feature_responseC
oding,train_variation_feature_responseCoding))
test_gene_var_responseCoding = np.hstack((test_gene_feature_responseCod
ing,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, trai
n_text_feature_responseCoding))
test_x_responseCoding = np.hstack((test_gene_var_responseCoding, test_t
ext_feature_responseCoding))
cv_x_responseCoding = np.hstack((cv_gene_var_responseCoding, cv_text_fe
ature_responseCoding))
```

```python
In [144]: 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 = ", t
est_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, 64
30)
(number of data points * number of features) in test data =  (665, 643
0)
(number of data points * number of features) in cross validation data =
(532, 6430)
```

In [145]:
```
print(" Response encoding features :")
print("(number of data points * number of features) in train data = ",
train_x_responseCoding.shape)
print("(number of data points * number of features) in test data = ", t
est_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, 2
7)
(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 [146]:
```
alpha = [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100,1000]
```

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

```
for alpha = 1e-05
Log Loss :  1.2696926448330919
for alpha = 0.0001
Log Loss :  1.2675173088326164
for alpha = 0.001
Log Loss :  1.2677552880301473
for alpha = 0.1
Log Loss :  1.2649840739954383
for alpha = 1
Log Loss :  1.3020587742514516
for alpha = 10
Log Loss :  1.5124417059265438
for alpha = 100
Log Loss :  1.4899333824751542
for alpha = 1000
Log Loss :  1.4680355636504934
```

```
For values of best alpha =  0.1 The train log loss is: 0.73083779120114
83
For values of best alpha =  0.1 The cross validation log loss is: 1.264
9840739954383
For values of best alpha =  0.1 The test log loss is: 1.194488541321399
3
```

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

In [147]:
```python
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-pro
bability estimates
print("Log Loss :",log_loss(cv_y, sig_clf_probs))
print("Number of missclassified point :", np.count_nonzero((sig_clf.pre
dict(cv_x_onehotCoding)- cv_y))/cv_y.shape[0])
plot_confusion_matrix(cv_y, sig_clf.predict(cv_x_onehotCoding.toarray
()))
```

```
Log Loss : 1.2649840739954383
Number of missclassified point : 0.38345864661654133
-------------------- Confusion matrix --------------------
```

| 70.000 | 1.000 | 0.000 | 8.000 | 7.000 | 2.000 | 3.000 | 0.000 | 0.000 |

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 2 | 8.000 | 32.000 | 0.000 | 4.000 | 0.000 | 1.000 | 27.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 4.000 | 4.000 | 2.000 | 0.000 | 4.000 | 0.000 | 0.000 |
| 4 | 39.000 | 1.000 | 1.000 | 58.000 | 3.000 | 1.000 | 6.000 | 0.000 | 1.000 |
| 5 | 11.000 | 1.000 | 0.000 | 4.000 | 13.000 | 3.000 | 7.000 | 0.000 | 0.000 |
| 6 | 7.000 | 4.000 | 1.000 | 3.000 | 4.000 | 23.000 | 2.000 | 0.000 | 0.000 |
| 7 | 3.000 | 23.000 | 2.000 | 0.000 | 0.000 | 0.000 | 124.000 | 0.000 | 1.000 |
| 8 | 2.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 |
| 9 | 1.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 4.000 |

Predicted Class

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



| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.496 | 0.016 | 0.000 | 0.098 | 0.241 | 0.067 | 0.017 | | 0.000 |
| 2 | 0.057 | 0.516 | 0.000 | 0.049 | 0.000 | 0.033 | 0.155 | | 0.000 |
| 3 | 0.000 | 0.000 | 0.500 | 0.049 | 0.069 | 0.000 | 0.023 | | 0.000 |
| 4 | 0.277 | 0.016 | 0.125 | 0.707 | 0.103 | 0.033 | 0.034 | | 0.167 |
| 5 | 0.078 | 0.016 | 0.000 | 0.049 | 0.448 | 0.100 | 0.040 | | 0.000 |
| 6 | 0.050 | 0.065 | 0.125 | 0.037 | 0.138 | 0.767 | 0.011 | | 0.000 |
| 7 | 0.021 | 0.371 | 0.250 | 0.000 | 0.000 | 0.000 | 0.713 | | 0.167 |
| 8 | 0.014 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.006 | | 0.000 |
| 9 | 0.007 | 0.000 | 0.000 | 0.012 | 0.000 | 0.000 | 0.000 | | 0.667 |

Predicted Class

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



| Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.769 | 0.011 | 0.000 | 0.088 | 0.077 | 0.022 | 0.033 | 0.000 | 0.000 |
| 2 | 0.111 | 0.444 | 0.000 | 0.056 | 0.000 | 0.014 | 0.375 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 0.286 | 0.286 | 0.143 | 0.000 | 0.286 | 0.000 | 0.000 |
| 4 | 0.355 | 0.009 | 0.009 | 0.527 | 0.027 | 0.009 | 0.055 | 0.000 | 0.009 |

### 4.1.1.3. Feature Importance, Correctly classified point

```
In [148]: 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)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[0.049  0.05   0.0152 0.0784 0.0411 0.0
369 0.7194 0.0055 0.0044]]
Actual Class : 7
--------------------------------------------------
```

### 4.1.1.4. Feature Importance, Incorrectly classified point

```
In [149]: 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)
```

```
Predicted Class : 1
Predicted Class Probabilities: [[0.6399 0.0545 0.0167 0.1034 0.0447 0.0
411 0.0887 0.0061 0.0049]]
Actual Class : 1
--------------------------------------------------
```

## 4.2. K Nearest Neighbour Classification

### 4.2.1. Hyper parameter tuning

In [150]:
```python
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.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
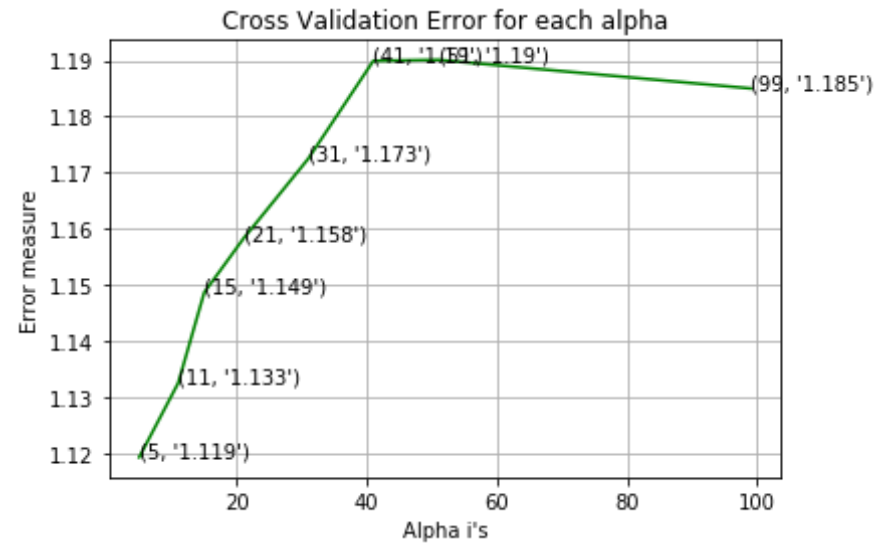plt.ylabel("Error measure")
```

```python
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 vali
dation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps
=1e-15))
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log l
oss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
for alpha = 5
Log Loss : 1.1193048504920426
for alpha = 11
Log Loss : 1.132535815490698
for alpha = 15
Log Loss : 1.1486649264772297
for alpha = 21
Log Loss : 1.1582325426147215
for alpha = 31
Log Loss : 1.1725179928027143
for alpha = 41
Log Loss : 1.1899448194991133
for alpha = 51
Log Loss : 1.1901175886177253
for alpha = 99
Log Loss : 1.1850194782310621
```

Cross Validation Error for each alpha

```
For values of best alpha =  5 The train log loss is: 0.4914817211218291
For values of best alpha =  5 The cross validation log loss is: 1.11930
48504920426
For values of best alpha =  5 The test log loss is: 1.0003934517367812
```

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

In [151]:
```python
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 : 1.1193048504920426
Number of mis-classified points : 0.3533834586466165
------------------- Confusion matrix -------------------
```

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



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

### 4.2.3.Sample Query point -1

```
In [152]:  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].resh
           ape(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 : 1
Actual Class : 7
The  5  nearest neighbours of the test points belongs to classes [7 2 7
5 7]
Fequency of nearest points : Counter({7: 3, 2: 1, 5: 1})
```

### 4.2.4. Sample Query Point-2

In [153]:
```python
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].resh
ape(1, -1), alpha[best_alpha])
print("the k value for knn is",alpha[best_alpha],"and the nearest neigh
bours 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 : 1
Actual Class : 1
the k value for knn is 5 and the nearest neighbours of the test points
belongs to classes [1 1 4 1 1]
Fequency of nearest points : Counter({1: 4, 4: 1})
```

## 4.3. Logistic Regression

### 4.3.1. With Class balancing

### 4.3.1.1. Hyper paramter tuning

```python
In [183]: 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], p
          enalty='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 vali
dation log loss is:",log_loss(y_cv, predict_y, labels=clf.classes_, eps
=1e-15))
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log l
oss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
for alpha = 1e-06
Log Loss : 1.284606477984866
for alpha = 1e-05
Log Loss : 1.2366084343258055
for alpha = 0.0001
Log Loss : 1.15187857657033992
for alpha = 0.001
Log Loss : 1.1223751580594532
for alpha = 0.01
Log Loss : 1.1602011950919846
for alpha = 0.1
Log Loss : 1.5315837758039985
for alpha = 1
Log Loss : 1.694512810475283
for alpha = 10
Log Loss : 1.7134826552379263
for alpha = 100
Log Loss : 1.7156265275846037
```

Cross Validation Error for each alpha

For values of best alpha =  0.001 The train log loss is: 0.5884672834247915
For values of best alpha =  0.001 The cross validation log loss is: 1.1223751580594532
For values of best alpha =  0.001 The test log loss is: 0.9653246965642219

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

In [184]:
```python
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y, clf)
```

Log loss : 1.1223751580594532
Number of mis-classified points : 0.37030075187969924
-------------------- Confusion matrix --------------------

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



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

### 4.3.1.3.1. Correctly Classified point

In [156]:
```python
# from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], p
enalty='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)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[7.300e-03 7.060e-02 1.200e-03 5.500e-0
3 1.180e-01 7.800e-03 7.869e-01
  2.100e-03 6.000e-04]]
Actual Class : 7
```

```
--------------------------------------------------
```

### *4.3.1.3.2. Incorrectly Classified point*

```
In [157]:  test_point_index = 101
           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)
```

```
Predicted Class : 1
Predicted Class Probabilities: [[9.118e-01 1.550e-02 2.160e-02 1.450e-0
2 1.530e-02 1.590e-02 2.800e-03
  1.700e-03 8.000e-04]]
Actual Class : 1
--------------------------------------------------
```

## 4.3.2. Without Class balancing

### 4.3.2.1. Hyper paramter tuning

```
In [158]:  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)
```

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

```
for alpha = 1e-06
Log Loss : 1.2568875500654497
for alpha = 1e-05
Log Loss : 1.2530247016024205
for alpha = 0.0001
Log Loss : 1.170492394600119
```

```
for alpha = 0.001
Log Loss : 1.1539978193390932
for alpha = 0.01
Log Loss : 1.3607347850105698
for alpha = 0.1
Log Loss : 1.7529384828524115
for alpha = 1
Log Loss : 1.833927720279602
```



Cross Validation Error for each alpha

```
For values of best alpha =  0.001 The train log loss is: 0.5838354508962968
For values of best alpha =  0.001 The cross validation log loss is: 1.1539978193390932
For values of best alpha =  0.001 The test log loss is: 1.0134600008369543
```

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

```
In [159]: 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_o
nehotCoding, cv_y, clf)
```

Log loss : 1.1539978193390932
Number of mis-classified points : 0.35150375939849626
------------------- Confusion matrix -------------------



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

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



### 4.3.2.3. Correctly Classified point

```
In [160]:  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)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[6.100e-03 5.450e-02 5.000e-04 4.200e-0
3 5.720e-02 4.500e-03 8.707e-01
  2.200e-03 1.000e-04]]
Actual Class : 7
--------------------------------------------------
```

**4.3.2.4.Inorrectly Classified point**

```
In [164]: test_point_index = 105
          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)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[0.1731 0.1498 0.0022 0.2882 0.0209 0.0
134 0.3422 0.008  0.0022]]
Actual Class : 8
--------------------------------------------------
```

# 4.4. Linear Support Vector Machines

### 4.4.1. Hyper paramter tuning

```
In [165]: 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='bal
          anced')
              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)
```

```python
    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='balance
d')
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], p
enalty='l2', loss='hinge', random_state=42)
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

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

```
for C = 1e-05
Log Loss : 1.245127653265587
for C = 0.0001
Log Loss : 1.2073056278013425
```

```
for C = 0.001
Log Loss : 1.1814023724440603
for C = 0.01
Log Loss : 1.3074369106026162
for C = 0.1
Log Loss : 1.530861498605193
for C = 1
Log Loss : 1.71621151127484
for C = 10
Log Loss : 1.716099864177392
for C = 100
Log Loss : 1.716099885072309
```



Cross Validation Error for each alpha

```
For values of best alpha =  0.001 The train log loss is: 0.508798749382
5432
For values of best alpha =  0.001 The cross validation log loss is: 1.1
814023724440603
For values of best alpha =  0.001 The test log loss is: 1.0360222881267
132
```

## 4.4.2. Testing model with best hyper parameters

```
In [166]: # 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_on
          ehotCoding,cv_y, clf)
```

```
Log loss : 1.1814023724440603
Number of mis-classified points : 0.35150375939849626
------------------- Confusion matrix -------------------
```



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

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



### 4.3.3.1. For Correctly classified point

```
In [167]:  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)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0793 0.0634 0.0049 0.0509 0.1673 0.0
215 0.6035 0.0044 0.0049]]
Actual Class : 7
--------------------------------------------------
```

### 4.3.3.2. For Incorrectly classified point

In [168]:
```
test_point_index = 101
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)
```

```
Predicted Class : 1
Predicted Class Probabilities: [[0.8379 0.0232 0.09   0.0058 0.0111 0.0
182 0.0102 0.0023 0.0014]]
Actual Class : 1
--------------------------------------------------
```

## 4.5 Random Forest Classifier

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

```python
In [169]: 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]).ra
          vel()
          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)), (featur
          es[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)], cri
          terion='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")
```

```python
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.cl
asses_, 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_, ep
s=1e-15))
```

```
for n_estimators = 100 and max depth =  5
Log Loss : 1.2710700902932879
for n_estimators = 100 and max depth =  10
Log Loss : 1.2513576468392098
for n_estimators = 200 and max depth =  5
Log Loss : 1.2651676997698436
for n_estimators = 200 and max depth =  10
Log Loss : 1.2417516365545933
for n_estimators = 500 and max depth =  5
Log Loss : 1.2570689488836735
for n_estimators = 500 and max depth =  10
Log Loss : 1.233012502544514
for n_estimators = 1000 and max depth =  5
Log Loss : 1.2579718008407204
for n_estimators = 1000 and max depth =  10
Log Loss : 1.230497002592308
for n_estimators = 2000 and max depth =  5
Log Loss : 1.2560017942074708
for n_estimators = 2000 and max depth =  10
Log Loss : 1.23058771496127
For values of best estimator =  1000 The train log loss is: 0.552309885
2930073
For values of best estimator =  1000 The cross validation log loss is:
1.230497002592308
```

For values of best estimator =  1000 The test log loss is: 1.1509158878
439314

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

In [170]:
```python
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], cri
terion='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_on
ehotCoding,cv_y, clf)
```

Log loss : 1.230497002592308
Number of mis-classified points : 0.39849624060150374
------------------- Confusion matrix -------------------



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

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 2 | 0.107 | 0.579 | | 0.045 | 0.000 | 0.000 | 0.147 | | 0.000 |
| 3 | 0.006 | 0.000 | | 0.056 | 0.071 | 0.000 | 0.037 | | 0.000 |
| 4 | 0.219 | 0.000 | | 0.685 | 0.071 | 0.038 | 0.047 | | 0.167 |
| 5 | 0.095 | 0.026 | | 0.034 | 0.500 | 0.077 | 0.053 | | 0.000 |
| 6 | 0.047 | 0.026 | | 0.034 | 0.214 | 0.885 | 0.032 | | 0.000 |
| 7 | 0.065 | 0.342 | | 0.011 | 0.071 | 0.000 | 0.668 | | 0.000 |
| 8 | 0.012 | 0.000 | | 0.000 | 0.000 | 0.000 | 0.005 | | 0.000 |
| 9 | 0.006 | 0.000 | | 0.000 | 0.000 | 0.000 | 0.000 | | 0.833 |

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

| Original Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.824 | 0.011 | 0.000 | 0.132 | 0.011 | 0.000 | 0.022 | 0.000 | 0.000 |
| 2 | 0.250 | 0.306 | 0.000 | 0.056 | 0.000 | 0.000 | 0.389 | 0.000 | 0.000 |
| 3 | 0.071 | 0.000 | 0.000 | 0.357 | 0.071 | 0.000 | 0.500 | 0.000 | 0.000 |
| 4 | 0.336 | 0.000 | 0.000 | 0.555 | 0.009 | 0.009 | 0.082 | 0.000 | 0.009 |
| 5 | 0.410 | 0.026 | 0.000 | 0.077 | 0.179 | 0.051 | 0.256 | 0.000 | 0.000 |
| 6 | 0.182 | 0.023 | 0.000 | 0.068 | 0.068 | 0.523 | 0.136 | 0.000 | 0.000 |
| 7 | 0.072 | 0.085 | 0.000 | 0.007 | 0.007 | 0.000 | 0.830 | 0.000 | 0.000 |
| 8 | 0.667 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.333 | 0.000 | 0.000 |
| 9 | 0.167 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.833 |

**4.5.3.1. Correctly Classified point**

In [171]:
```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)
```

```python
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)
```

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0741 0.1247 0.0179 0.0634 0.0519 0.0
418 0.6112 0.0065 0.0084]]
Actual Class : 7
--------------------------------------------------
```

### 4.5.3.2. Inorrectly Classified point

In [172]:
```python
test_point_index = 101
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)
```

```
Predicted Class : 1
Predicted Class Probabilities: [[0.4398 0.1    0.0237 0.1882 0.0793 0.0
88  0.0626 0.0076 0.0109]]
Actuall Class : 1
--------------------------------------------------
```

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

In [173]:
```python
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]).ra
vel()
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)), (featur
es[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)], cri
terion='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)
```

```python
predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The tra
in 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 cro
ss validation log loss is:",log_loss(y_cv, predict_y, labels=clf.classe
s_, eps=1e-15))
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The tes
t 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 : 2.37702802739506
for n_estimators = 10 and max depth =  3
Log Loss : 1.7125833203965626
for n_estimators = 10 and max depth =  5
Log Loss : 1.4417264574873057
for n_estimators = 10 and max depth =  10
Log Loss : 2.0544772021374285
for n_estimators = 50 and max depth =  2
Log Loss : 1.843630567470931
for n_estimators = 50 and max depth =  3
Log Loss : 1.5311001033017848
for n_estimators = 50 and max depth =  5
Log Loss : 1.4233533953969322
for n_estimators = 50 and max depth =  10
Log Loss : 1.846096039217289
for n_estimators = 100 and max depth =  2
Log Loss : 1.6696641296726993
for n_estimators = 100 and max depth =  3
Log Loss : 1.5650806815069216
for n_estimators = 100 and max depth =  5
Log Loss : 1.3791983080100536
for n_estimators = 100 and max depth =  10
Log Loss : 1.8153578198236895
for n_estimators = 200 and max depth =  2
Log Loss : 1.6989302743638415
```

```
Log Loss : 1.6989302743638415
for n_estimators = 200 and max depth =  3
Log Loss : 1.552698434184467
for n_estimators = 200 and max depth =  5
Log Loss : 1.3986638611030189
for n_estimators = 200 and max depth =  10
Log Loss : 1.808253869978153
for n_estimators = 500 and max depth =  2
Log Loss : 1.763224675915781
for n_estimators = 500 and max depth =  3
Log Loss : 1.5976337610317257
for n_estimators = 500 and max depth =  5
Log Loss : 1.4200065061201168
for n_estimators = 500 and max depth =  10
Log Loss : 1.7927790429098296
for n_estimators = 1000 and max depth =  2
Log Loss : 1.7398202525745232
for n_estimators = 1000 and max depth =  3
Log Loss : 1.5971012396482984
for n_estimators = 1000 and max depth =  5
Log Loss : 1.4117519115794284
for n_estimators = 1000 and max depth =  10
Log Loss : 1.7943005139062507
For values of best alpha =  100 The train log loss is: 0.05500270932277
0904
For values of best alpha =  100 The cross validation log loss is: 1.379
1983080100536
For values of best alpha =  100 The test log loss is: 1.2549000757519
```

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

In [174]:
```python
clf = RandomForestClassifier(max_depth=max_depth[int(best_alpha%4)], n_
estimators=alpha[int(best_alpha/4)], criterion='gini', max_features='au
to',random_state=42)
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y,cv_x_
responseCoding,cv_y, clf)
```

```
Log loss : 1.3791983080100536
```

Number of mis-classified points : 0.5018796992481203
-------------------- Confusion matrix --------------------



| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 34.000 | 3.000 | 2.000 | 33.000 | 12.000 | 3.000 | 1.000 | 3.000 | 0.000 |
| 2 | 0.000 | 46.000 | 1.000 | 6.000 | 0.000 | 3.000 | 10.000 | 6.000 | 0.000 |
| 3 | 0.000 | 3.000 | 4.000 | 4.000 | 2.000 | 0.000 | 1.000 | 0.000 | 0.000 |
| 4 | 6.000 | 4.000 | 1.000 | 82.000 | 9.000 | 2.000 | 0.000 | 6.000 | 0.000 |
| 5 | 3.000 | 5.000 | 1.000 | 4.000 | 18.000 | 3.000 | 4.000 | 1.000 | 0.000 |
| 6 | 3.000 | 6.000 | 0.000 | 4.000 | 8.000 | 23.000 | 0.000 | 0.000 | 0.000 |
| 7 | 0.000 | 85.000 | 9.000 | 4.000 | 1.000 | 1.000 | 51.000 | 2.000 | 0.000 |
| 8 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 2.000 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 5.000 |

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



| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.739 | 0.020 | 0.111 | 0.241 | 0.240 | 0.086 | 0.015 | 0.143 | 0.000 |
| 2 | 0.000 | 0.301 | 0.056 | 0.044 | 0.000 | 0.086 | 0.149 | 0.286 | 0.000 |
| 3 | 0.000 | 0.020 | 0.222 | 0.029 | 0.040 | 0.000 | 0.015 | 0.000 | 0.000 |
| 4 | 0.130 | 0.026 | 0.056 | 0.599 | 0.180 | 0.057 | 0.000 | 0.286 | 0.000 |
| 5 | 0.065 | 0.033 | 0.056 | 0.029 | 0.360 | 0.086 | 0.060 | 0.048 | 0.000 |
| 6 | 0.065 | 0.039 | 0.000 | 0.029 | 0.160 | 0.657 | 0.000 | 0.000 | 0.000 |
| 7 | 0.000 | 0.556 | 0.500 | 0.029 | 0.020 | 0.029 | 0.761 | 0.095 | 0.000 |

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 8 | 0.000 | 0.007 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.095 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.048 | 1.000 |

Predicted Class

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



Original Class / Predicted Class

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.374 | 0.033 | 0.022 | 0.363 | 0.132 | 0.033 | 0.011 | 0.033 | 0.000 |
| 2 | 0.000 | 0.639 | 0.014 | 0.083 | 0.000 | 0.042 | 0.139 | 0.083 | 0.000 |
| 3 | 0.000 | 0.214 | 0.286 | 0.286 | 0.143 | 0.000 | 0.071 | 0.000 | 0.000 |
| 4 | 0.055 | 0.036 | 0.009 | 0.745 | 0.082 | 0.018 | 0.000 | 0.055 | 0.000 |
| 5 | 0.077 | 0.128 | 0.026 | 0.103 | 0.462 | 0.077 | 0.103 | 0.026 | 0.000 |
| 6 | 0.068 | 0.136 | 0.000 | 0.091 | 0.182 | 0.523 | 0.000 | 0.000 | 0.000 |
| 7 | 0.000 | 0.556 | 0.059 | 0.026 | 0.007 | 0.007 | 0.333 | 0.013 | 0.000 |
| 8 | 0.000 | 0.333 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.667 | 0.000 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.167 | 0.833 |

Predicted Class

**4.5.5.1. Correctly Classified point**

In [175]:
```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(
```

```python
            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: [[0.0219 0.3347 0.1358 0.0261 0.0608 0.0
52  0.2993 0.0532 0.0162]]
Actual Class : 7
--------------------------------------------------
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Variation is important feature
Text is important feature
Text is important feature
Gene is important feature
Text is important feature
Text is important feature
Text is important feature
Gene is important feature
Gene is important feature
Variation is important feature
Text is important feature
Gene is important feature
Variation is important feature
Gene is important feature
Text is important feature
Text is important feature

Gene is important feature
```

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

**4.5.5.2. Incorrectly Classified point**

In [176]:
```python
test_point_index = 101
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 : 1
Predicted Class Probabilities: [[0.964  0.0031 0.0033 0.0077 0.0023 0.0
076 0.0019 0.0073 0.0028]]
Actual Class : 1
--------------------------------------------------
Variation is important feature
Variation is important feature
Variation is important feature
Variation is important feature
Gene is important feature
Variation is important feature
Variation is important feature
Text is important feature
Text is important feature
```

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

## 4.7 Stack the models

### 4.7.1 testing with hyper parameter tuning

```python
In [177]: clf1 = SGDClassifier(alpha=0.001, penalty='l2', loss='log', class_weigh
          t='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)
```

```python
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_cl
f1.predict_proba(cv_x_onehotCoding))))
sig_clf2.fit(train_x_onehotCoding, train_y)
print("Support vector machines : Log Loss: %0.2f" % (log_loss(cv_y, sig
_clf2.predict_proba(cv_x_onehotCoding))))
sig_clf3.fit(train_x_onehotCoding, train_y)
print("Naive Bayes : Log Loss: %0.2f" % (log_loss(cv_y, sig_clf3.predic
t_proba(cv_x_onehotCoding))))
print("-"*50)
alpha = [0.0001,0.001,0.01,0.1,1,10]
best_alpha = 999
for i in alpha:
    lr = LogisticRegression(C=i)
    sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3
], meta_classifier=lr, use_probas=True)
    sclf.fit(train_x_onehotCoding, train_y)
    print("Stacking Classifer : for the value of alpha: %f Log Loss: %
0.3f" % (i, log_loss(cv_y, sclf.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: 1.13
Support vector machines : Log Loss: 1.72
Naive Bayes : Log Loss: 1.27
--------------------------------------------------
Stacking Classifer : for the value of alpha: 0.000100 Log Loss: 2.178
Stacking Classifer : for the value of alpha: 0.001000 Log Loss: 2.034
Stacking Classifer : for the value of alpha: 0.010000 Log Loss: 1.503
Stacking Classifer : for the value of alpha: 0.100000 Log Loss: 1.163
Stacking Classifer : for the value of alpha: 1.000000 Log Loss: 1.367
Stacking Classifer : for the value of alpha: 10.000000 Log Loss: 1.774
```

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

```
In [178]: lr = LogisticRegression(C=0.1)
          sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], m
          eta_classifier=lr, use_probas=True)
          sclf.fit(train_x_onehotCoding, train_y)
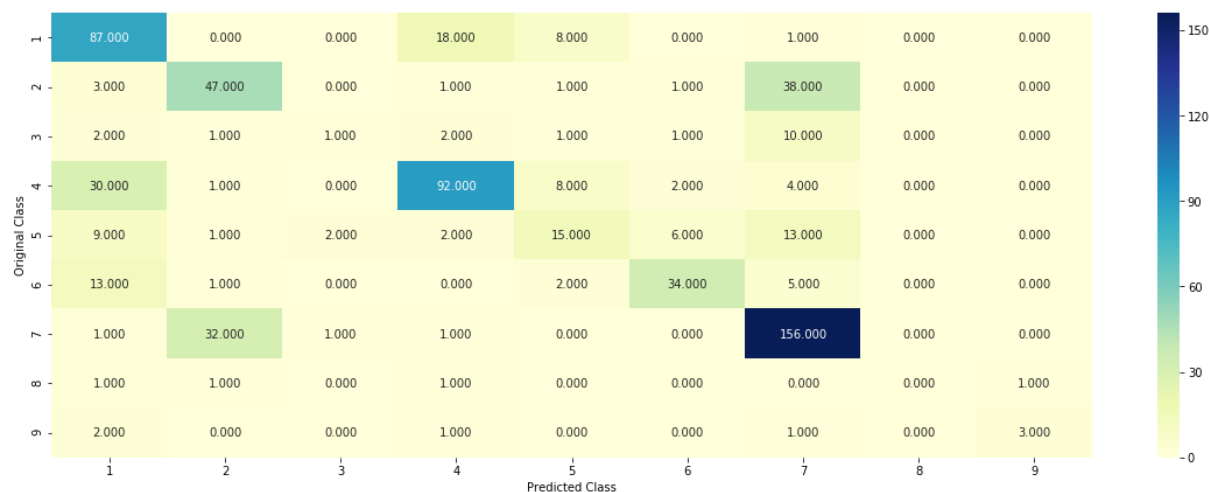
          log_error = log_loss(train_y, sclf.predict_proba(train_x_onehotCoding))
          print("Log loss (train) on the stacking classifier :",log_error)

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

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

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

```
Log loss (train) on the stacking classifier : 0.5976972830026743
Log loss (CV) on the stacking classifier : 1.163490427360028
Log loss (test) on the stacking classifier : 1.0598193030112681
Number of missclassified point : 0.3458646616541353
------------------- Confusion matrix -------------------
```

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



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



### 4.7.3 Maximum Voting classifier

```
In [179]: #Refer:http://scikit-learn.org/stable/modules/generated/sklearn.ensembl
          e.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, v
          clf.predict_proba(train_x_onehotCoding)))
          print("Log loss (CV) on the VotingClassifier :", log_loss(cv_y, vclf.pr
          edict_proba(cv_x_onehotCoding)))
          print("Log loss (test) on the VotingClassifier :", log_loss(test_y, vcl
          f.predict_proba(test_x_onehotCoding)))
          print("Number of missclassified point :", np.count_nonzero((vclf.predic
          t(test_x_onehotCoding)- test_y))/test_y.shape[0])
          plot_confusion_matrix(test_y=test_y, predict_y=vclf.predict(test_x_oneh
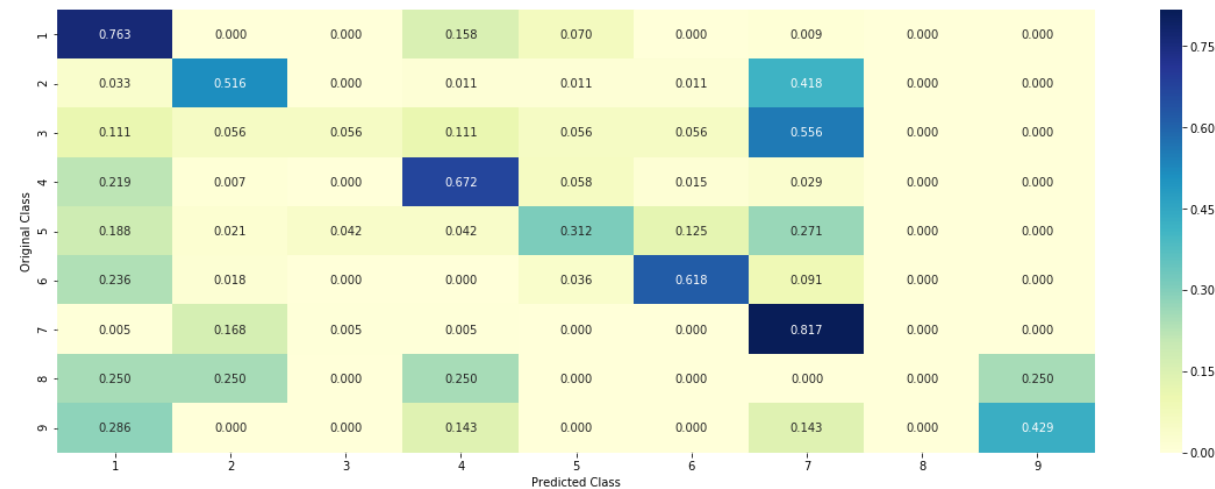          otCoding))
```

```
Log loss (train) on the VotingClassifier : 0.8298191492016891
Log loss (CV) on the VotingClassifier : 1.206389010943813
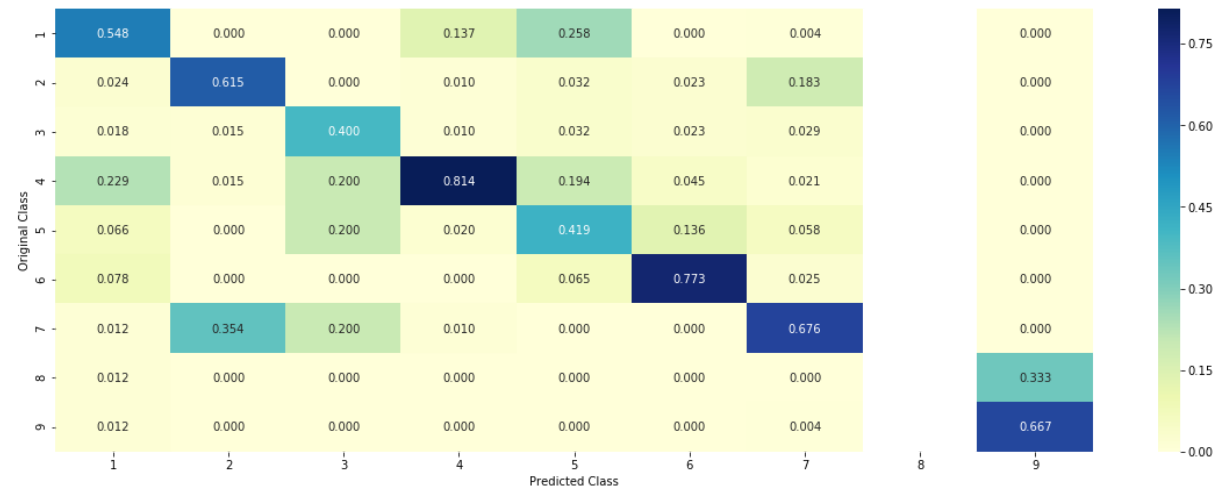Log loss (test) on the VotingClassifier : 1.125842503798544
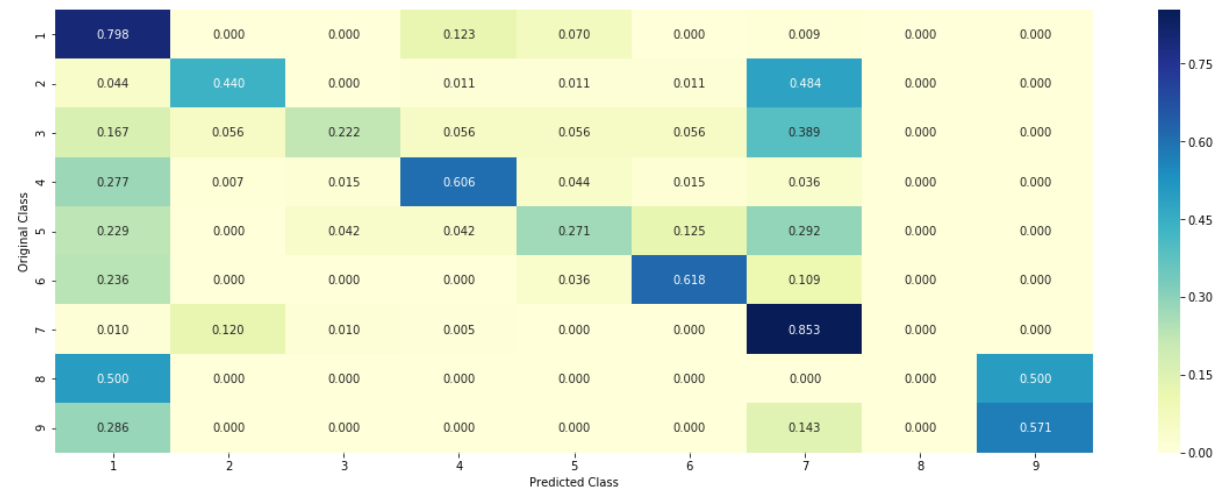Number of missclassified point : 0.35037593984962406
------------------- Confusion matrix -------------------
```



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

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.548 | 0.000 | 0.000 | 0.137 | 0.258 | 0.000 | 0.004 |  | 0.000 |
| 2 | 0.024 | 0.615 | 0.000 | 0.010 | 0.032 | 0.023 | 0.183 |  | 0.000 |
| 3 | 0.018 | 0.015 | 0.400 | 0.010 | 0.032 | 0.023 | 0.029 |  | 0.000 |
| 4 | 0.229 | 0.015 | 0.200 | 0.814 | 0.194 | 0.045 | 0.021 |  | 0.000 |
| 5 | 0.066 | 0.000 | 0.200 | 0.020 | 0.419 | 0.136 | 0.058 |  | 0.000 |
| 6 | 0.078 | 0.000 | 0.000 | 0.000 | 0.065 | 0.773 | 0.025 |  | 0.000 |
| 7 | 0.012 | 0.354 | 0.200 | 0.010 | 0.000 | 0.000 | 0.676 |  | 0.000 |
| 8 | 0.012 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |  | 0.333 |
| 9 | 0.012 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.004 |  | 0.667 |

Original Class (rows), Predicted Class (columns)

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



|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.798 | 0.000 | 0.000 | 0.123 | 0.070 | 0.000 | 0.009 | 0.000 | 0.000 |
| 2 | 0.044 | 0.440 | 0.000 | 0.011 | 0.011 | 0.011 | 0.484 | 0.000 | 0.000 |
| 3 | 0.167 | 0.056 | 0.222 | 0.056 | 0.056 | 0.056 | 0.389 | 0.000 | 0.000 |
| 4 | 0.277 | 0.007 | 0.015 | 0.606 | 0.044 | 0.015 | 0.036 | 0.000 | 0.000 |
| 5 | 0.229 | 0.000 | 0.042 | 0.042 | 0.271 | 0.125 | 0.292 | 0.000 | 0.000 |
| 6 | 0.236 | 0.000 | 0.000 | 0.000 | 0.036 | 0.618 | 0.109 | 0.000 | 0.000 |
| 7 | 0.010 | 0.120 | 0.010 | 0.005 | 0.000 | 0.000 | 0.853 | 0.000 | 0.000 |
| 8 | 0.500 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.500 |
| 9 | 0.286 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.143 | 0.000 | 0.571 |

Original Class (rows), Predicted Class (columns)

In [185]:
```python
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["S.No","Model","Train logloss","Cv logloss","Test logl
oss", "Misclassified error"]
```

```
x.add_row(["1","Naive Bayes","0.73","1.2","1.19","0.38"])
x.add_row(["2","KNN","0.49","1.1","1.0","0.353"])
x.add_row(["3","Logistic regression with C.Bal","0.56","1.12","0.98",
"0.36"])
x.add_row(["4","Logistic regression without C.Bal","0.5","1.15","1.08",
"0.35"])
x.add_row(["5","Linear svm(with one hot encoding)","0.58","1.01","1.03"
,"0.35"])
x.add_row(["6","Random Forest(with one hot encoding)","0.55","1.2","1.1
5","0.398"])
x.add_row(["7","Random Forest(with response coding)","0.55","1.3","1.2"
,"0.51"])
x.add_row(["8","Stacking classifier","0.59","1.16","1.05","0.345"])
x.add_row(["9","Maximum voting classifier","0.82","1.2","1.12","0.35"])


print(x)
```

```
+------+-------------------------------------------+---------------+--------
----+-------------+---------------------+
| S.No |                 Model                     | Train logloss | Cv logl
oss | Test logloss | Misclassified error |
+------+-------------------------------------------+---------------+--------
----+-------------+---------------------+
|  1   |                Naive Bayes                |      0.73      |    1.2
    |     1.19      |         0.38        |
|  2   |                   KNN                     |      0.49      |    1.1
    |     1.0       |        0.353        |
|  3   |      Logistic regression with C.Bal       |      0.56      |    1.12
    |     0.98      |         0.36        |
|  4   |    Logistic regression without C.Bal      |      0.5       |    1.15
    |     1.08      |         0.35        |
|  5   |    Linear svm(with one hot encoding)      |      0.58      |    1.01
    |     1.03      |         0.35        |
|  6   |  Random Forest(with one hot encoding)     |      0.55      |    1.2
    |     1.15      |        0.398        |
|  7   |  Random Forest(with response coding)      |      0.55      |    1.3
    |     1.2       |         0.51        |
|  8   |            Stacking classifier            |      0.59      |    1.16
```

```
        |      1.05       |         0.345          |
|  9    |       Maximum voting classifier          |        0.82       |    1.2
        |      1.12       |         0.35           |
+------+------------------------------------------+---------------+--------
----+-------------+--------------------+
```

# Conclusion

Here I replaced everything with Tfidf vectorizer and top 1000 features and we also apply feature engineering when we observe this Logistic Regression with class balancing has decent logloss values and reasonable Misclassification error compared to all other.. It seems slightly unstable but better than all other ...

Here ,Due to feature engineering we can observe some of the missclassification error gets decent value but slightly increases and also many tain logloss reduces to less than 1 for Logistic regression we also have test loss less than 1..

While due to this feature engineering some of them are getting good logloss values while some are increasing