

# Personalized cancer diagnosis

## 1. Business Problem ¶

### 1.1. Description

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

#### **Context:**

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

#### **Problem statement :**

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

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

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

## 2. Machine Learning Problem Formulation

### 2.1. Data

#### 2.1.1. Data Overview

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

## 2.1.2. Example Data Point

### *training\_variants*

---

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

### *training\_text*

---

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

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

### 2.2.1. Type of Machine Learning Problem

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

### 2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

### 2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

## 2.3. Train, CV and Test Datasets

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

## 3. Exploratory Data Analysis

```
In [4]: import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import numpy as np
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import SGDClassifier
from imblearn.over_sampling import SMOTE
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC

from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier
warnings.filterwarnings("ignore")

from mlxtend.classifier import StackingClassifier

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

## 3.1. Reading Data

### 3.1.1. Reading Gene and Variation Data

```
In [5]: data = pd.read_csv('training_variants')
print('Number of data points : ', data.shape[0])
print('Number of features : ', data.shape[1])
print('Features : ', data.columns.values)
data.head()
```

Number of data points : 3321

Number of features : 4

Features : ['ID' 'Gene' 'Variation' 'Class']

Out[5]:

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

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

Fields are

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

### 3.1.2. Reading Text Data

```
In [7]: # note the separator in this file
data_text = pd.read_csv("training_text", sep="\|", engine="python", names=["ID", "TEXT"])
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[7]:

	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 [8]: # Loading stop words from nltk library
stop_words = set(stopwords.words('english'))

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

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

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

```
In [9]: #text processing stage.
start_time = time.clock()
for index, row in data_text.iterrows():
    if type(row['TEXT']) is str:
        nlp_preprocessing(row['TEXT'], index, 'TEXT')
    else:
        print("there is no text description for id:",index)
print('Time took for preprocessing the text :',time.clock() - start_time, "second
```

```
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 : 53.09783729893066 seconds
```

```
In [10]: #merging both gene_variations and text data based on ID
result = pd.merge(data, data_text,on='ID', how='left')
result.head()
```

Out[10]:

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

Out[11]:

	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 [12]: result.loc[result['TEXT'].isnull(), 'TEXT'] = result['Gene'] + ' '+result['Variatio
```

```
In [13]: result[result['ID']==1109]
```

Out[13]:

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

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

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

```
In [15]: 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 [21]: # it returns a dict, keys as class labels and values as the number of data points
train_class_distribution = train_df['Class'].value_counts().sort_values()
test_class_distribution = test_df['Class'].value_counts().sort_values()
cv_class_distribution = cv_df['Class'].value_counts().sort_values()

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

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

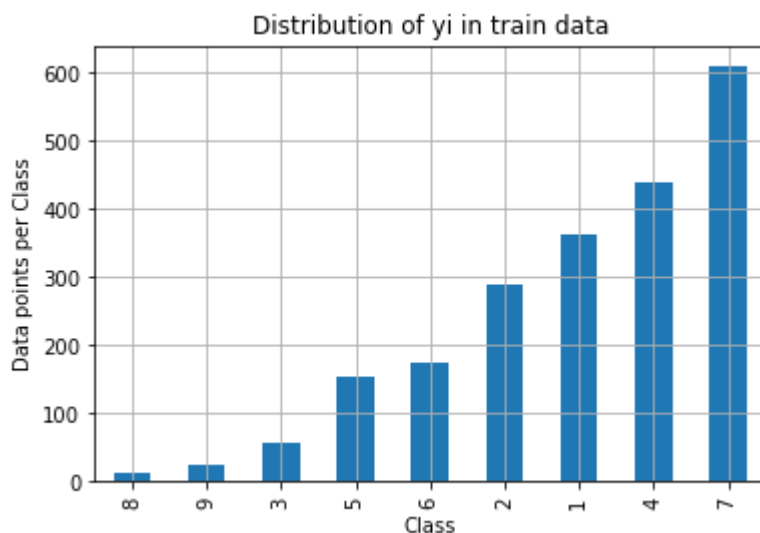
print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

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

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

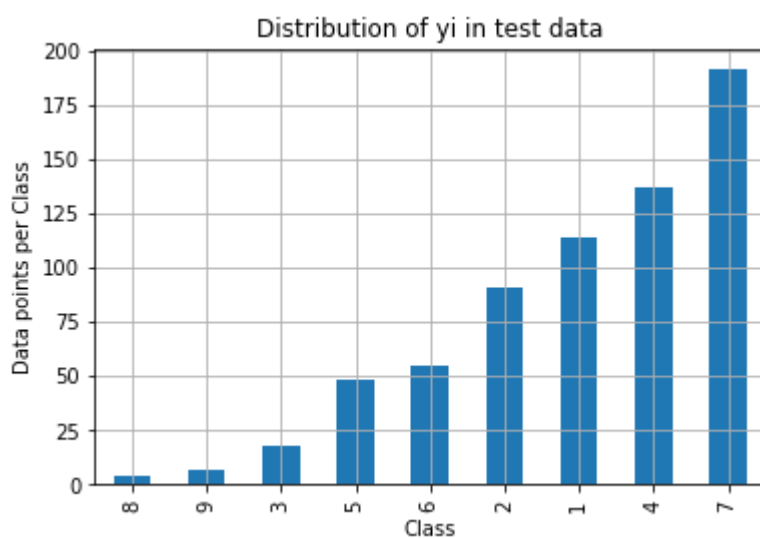
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':', cv_class_distribution.values[i])

```



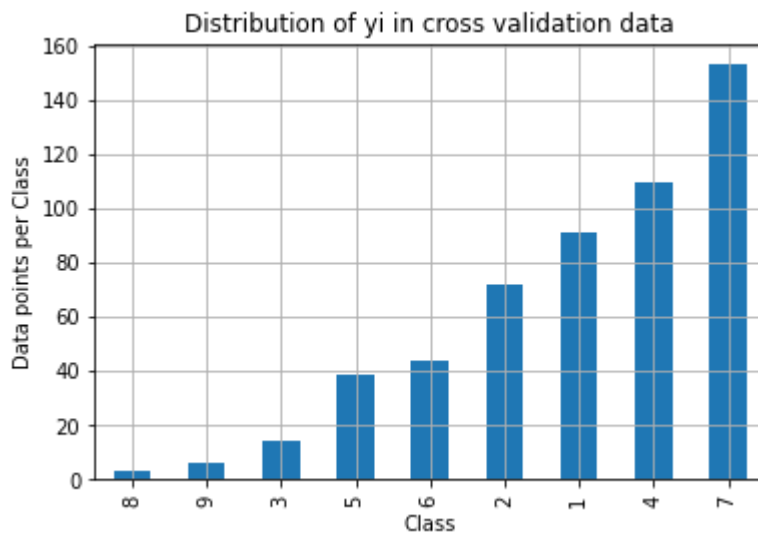
Number of data points in class 9 : 609 ( 28.672 %)  
 Number of data points in class 8 : 439 ( 20.669 %)  
 Number of data points in class 7 : 363 ( 17.09 %)  
 Number of data points in class 6 : 289 ( 13.606 %)  
 Number of data points in class 5 : 176 ( 8.286 %)  
 Number of data points in class 4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

---



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

---



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

## 3.2 Prediction using a 'Random' Model

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

```
In [22]: # 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)

    A = (((C.T)/(C.sum(axis=1))).T)

    B = (C/C.sum(axis=0))

    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, ytic
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()

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

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

In [23]:

```
# 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_pre

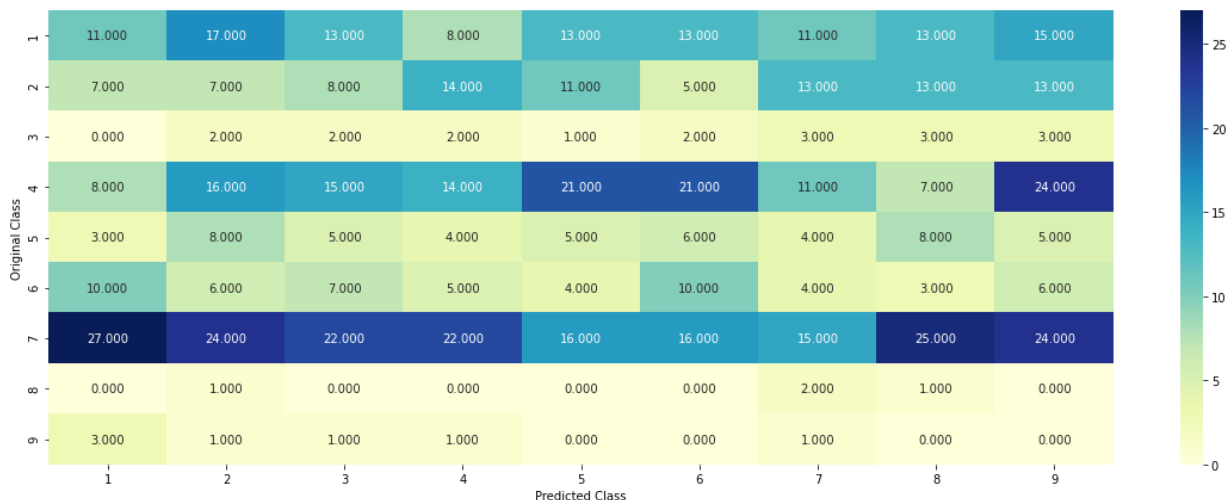
# Test-Set error.
#we create a output array that has exactly same as the test data
test_predicted_y = np.zeros((test_data_len,9))
for i in range(test_data_len):
    rand_probs = np.random.rand(1,9)
    test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs))))[0])
print("Log loss on Test Data using Random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

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 2.4269863241695178

Log loss on Test Data using Random Model 2.41589352708019

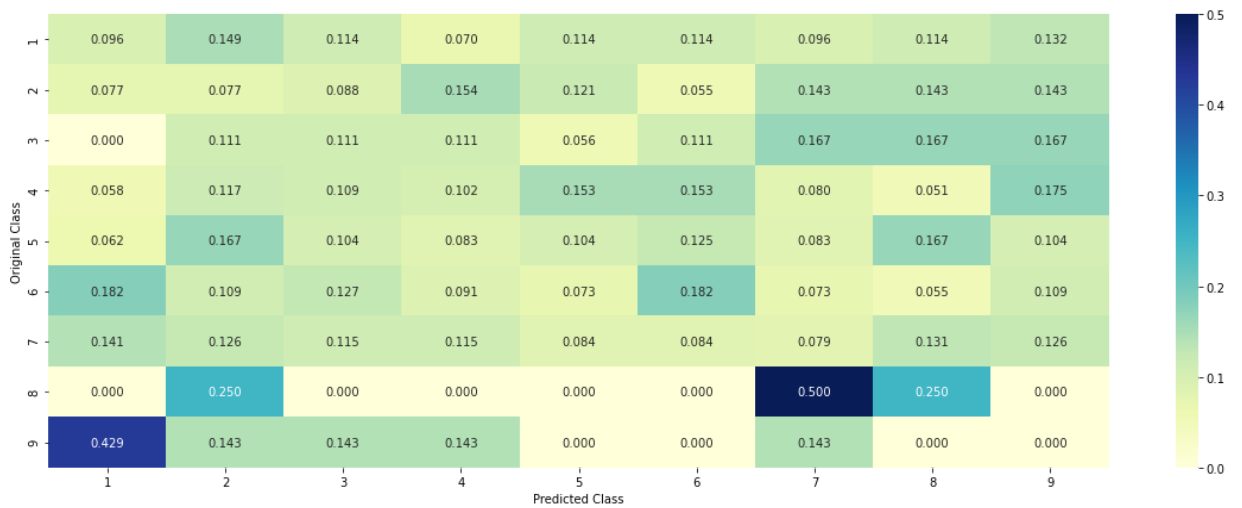
----- Confusion matrix -----



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



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



### 3.3 Univariate Analysis

In [24]:

```

# get_gv_fea_dict: Get Gene variation Feature Dict
def get_gv_fea_dict(alpha, feature, df):

    value_count = train_df[feature].value_counts()

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

    # denominator will contain the number of time that particular feature occurred
    for i, denominator in value_count.items():
        # vec will contain (p(yi==1/Gi) probability of gene/variation belongs to particular class)
        # vec is 9 dimensional vector
        vec = []
        for k in range(1,10):
            cls_cnt = train_df.loc[(train_df['Class']==k) & (train_df[feature]==i)].shape[0]

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

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

# Get Gene variation feature
def get_gv_feature(alpha, feature, df):

    gv_dict = get_gv_fea_dict(alpha, feature, df)
    # value_count is similar in get_gv_fea_dict
    value_count = train_df[feature].value_counts()

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

    return gv_fea

```

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

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

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

Number of Unique Genes : 237

BRCA1 179

TP53 103

EGFR 82

PTEN 82

BRCA2 72

KIT 63

BRAF 56

ERBB2 45

ALK 45

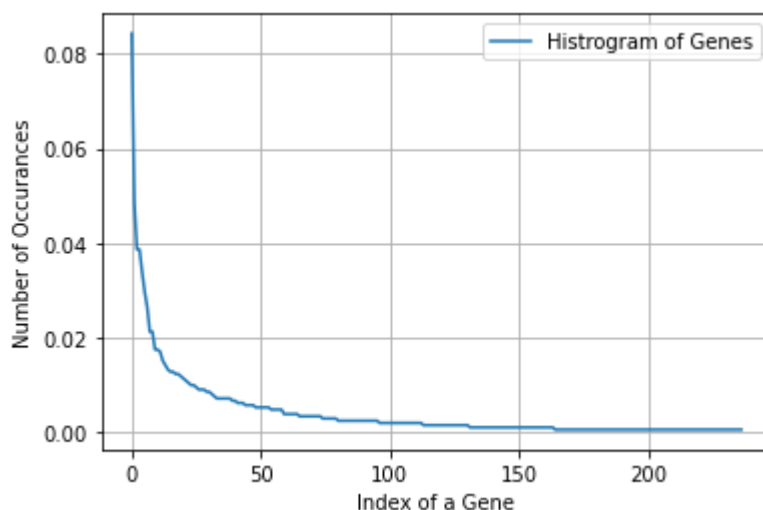
CDKN2A 37

Name: Gene, dtype: int64

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

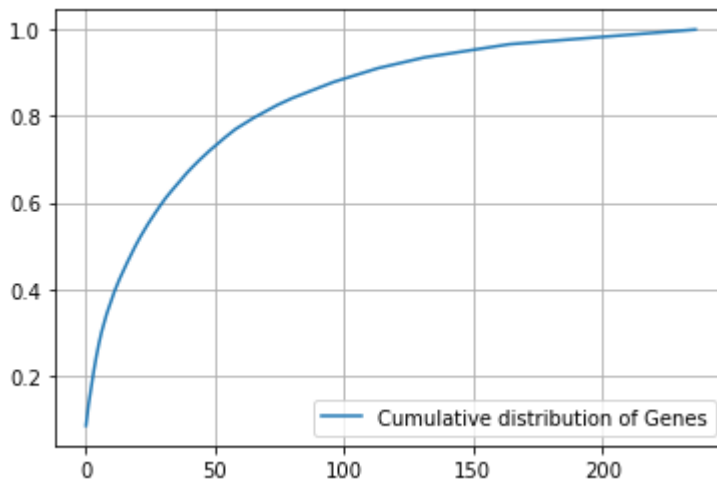
Ans: There are 229 different categories of genes in the train data, and they are distributed as follows

```
In [26]: s = sum(unique_genes.values);
h = unique_genes.values/s;
plt.plot(h, label="Histogram of Genes")
plt.xlabel('Index of a Gene')
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```





```
In [27]: 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 ?

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 [28]: #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 [29]: print("train_gene_feature_responseCoding is converted feature using response coding")
```

train\_gene\_feature\_responseCoding is converted feature using response coding method. The shape of gene feature: (2124, 9)

```
In [30]: # one-hot encoding of Gene feature.
gene_vectorizer = CountVectorizer()
train_gene_feature_onehotCoding = gene_vectorizer.fit_transform(train_df['Gene'])
test_gene_feature_onehotCoding = gene_vectorizer.transform(test_df['Gene'])
cv_gene_feature_onehotCoding = gene_vectorizer.transform(cv_df['Gene'])
```

```
In [31]: train_df['Gene'].head()
```

```
Out[31]: 413      TP53  
         5       CBL  
         88      RYBP  
        1195    PIK3CA  
        3062    MED12  
         Name: Gene, dtype: object
```

```
In [33]: print("train_gene_feature_onehotCoding is converted feature using one-hot encoding")
```

```
train_gene_feature_onehotCoding is converted feature using one-hot encoding method. The shape of gene feature: (2124, 237)
```

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

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

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

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

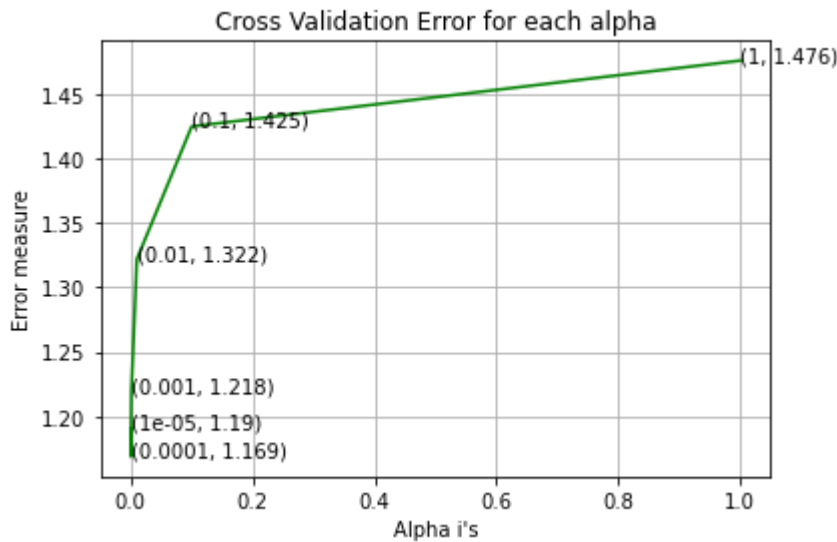
predict_y = sig_clf.predict_proba(train_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(cv_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_gene_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",

```

```

For values of alpha = 1e-05 The log loss is: 1.1901198329691247
For values of alpha = 0.0001 The log loss is: 1.1685824305462422
For values of alpha = 0.001 The log loss is: 1.2184519609886837
For values of alpha = 0.01 The log loss is: 1.3217444834733574
For values of alpha = 0.1 The log loss is: 1.4247488195515197
For values of alpha = 1 The log loss is: 1.4758292218934486

```



For values of best alpha = 0.0001 The train log loss is: 0.9812520083686037  
 For values of best alpha = 0.0001 The cross validation log loss is: 1.1685824305462422  
 For values of best alpha = 0.0001 The test log loss is: 1.2027569957312394

**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 [35]: print("Q6. How many data points in Test and CV datasets are covered by the ", uni
test_coverage=test_df[test_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]
cv_coverage=cv_df[cv_df['Gene'].isin(list(set(train_df['Gene'])))].shape[0]

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

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

Ans

1. In test data 646 out of 665 : 97.14285714285714
2. In cross validation data 515 out of 532 : 96.80451127819549

### 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?

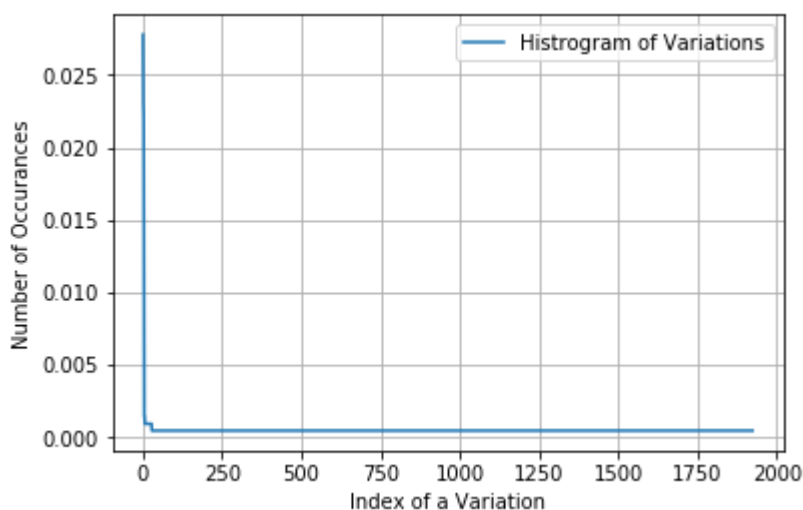
```
In [36]: unique_variations = train_df['Variation'].value_counts()
print('Number of Unique Variations :', unique_variations.shape[0])
# the top 10 variations that occurred most
print(unique_variations.head(10))
```

```
Number of Unique Variations : 1928
Truncating_Mutations      60
Deletion                  48
Amplification             43
Fusions                   25
Overexpression            4
G12V                      4
Q61R                      3
Q61H                      3
Q61L                      3
EWSR1-ETV1_Fusion        2
Name: Variation, dtype: int64
```

```
In [0]: print("Ans: There are", unique_variations.shape[0] , "different categories of vari
```

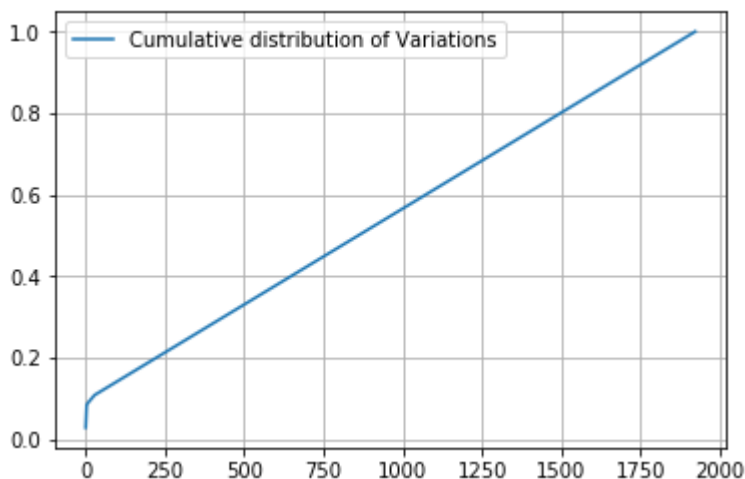
Ans: There are 1924 different categories of variations in the train data, and they are distributed as follows

```
In [0]: s = sum(unique_variations.values);
h = unique_variations.values/s;
plt.plot(h, label="Histogram of Variations")
plt.xlabel('Index of a Variation')
plt.ylabel('Number of Occurances')
plt.legend()
plt.grid()
plt.show()
```



```
In [0]: c = np.cumsum(h)
print(c)
plt.plot(c,label='Cumulative distribution of Variations')
plt.grid()
plt.legend()
plt.show()
```

```
[0.02777778 0.05084746 0.07297552 ... 0.99905838 0.99952919 1.          ]
```



## Q9. How to featurize this Variation feature ?

**Ans.** There are two ways we can featurize this variable

1. One hot Encoding
2. Response coding

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

```
In [37]: # alpha is used for Laplace smoothing
alpha = 1
# train gene feature
train_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variatio
# test gene feature
test_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation
# cross validation gene feature
cv_variation_feature_responseCoding = np.array(get_gv_feature(alpha, "Variation",
```

```
In [0]: print("train_variation_feature_responseCoding is a converted feature using the re
```

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

```
In [39]: # one-hot encoding of variation feature.  
variation_vectorizer = CountVectorizer()  
train_variation_feature_onehotCoding = variation_vectorizer.fit_transform(train_d  
test_variation_feature_onehotCoding = variation_vectorizer.transform(test_df['Var  
cv_variation_feature_onehotCoding = variation_vectorizer.transform(cv_df['Variati
```

```
In [0]: print("train_variation_feature_onehotEncoded is converted feature using the onne-
```

train\_variation\_feature\_onehotEncoded is converted feature using the onne-hot e  
ncoding method. The shape of Variation feature: (2124, 1960)

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

Let's build a model just like the earlier!

```

In [40]: alpha = [10 ** x for x in range(-5, 1)]

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

    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(train_variation_feature_onehotCoding, y_train)
    predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)

    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=clf.classes_, eps=
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_cv, predict_

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

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

predict_y = sig_clf.predict_proba(train_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(cv_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_variation_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", 1

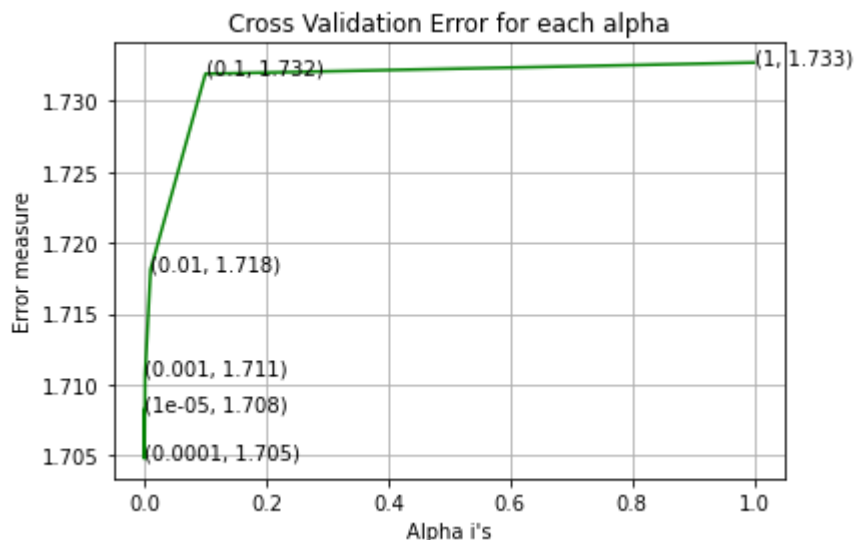
```

```

For values of alpha = 1e-05 The log loss is: 1.7082431858288039
For values of alpha = 0.0001 The log loss is: 1.7047589176085123
For values of alpha = 0.001 The log loss is: 1.7106904073827571
For values of alpha = 0.01 The log loss is: 1.71812568029598
For values of alpha = 0.1 The log loss is: 1.731925607343552
For values of alpha = 1 The log loss is: 1.7327170913083907

```





For values of best alpha = 0.0001 The train log loss is: 0.6472958168887896  
 For values of best alpha = 0.0001 The cross validation log loss is: 1.7047589176085123  
 For values of best alpha = 0.0001 The test log loss is: 1.709905394636484

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

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

```
In [41]: print("Q12. How many data points are covered by total ", unique_variations.shape[0], " genes in test and cross v
          test_coverage=test_df[test_df['Variation'].isin(list(set(train_df['Variation'])))].shape[0]
          cv_coverage=cv_df[cv_df['Variation'].isin(list(set(train_df['Variation'])))].shape[0]
          print('Ans\n1. In test data', test_coverage, 'out of', test_df.shape[0], ":", (test_coverage/test_df.shape[0])*100, "%")
          print('2. In cross validation data', cv_coverage, 'out of ', cv_df.shape[0], ":", (cv_coverage/cv_df.shape[0])*100, "%")
```

Q12. How many data points are covered by total 1928 genes in test and cross validation data sets?

Ans

1. In test data 64 out of 665 : 9.624060150375941
2. In cross validation data 57 out of 532 : 10.714285714285714

### 3.2.3 Univariate Analysis on Text Feature

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

```
In [42]: # cls_text is a data frame
# for every row in data fram consider the 'TEXT'
# split the words by space
# make a dict with those words
# increment its count whenever we see that word

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

```
In [47]: import math
#https://stackoverflow.com/a/1602964
def get_text_responsecoding(df):
    text_feature_responseCoding = np.zeros((df.shape[0],9))
    for i in range(0,9):
        row_index = 0
        for index, row in df.iterrows():
            sum_prob = 0
            for word in row['TEXT'].split():
                sum_prob += math.log(((dict_list[i].get(word,0)+10 ))/(total_dict.
            text_feature_responseCoding[row_index][i] = math.exp(sum_prob/len(row
            row_index += 1
    return text_feature_responseCoding
```

```
In [48]: # building a CountVectorizer with all the words that occurred minimum 3 times in t
text_vectorizer = CountVectorizer(min_df=3)
train_text_feature_onehotCoding = text_vectorizer.fit_transform(train_df['TEXT'])
# getting all the feature names (words)
train_text_features= text_vectorizer.get_feature_names()

# train_text_feature_onehotCoding.sum(axis=0).A1 will sum every row and returns (
train_text_fea_counts = train_text_feature_onehotCoding.sum(axis=0).A1

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

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

Total number of unique words in train data : 53968

```
In [51]: dict_list = []
# dict_list =[] contains 9 dictionaries each corresponds to a class
for i in range(1,10):
    cls_text = train_df[train_df['Class']==i]
    # build a word dict based on the words in that class
    dict_list.append(extract_dictionary_paddle(cls_text))
    # append it to dict_list

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

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

```
In [52]: #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 [53]: # 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.T).T
test_text_feature_responseCoding = (test_text_feature_responseCoding.T/test_text_feature_responseCoding.T).T
cv_text_feature_responseCoding = (cv_text_feature_responseCoding.T/cv_text_feature_responseCoding.T).T
```

```
In [54]: # normalizing every feature
train_text_feature_onehotCoding = normalize(train_text_feature_onehotCoding, axis=0)

# we use the same vectorizer that was trained on train data
test_text_feature_onehotCoding = text_vectorizer.transform(test_df['TEXT'])
# normalizing every feature
test_text_feature_onehotCoding = normalize(test_text_feature_onehotCoding, axis=0)

# we use the same vectorizer that was trained on train data
cv_text_feature_onehotCoding = text_vectorizer.transform(cv_df['TEXT'])
# normalizing every feature
cv_text_feature_onehotCoding = normalize(cv_text_feature_onehotCoding, axis=0)
```

```
In [76]: #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 [57]: *# Number of words for a given frequency.*

```
print(Counter(sorted_text_occur))
```

```
0, 33: 40, 104: 43, 117: 44, 121: 42, 103: 42, 148: 41, 133: 40, 133: 40, 11
6: 40, 109: 40, 127: 39, 124: 39, 123: 39, 156: 38, 136: 38, 151: 37, 130: 3
7, 122: 36, 160: 35, 132: 35, 142: 34, 138: 34, 106: 34, 155: 33, 149: 33, 14
1: 33, 125: 33, 118: 33, 200: 32, 169: 32, 137: 32, 134: 32, 126: 32, 129: 3
1, 174: 30, 145: 30, 168: 29, 139: 29, 187: 28, 164: 28, 161: 28, 159: 28, 15
4: 28, 131: 28, 128: 28, 180: 27, 163: 27, 158: 27, 234: 26, 208: 26, 197: 2
6, 194: 26, 181: 26, 166: 26, 152: 26, 146: 26, 143: 26, 266: 25, 186: 25, 17
8: 25, 150: 25, 144: 25, 222: 24, 204: 24, 175: 24, 223: 23, 218: 23, 211: 2
3, 201: 23, 195: 23, 192: 23, 189: 23, 177: 23, 173: 23, 162: 23, 203: 22, 19
6: 22, 171: 22, 165: 22, 213: 21, 198: 21, 176: 21, 167: 21, 267: 20, 231: 2
0, 228: 20, 216: 20, 207: 20, 191: 20, 184: 20, 182: 20, 170: 20, 236: 19, 22
5: 19, 212: 19, 210: 19, 206: 19, 270: 18, 246: 18, 230: 18, 229: 18, 221: 1
8, 209: 18, 205: 18, 157: 18, 287: 17, 272: 17, 258: 17, 248: 17, 243: 17, 22
0: 17, 214: 17, 309: 16, 245: 16, 232: 16, 217: 16, 190: 16, 188: 16, 183: 1
6, 341: 15, 294: 15, 291: 15, 282: 15, 280: 15, 279: 15, 276: 15, 264: 15, 26
3: 15, 261: 15, 254: 15, 250: 15, 227: 15, 199: 15, 179: 15, 172: 15, 330: 1
4, 300: 14, 289: 14, 284: 14, 268: 14, 259: 14, 249: 14, 215: 14, 370: 13, 35
9: 13, 316: 13, 315: 13, 312: 13, 307: 13, 298: 13, 295: 13, 273: 13, 251: 1
3, 247: 13, 244: 13, 202: 13, 469: 12, 412: 12, 361: 12, 353: 12, 351: 12, 30
5: 12, 303: 12, 299: 12, 269: 12, 265: 12, 257: 12, 252: 12, 242: 12, 241: 1
0, 225: 10, 221: 10, 220: 10, 219: 10, 218: 10, 217: 10, 216: 10, 215: 10, 214: 10, 213: 10, 212: 10, 211: 10, 210: 10, 209: 10, 208: 10, 207: 10, 206: 10, 205: 10, 204: 10, 203: 10, 202: 10, 201: 10, 200: 10, 199: 10, 198: 10, 197: 10, 196: 10, 195: 10, 194: 10, 193: 10, 192: 10, 191: 10, 190: 10, 189: 10, 188: 10, 187: 10, 186: 10, 185: 10, 184: 10, 183: 10, 182: 10, 181: 10, 180: 10, 179: 10, 178: 10, 177: 10, 176: 10, 175: 10, 174: 10, 173: 10, 172: 10, 171: 10, 170: 10, 169: 10, 168: 10, 167: 10, 166: 10, 165: 10, 164: 10, 163: 10, 162: 10, 161: 10, 160: 10, 159: 10, 158: 10, 157: 10, 156: 10, 155: 10, 154: 10, 153: 10, 152: 10, 151: 10, 150: 10, 149: 10, 148: 10, 147: 10, 146: 10, 145: 10, 144: 10, 143: 10, 142: 10, 141: 10, 140: 10, 139: 10, 138: 10, 137: 10, 136: 10, 135: 10, 134: 10, 133: 10, 132: 10, 131: 10, 130: 10, 129: 10, 128: 10, 127: 10, 126: 10, 125: 10, 124: 10, 123: 10, 122: 10, 121: 10, 120: 10, 119: 10, 118: 10, 117: 10, 116: 10, 115: 10, 114: 10, 113: 10, 112: 10, 111: 10, 110: 10, 109: 10, 108: 10, 107: 10, 106: 10, 105: 10, 104: 10, 103: 10, 102: 10, 101: 10, 100: 10, 99: 10, 98: 10, 97: 10, 96: 10, 95: 10, 94: 10, 93: 10, 92: 10, 91: 10, 90: 10, 89: 10, 88: 10, 87: 10, 86: 10, 85: 10, 84: 10, 83: 10, 82: 10, 81: 10, 80: 10, 79: 10, 78: 10, 77: 10, 76: 10, 75: 10, 74: 10, 73: 10, 72: 10, 71: 10, 70: 10, 69: 10, 68: 10, 67: 10, 66: 10, 65: 10, 64: 10, 63: 10, 62: 10, 61: 10, 60: 10, 59: 10, 58: 10, 57: 10, 56: 10, 55: 10, 54: 10, 53: 10, 52: 10, 51: 10, 50: 10, 49: 10, 48: 10, 47: 10, 46: 10, 45: 10, 44: 10, 43: 10, 42: 10, 41: 10, 40: 10, 39: 10, 38: 10, 37: 10, 36: 10, 35: 10, 34: 10, 33: 10, 32: 10, 31: 10, 30: 10, 29: 10, 28: 10, 27: 10, 26: 10, 25: 10, 24: 10, 23: 10, 22: 10, 21: 10, 20: 10, 19: 10, 18: 10, 17: 10, 16: 10, 15: 10, 14: 10, 13: 10, 12: 10, 11: 10, 10: 10, 9: 10, 8: 10, 7: 10, 6: 10, 5: 10, 4: 10, 3: 10, 2: 10, 1: 10, 0: 10
```

```

In [58]: # Train a Logistic regression+Calibration model using text features which are on-
alpha = [10 ** x for x in range(-5, 1)]

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

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

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

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

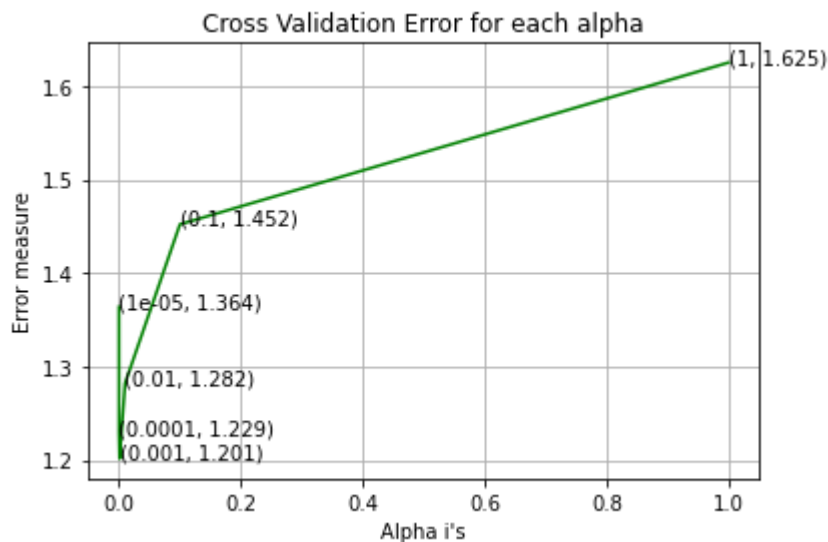
predict_y = sig_clf.predict_proba(train_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(cv_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_text_feature_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", 1

```

```

For values of alpha = 1e-05 The log loss is: 1.36361870374097
For values of alpha = 0.0001 The log loss is: 1.2285895647516927
For values of alpha = 0.001 The log loss is: 1.2014457912100782
For values of alpha = 0.01 The log loss is: 1.281986991738801
For values of alpha = 0.1 The log loss is: 1.451908981735228
For values of alpha = 1 The log loss is: 1.6247969033234715

```



For values of best alpha = 0.001 The train log loss is: 0.625889574286602

For values of best alpha = 0.001 The cross validation log loss is: 1.2014457912100782

For values of best alpha = 0.001 The test log loss is: 1.177458263252274

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

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

```
In [59]: def get_intersec_text(df):
df_text_vec = CountVectorizer(min_df=3)
df_text_fea = df_text_vec.fit_transform(df['TEXT'])
df_text_features = df_text_vec.get_feature_names()

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

```
In [60]: len1, len2 = get_intersec_text(test_df)
print(np.round((len2/len1)*100, 3), "% of word of test data appeared in train data")
len1, len2 = get_intersec_text(cv_df)
print(np.round((len2/len1)*100, 3), "% of word of Cross Validation appeared in train data")
```

97.738 % of word of test data appeared in train data

98.281 % of word of Cross Validation appeared in train data

## 4. Machine Learning Models

In [61]: *#Data preparation for ML models.*

*#Misc. functions for ML models*

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

    # for calculating log_loss we will provide the array of probabilities belonging 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)))
    plot_confusion_matrix(test_y, pred_y)
```

In [62]: **def** report\_log\_loss(train\_x, train\_y, test\_x, test\_y, clf):

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

```

In [85]: # this function will be used just for naive bayes
# for the given indices, we will print the name of the features
# and we will check whether the feature present in the test point text or not
def get_impfeature_names(indices, text, gene, var, no_features):
    gene_count_vec = CountVectorizer()
    var_count_vec = CountVectorizer()
    text_count_vec = CountVectorizer(min_df=3)

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

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

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

    print("Out of the top ",no_features," features ", word_present, "are present

```

## Stacking the three types of features



In [64]: *# merging gene, variance and text features*

```
train_gene_var_onehotCoding = hstack((train_gene_feature_onehotCoding, train_varia
test_gene_var_onehotCoding = hstack((test_gene_feature_onehotCoding, test_variatio
cv_gene_var_onehotCoding = hstack((cv_gene_feature_onehotCoding, cv_variation_feat

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

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

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

train_gene_var_responseCoding = np.hstack((train_gene_feature_responseCoding, trai
test_gene_var_responseCoding = np.hstack((test_gene_feature_responseCoding, test_v
cv_gene_var_responseCoding = np.hstack((cv_gene_feature_responseCoding, cv_variatio

train_x_responseCoding = np.hstack((train_gene_var_responseCoding, train_text_fea
test_x_responseCoding = np.hstack((test_gene_var_responseCoding, test_text_featur
cv_x_responseCoding = np.hstack((cv_gene_var_responseCoding, cv_text_feature_resp
```

In [65]: `print("One hot encoding features :")`  
`print("(number of data points * number of features) in train data = ", train_x_on`  
`print("(number of data points * number of features) in test data = ", test_x_oneh`  
`print("(number of data points * number of features) in cross validation data =",`

```
One hot encoding features :
(number of data points * number of features) in train data = (2124, 56155)
(number of data points * number of features) in test data = (665, 56155)
(number of data points * number of features) in cross validation data = (532, 5
6155)
```

In [66]: `print(" Response encoding features :")`  
`print("(number of data points * number of features) in train data = ", train_x_re`  
`print("(number of data points * number of features) in test data = ", test_x_resp`  
`print("(number of data points * number of features) in cross validation data =",`

```
Response encoding features :
(number of data points * number of features) in train data = (2124, 27)
(number of data points * number of features) in test data = (665, 27)
(number of data points * number of features) in cross validation data = (532, 2
7)
```

## 4.1. Base Line Model

### 4.1.1. Naive Bayes

#### 4.1.1.1. Hyper parameter tuning

In [67]:

```

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

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

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

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1

```

```

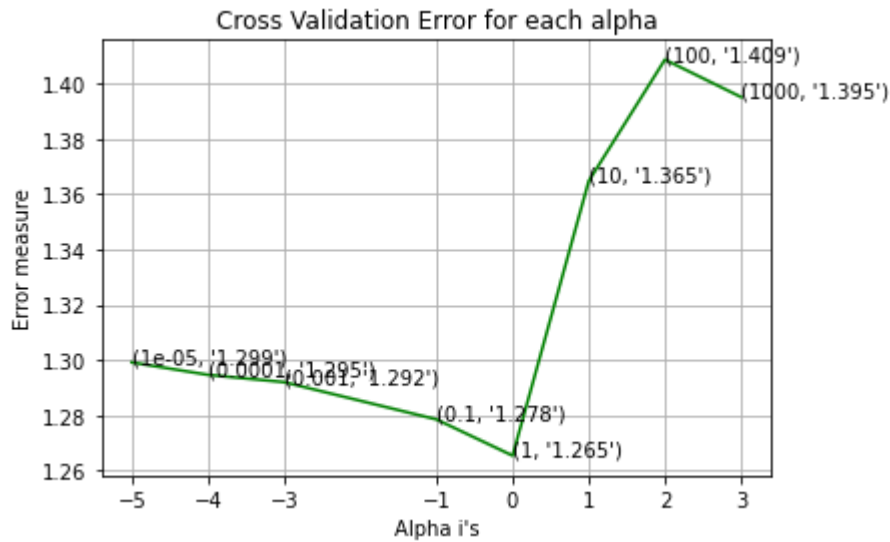
for alpha = 1e-05
Log Loss : 1.299124812093323
for alpha = 0.0001
Log Loss : 1.2946307612790837
for alpha = 0.001
Log Loss : 1.291985139550523
for alpha = 0.1
Log Loss : 1.2784862241644395
for alpha = 1
Log Loss : 1.2653570618225016
for alpha = 10
Log Loss : 1.3648165427218035
for alpha = 100

```

Log Loss : 1.4087839749711353

for alpha = 1000

Log Loss : 1.3952182404724107



For values of best alpha = 1 The train log loss is: 0.8857842350022214

For values of best alpha = 1 The cross validation log loss is: 1.2653570618225016

For values of best alpha = 1 The test log loss is: 1.3288125990526927

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

In [68]:

```

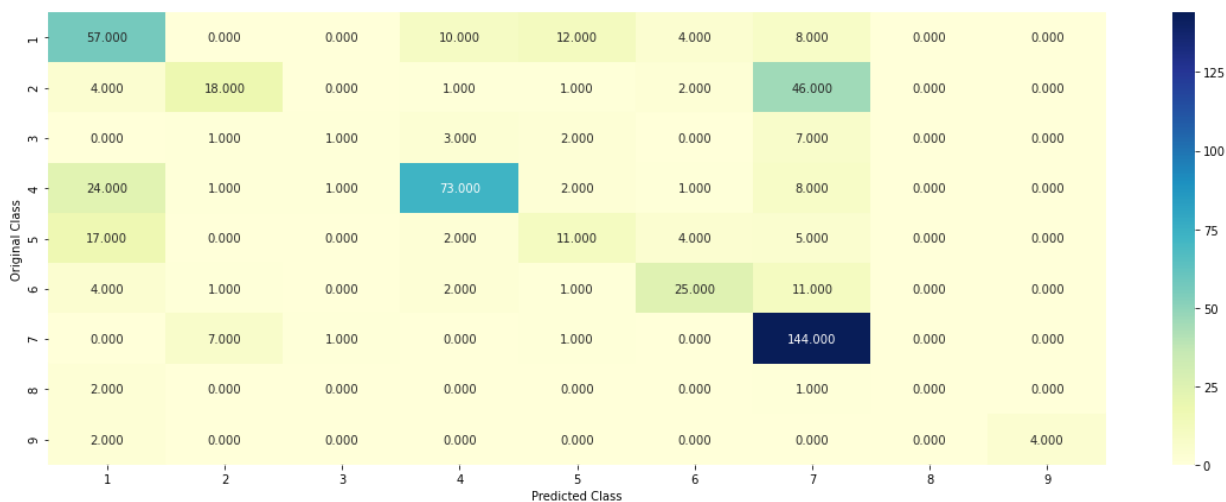
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)
sig_clf_probs = sig_clf.predict_proba(cv_x_onehotCoding)
# to avoid rounding error while multiplying probabilities we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

```

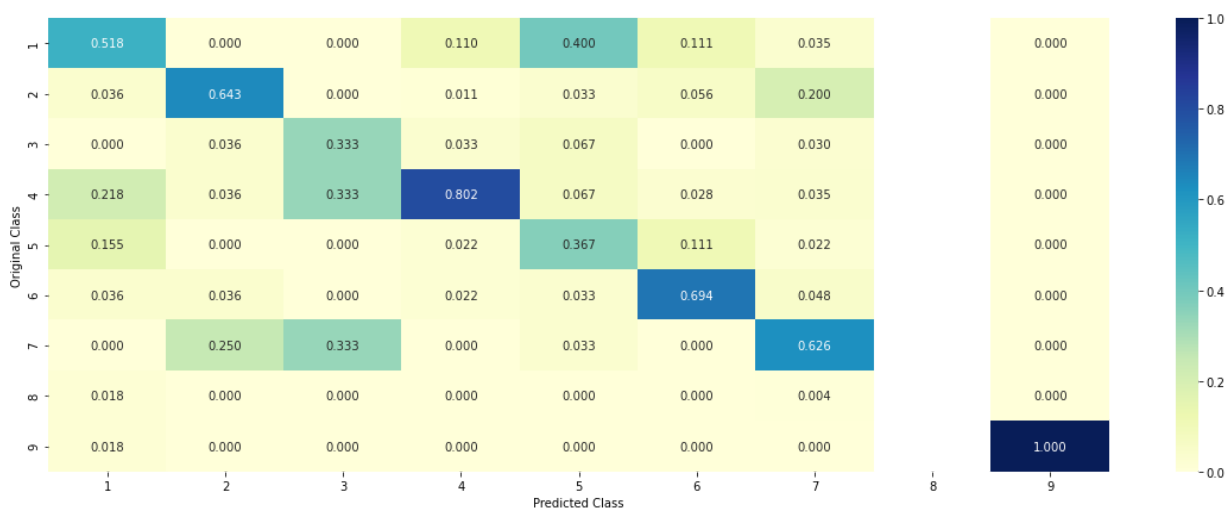
Log Loss : 1.2653570618225016

Number of missclassified point : 0.37406015037593987

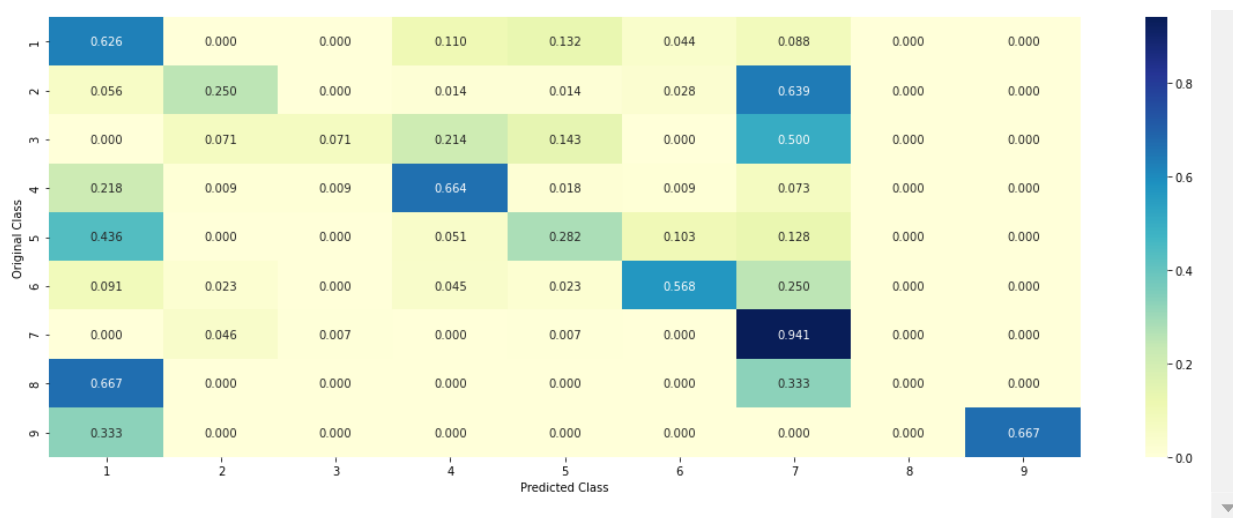
----- Confusion matrix -----



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



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



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

```
In [86]: 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])[0], 3))
print("Actual Class :", test_y[test_point_index])
indices=np.argsort(-1*abs(clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['Gender'].iloc[test_point_index])
```

```
Predicted Class : 6
Predicted Class Probabilities: [[0.0011 0.0037 0.0022 0.153 0.0015 0.8317 0.0012 0.0044 0.0013]]
Actual Class : 6
```

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

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

```
In [84]: 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_one
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(clf.coef_))[predicted_cls-1][:,no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['
```

Predicted Class : 3

Predicted Class Probabilities: [[0.0068 0.0184 0.5461 0.0086 0.2782 0.0041 0.129 0.0054 0.0032]]

Actual Class : 7

-----

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

## 4.2. K Nearest Neighbour Classification

### 4.2.1. Hyper parameter tuning

In [79]:

```

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_,
    # to avoid rounding error while multiplying probabilities we use log-probabilities
    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

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

predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", 1

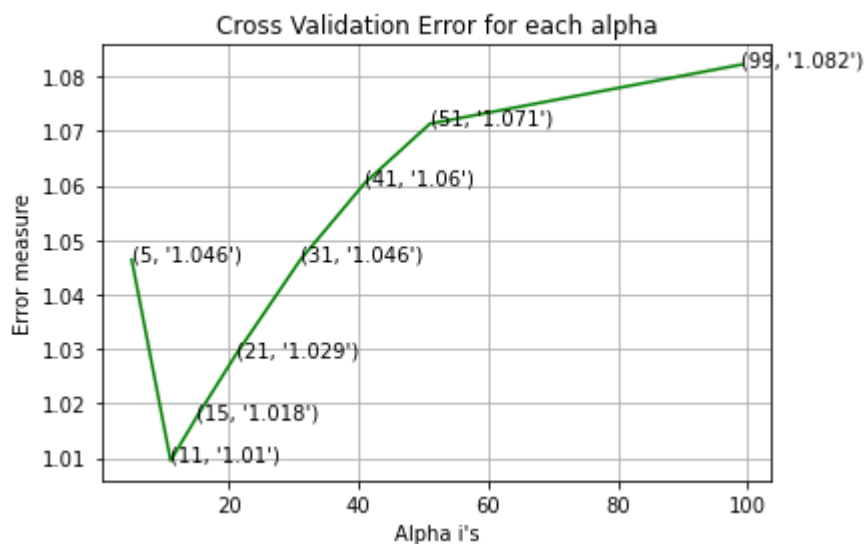
```

```

for alpha = 5
Log Loss : 1.0463826255378885
for alpha = 11
Log Loss : 1.009521670795516
for alpha = 15
Log Loss : 1.017500535784801
for alpha = 21
Log Loss : 1.0288543571290953
for alpha = 31
Log Loss : 1.046346343428731
for alpha = 41
Log Loss : 1.0604966469127457
for alpha = 51
Log Loss : 1.0713145794669074
for alpha = 99
Log Loss : 1.0821611734429417

```





For values of best alpha = 11 The train log loss is: 0.6038224111225744

For values of best alpha = 11 The cross validation log loss is: 1.009521670795516

For values of best alpha = 11 The test log loss is: 1.1057287985411093

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

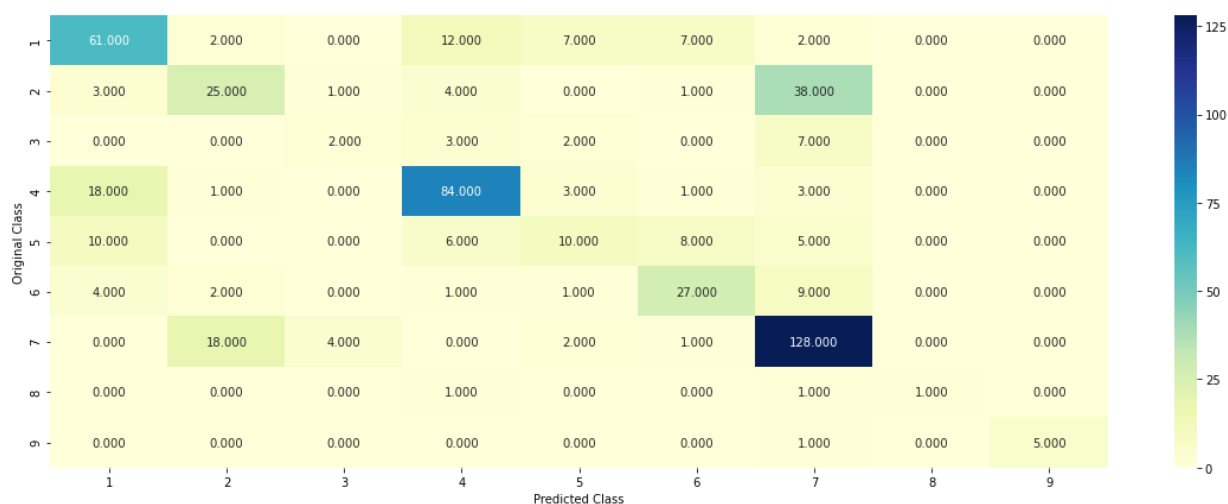
In [80]:

```
clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
predict_and_plot_confusion_matrix(train_x_responseCoding, train_y, cv_x_responseC
```

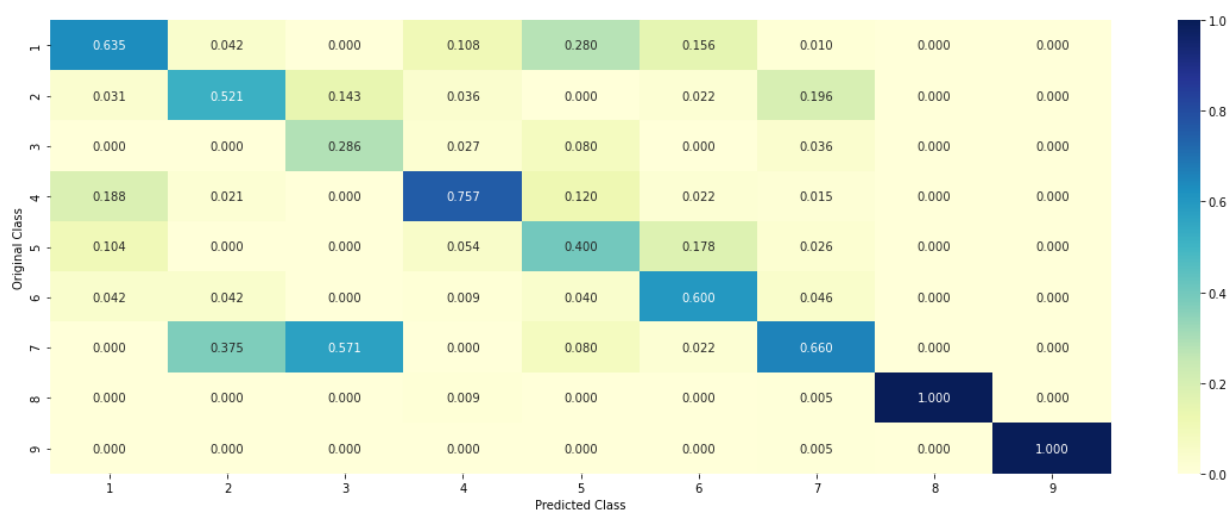
Log loss : 1.009521670795516

Number of mis-classified points : 0.35526315789473684

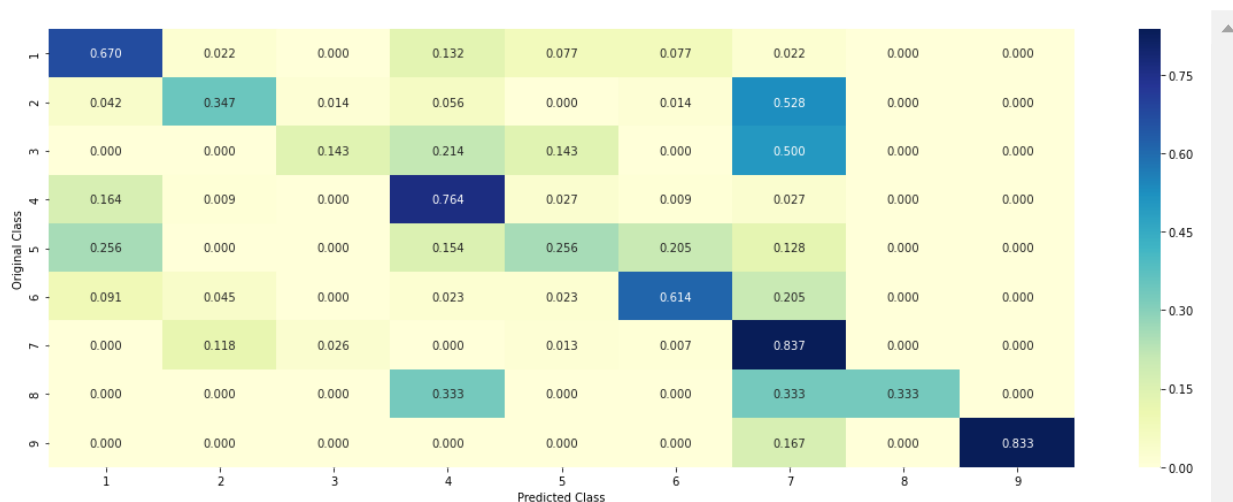
----- Confusion matrix -----



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



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



### 4.2.3. Sample Query point -1

```
In [81]: clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

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

Predicted Class : 1

Actual Class : 6

The 11 nearest neighbours of the test points belongs to classes [6 6 6 6 6 6 6 6 6 6 6]

Frequency of nearest points : Counter({6: 11})

### 4.2.4. Sample Query Point-2

```
In [82]: clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
clf.fit(train_x_responseCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_responseCoding, train_y)

test_point_index = 100

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

Predicted Class : 4

Actual Class : 4

the k value for knn is 11 and the nearest neighbours of the test points belongs to classes [4 4 4 4 4 4 4 4 4 7 4]

Frequency of nearest points : Counter({4: 10, 7: 1})

## **4.3. Logistic Regression**

### **4.3.1. With Class balancing**

#### **4.3.1.1. Hyper paramter tuning**

In [83]:

```

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')
    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_,
    # to avoid rounding error while multiplying probabilities we use log-probabilities
    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

best_alpha = np.argmin(cv_log_error_array)
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2')
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:",
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", 1

```

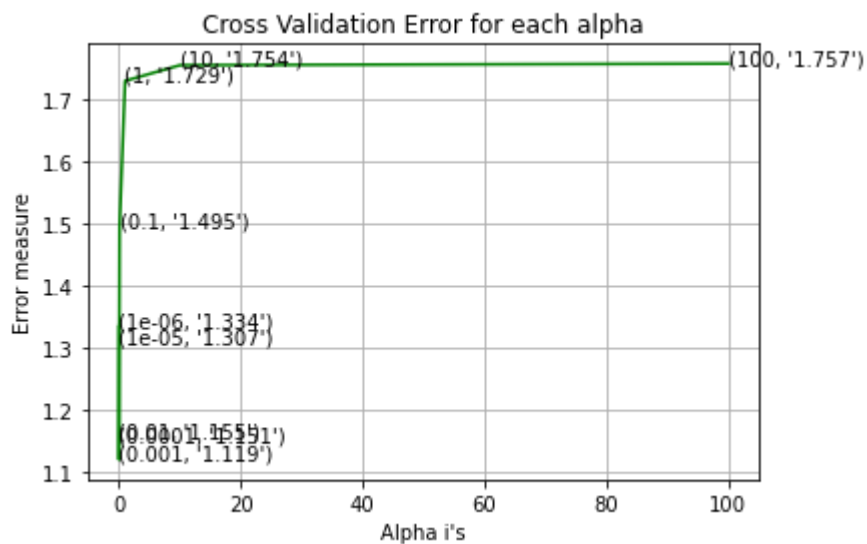
```

for alpha = 1e-06
Log Loss : 1.3335476221684441
for alpha = 1e-05
Log Loss : 1.3065902646122092
for alpha = 0.0001
Log Loss : 1.150678356021762
for alpha = 0.001
Log Loss : 1.1192692412260004
for alpha = 0.01
Log Loss : 1.1554949895981694
for alpha = 0.1
Log Loss : 1.4951787419953513
for alpha = 1
Log Loss : 1.7289437212722742
for alpha = 10
Log Loss : 1.7542144930759165

```

for alpha = 100

Log Loss : 1.7567948772379163



For values of best alpha = 0.001 The train log loss is: 0.49521010275728744

For values of best alpha = 0.001 The cross validation log loss is: 1.1192692412260004

For values of best alpha = 0.001 The test log loss is: 1.0926946004956906

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

In [87]:

```
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2')
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding)
```

Log loss : 1.1192692412260004

Number of mis-classified points : 0.3082706766917293

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



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



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



#### 4.3.1.3. Feature Importance

```
In [88]: def get_imp_feature_names(text, indices, removed_ind = []):
word_present = 0
tabulte_list = []
incresingorder_ind = 0
for i in indices:
    if i < train_gene_feature_onehotCoding.shape[1]:
        tabulte_list.append([incresingorder_ind, "Gene", "Yes"])
    elif i < 18:
        tabulte_list.append([incresingorder_ind, "Variation", "Yes"])
    if ((i > 17) & (i not in removed_ind)) :
        word = train_text_features[i]
        yes_no = True if word in text.split() else False
        if yes_no:
            word_present += 1
            tabulte_list.append([incresingorder_ind, train_text_features[i], yes_no])
        incresingorder_ind += 1
print(word_present, "most important features are present in our query point")
print("-"*50)
print("The features that are most important of the ", predicted_cls[0], " class")
print(tabulate(tabulte_list, headers=["Index", "Feature name", "Present or No"])
```

##### 4.3.1.3.1. Correctly Classified point



```
In [89]: # from tabulate import tabulate
clf = SGDClassifier(class_weight='balanced', alpha=alpha[best_alpha], penalty='l2')
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(-1*abs(clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['TEXT'].iloc[test_point_index])
```

Predicted Class : 6

Predicted Class Probabilities: [[0.0011 0.0037 0.0022 0.153 0.0015 0.8317 0.0012 0.0044 0.0013]]

Actual Class : 6

```
-----
112 Text feature [simplex] present in test data point [True]
232 Text feature [encoding] present in test data point [True]
316 Text feature [2b] present in test data point [True]
323 Text feature [hospitals] present in test data point [True]
328 Text feature [3b] present in test data point [True]
336 Text feature [author] present in test data point [True]
387 Text feature [3a] present in test data point [True]
400 Text feature [nih] present in test data point [True]
437 Text feature [mutants] present in test data point [True]
438 Text feature [2c] present in test data point [True]
466 Text feature [previously] present in test data point [True]
469 Text feature [constructs] present in test data point [True]
475 Text feature [weakened] present in test data point [True]
483 Text feature [tagged] present in test data point [True]
487 Text feature [ccdc98] present in test data point [True]
Out of the top 500 features 15 are present in query point
```

#### 4.3.1.3.2. Incorrectly Classified point

```
In [91]: 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_one
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(clf.coef_))[predicted_cls-1][:,no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['
```

Predicted Class : 6

Predicted Class Probabilities: [[0.0184 0.0052 0.0016 0.0051 0.4195 0.5409 0.0015 0.0066 0.0012]]

Actual Class : 5

```
-----
182 Text feature [v1804d] present in test data point [True]
400 Text feature [nih] present in test data point [True]
424 Text feature [assays] present in test data point [True]
466 Text feature [previously] present in test data point [True]
468 Text feature [e2663v] present in test data point [True]
474 Text feature [similarly] present in test data point [True]
476 Text feature [res] present in test data point [True]
489 Text feature [mayo] present in test data point [True]
493 Text feature [analyzed] present in test data point [True]
Out of the top 500 features 9 are present in query point
```

## 4.3.2. Without Class balancing

### 4.3.2.1. Hyper paramter tuning

In [0]:

```

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

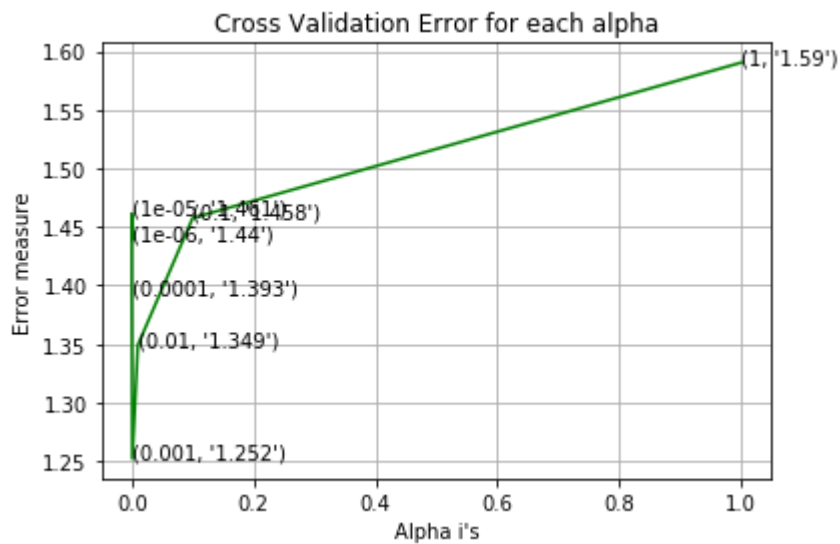
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(test_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", 1

```

```

for alpha = 1e-06
Log Loss : 1.4395190222240433
for alpha = 1e-05
Log Loss : 1.4613951945118617
for alpha = 0.0001
Log Loss : 1.392640595913179
for alpha = 0.001
Log Loss : 1.2521811628755943
for alpha = 0.01
Log Loss : 1.349151219922669
for alpha = 0.1
Log Loss : 1.457591708320943
for alpha = 1
Log Loss : 1.5902258764770603

```



For values of best alpha = 0.001 The train log loss is: 0.6257422677412771  
 For values of best alpha = 0.001 The cross validation log loss is: 1.2521811628755943  
 For values of best alpha = 0.001 The test log loss is: 1.1306020069615057

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

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

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

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

```
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=None)
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding, cv_y)
```

Log loss : 1.2521811628755943

Number of mis-classified points : 0.37218045112781956

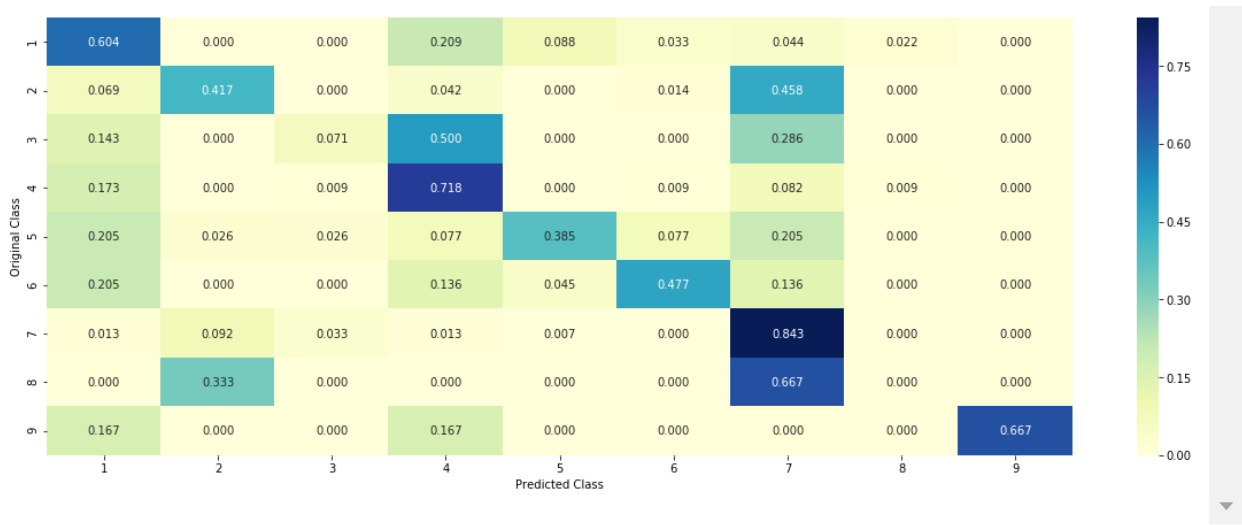
----- Confusion matrix -----



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



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



4.3.2.3. Feature Importance, Correctly Classified point

```
In [0]: clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=0)
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]), 5))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['TEXT'].iloc[test_point_index])
```

Predicted Class : 7

Predicted Class Probabilities: [[5.100e-03 1.255e-01 2.000e-04 1.300e-03 2.300e-03 1.400e-03 8.556e-01 8.500e-03 1.000e-04]]

Actual Class : 7

```
-----
60 Text feature [constitutively] present in test data point [True]
107 Text feature [flt1] present in test data point [True]
124 Text feature [cysteine] present in test data point [True]
157 Text feature [oncogenes] present in test data point [True]
158 Text feature [inhibited] present in test data point [True]
195 Text feature [activating] present in test data point [True]
200 Text feature [ligand] present in test data point [True]
203 Text feature [oncogene] present in test data point [True]
204 Text feature [technology] present in test data point [True]
257 Text feature [gaiix] present in test data point [True]
260 Text feature [concentrations] present in test data point [True]
265 Text feature [downstream] present in test data point [True]
314 Text feature [hki] present in test data point [True]
316 Text feature [dramatic] present in test data point [True]
323 Text feature [expressing] present in test data point [True]
371 Text feature [cdnas] present in test data point [True]
380 Text feature [viability] present in test data point [True]
412 Text feature [thyroid] present in test data point [True]
459 Text feature [activation] present in test data point [True]
461 Text feature [manageable] present in test data point [True]
462 Text feature [ser473] present in test data point [True]
468 Text feature [axilla] present in test data point [True]
495 Text feature [extracellular] present in test data point [True]
Out of the top 500 features 23 are present in query point
```

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

```
In [0]: 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(-1*abs(clf.coef_))[predicted_cls-1][:,no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['TEXT'].iloc[test_point_index])
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0485 0.1851 0.0052 0.0442 0.0617 0.0143 0.6317 0.0072 0.0022]]

Actual Class : 7

```
-----
60 Text feature [constitutively] present in test data point [True]
89 Text feature [constitutive] present in test data point [True]
116 Text feature [activated] present in test data point [True]
158 Text feature [inhibited] present in test data point [True]
159 Text feature [transforming] present in test data point [True]
193 Text feature [receptors] present in test data point [True]
195 Text feature [activating] present in test data point [True]
203 Text feature [oncogene] present in test data point [True]
226 Text feature [transform] present in test data point [True]
241 Text feature [isozyme] present in test data point [True]
265 Text feature [downstream] present in test data point [True]
377 Text feature [agar] present in test data point [True]
442 Text feature [interatomic] present in test data point [True]
459 Text feature [activation] present in test data point [True]
Out of the top 500 features 14 are present in query point
```

## 4.4. Linear Support Vector Machines

### 4.4.1. Hyper paramter tuning



In [0]:

```

alpha = [10 ** x for x in range(-5, 3)]
cv_log_error_array = []
for i in alpha:
    print("for C =", i)
    # clf = SVC(C=i, kernel='linear', probability=True, class_weight='balanced')
    clf = SGDClassifier(class_weight='balanced', alpha=i, penalty='l2', loss='hinge')
    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_))
    print("Log Loss :", log_loss(cv_y, sig_clf_probs))

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

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

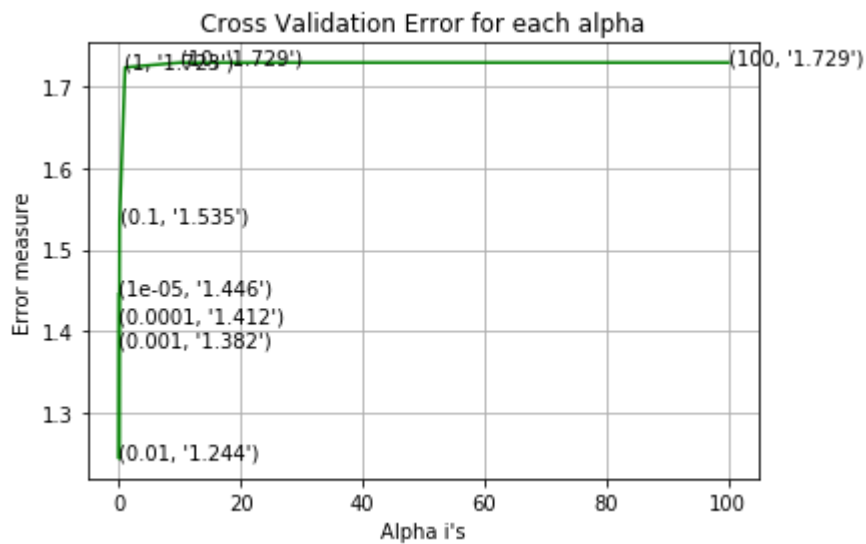
predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
      log_loss(train_y, predict_y, labels=clf.classes_))
predict_y = sig_clf.predict_proba(cv_x_onehotCoding)
print('For values of best alpha = ', alpha[best_alpha],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

```

```

for C = 1e-05
Log Loss : 1.4456349250609233
for C = 0.0001
Log Loss : 1.4117883301099556
for C = 0.001
Log Loss : 1.3818342037841624
for C = 0.01
Log Loss : 1.2442964974823838
for C = 0.1
Log Loss : 1.5346828298587332
for C = 1
Log Loss : 1.722800653929441
for C = 10
Log Loss : 1.7286360420759161
for C = 100
Log Loss : 1.7286184454094997

```



For values of best alpha = 0.01 The train log loss is: 0.7628309867716067

For values of best alpha = 0.01 The cross validation log loss is: 1.2442964974823838

For values of best alpha = 0.01 The test log loss is: 1.1541891969863685

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

```
In [0]: # read more about support vector machines with linear kernal here http://scikit-
# -----
# default parameters
# SVC(C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, proba
# cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function

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

```
# clf = SVC(C=alpha[best_alpha],kernel='linear',probability=True, class_weight='b
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_s
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y,cv_x_onehotCoding
```

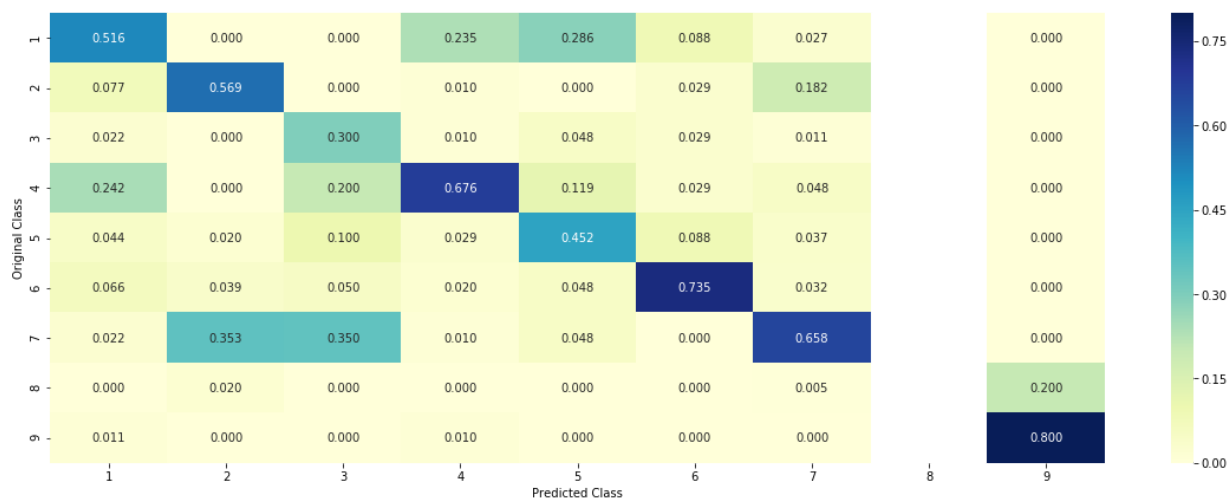
Log loss : 1.2442964974823838

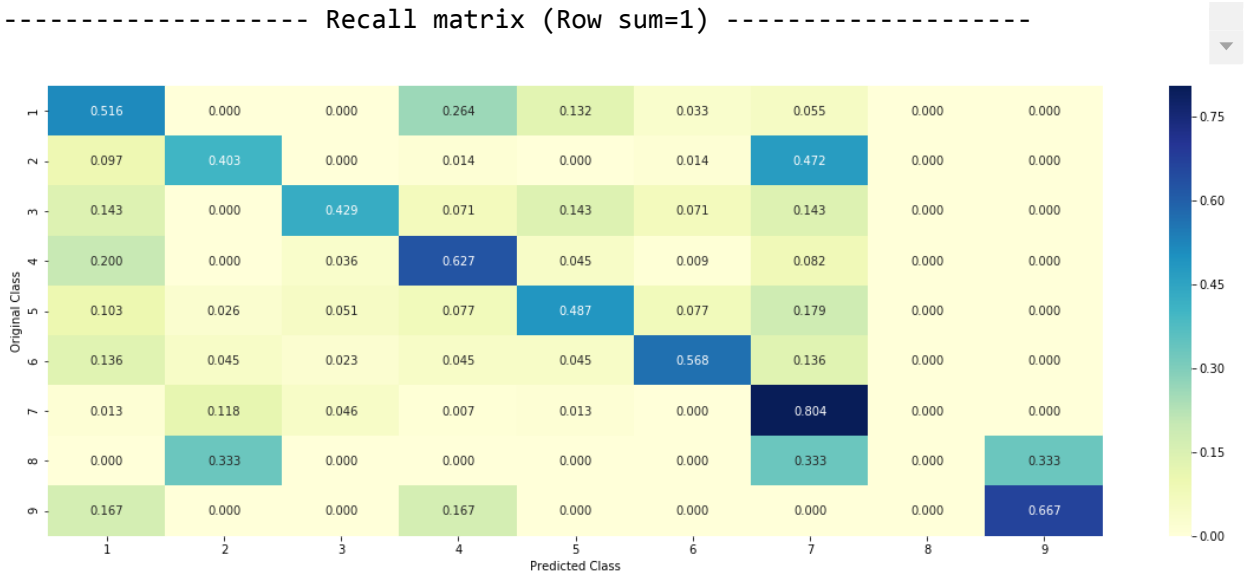
Number of mis-classified points : 0.39473684210526316

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



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





4.3.3. Feature Importance

4.3.3.1. For Correctly classified point

```
In [0]: clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='hinge', random_s
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_one
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-1*abs(clf.coef_))[predicted_cls-1][:,:no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index],test_df['
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0153 0.1199 0.0029 0.0151 0.0121 0.0075 0.8104 0.0129 0.0039]]

Actual Class : 7

```
-----
28 Text feature [constitutively] present in test data point [True]
29 Text feature [cysteine] present in test data point [True]
49 Text feature [cdnas] present in test data point [True]
76 Text feature [flt1] present in test data point [True]
79 Text feature [concentrations] present in test data point [True]
82 Text feature [gaiix] present in test data point [True]
96 Text feature [technology] present in test data point [True]
101 Text feature [inhibited] present in test data point [True]
104 Text feature [activating] present in test data point [True]
114 Text feature [oncogenes] present in test data point [True]
147 Text feature [expressing] present in test data point [True]
150 Text feature [mapk] present in test data point [True]
151 Text feature [oncogene] present in test data point [True]
169 Text feature [thyroid] present in test data point [True]
171 Text feature [inhibitor] present in test data point [True]
205 Text feature [transduced] present in test data point [True]
211 Text feature [seeded] present in test data point [True]
230 Text feature [ligand] present in test data point [True]
255 Text feature [activation] present in test data point [True]
279 Text feature [downstream] present in test data point [True]
314 Text feature [doses] present in test data point [True]
351 Text feature [subcutaneous] present in test data point [True]
366 Text feature [atcc] present in test data point [True]
405 Text feature [melanocyte] present in test data point [True]
436 Text feature [hours] present in test data point [True]
445 Text feature [selleck] present in test data point [True]
446 Text feature [dramatic] present in test data point [True]
454 Text feature [chemiluminescence] present in test data point [True]
487 Text feature [viability] present in test data point [True]
489 Text feature [ser473] present in test data point [True]
Out of the top 500 features 30 are present in query point
```

#### 4.3.3.2. For Incorrectly classified point

```
In [0]: 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(-1*abs(clf.coef_))[predicted_cls-1][:, :no_feature]
print("-"*50)
get_impfeature_names(indices[0], test_df['TEXT'].iloc[test_point_index], test_df['TEXT'].iloc[test_point_index])
```

Predicted Class : 7

Predicted Class Probabilities: [[0.0786 0.1516 0.0146 0.1064 0.1105 0.0323 0.4839 0.0128 0.0094]]

Actual Class : 7

```
-----
28 Text feature [constitutively] present in test data point [True]
40 Text feature [constitutive] present in test data point [True]
73 Text feature [activated] present in test data point [True]
75 Text feature [transforming] present in test data point [True]
94 Text feature [receptors] present in test data point [True]
97 Text feature [exchange] present in test data point [True]
101 Text feature [inhibited] present in test data point [True]
104 Text feature [activating] present in test data point [True]
151 Text feature [oncogene] present in test data point [True]
231 Text feature [transform] present in test data point [True]
255 Text feature [activation] present in test data point [True]
279 Text feature [downstream] present in test data point [True]
440 Text feature [doubled] present in test data point [True]
470 Text feature [substituting] present in test data point [True]
Out of the top 500 features 14 are present in query point
```

## 4.5 Random Forest Classifier

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

In [0]:

```

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)
        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_))
        print("Log Loss :",log_loss(cv_y, sig_clf_probs))

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

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

predict_y = sig_clf.predict_proba(train_x_onehotCoding)
print('For values of best estimator = ', alpha[int(best_alpha/2)], "The train log loss is: ", log_loss(train_y, predict_y))
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(cv_y, predict_y))
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(test_y, predict_y))

```

```

for n_estimators = 100 and max depth = 5
Log Loss : 1.2572535683354957
for n_estimators = 100 and max depth = 10
Log Loss : 1.1868414223711878
for n_estimators = 200 and max depth = 5
Log Loss : 1.2378734502517341
for n_estimators = 200 and max depth = 10
Log Loss : 1.1811031780258958
for n_estimators = 500 and max depth = 5
Log Loss : 1.2368241894319212
for n_estimators = 500 and max depth = 10
Log Loss : 1.176754594516683
for n_estimators = 1000 and max depth = 5
Log Loss : 1.2357829533963691

```

```
for n_estimators = 1000 and max depth = 10
Log Loss : 1.174993079576866
for n_estimators = 2000 and max depth = 5
Log Loss : 1.236042392554891
for n_estimators = 2000 and max depth = 10
Log Loss : 1.1759745074379755
For values of best estimator = 1000 The train log loss is: 0.709539673208275
2
For values of best estimator = 1000 The cross validation log loss is: 1.1749
93079576866
For values of best estimator = 1000 The test log loss is: 1.1630923149103904
```

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



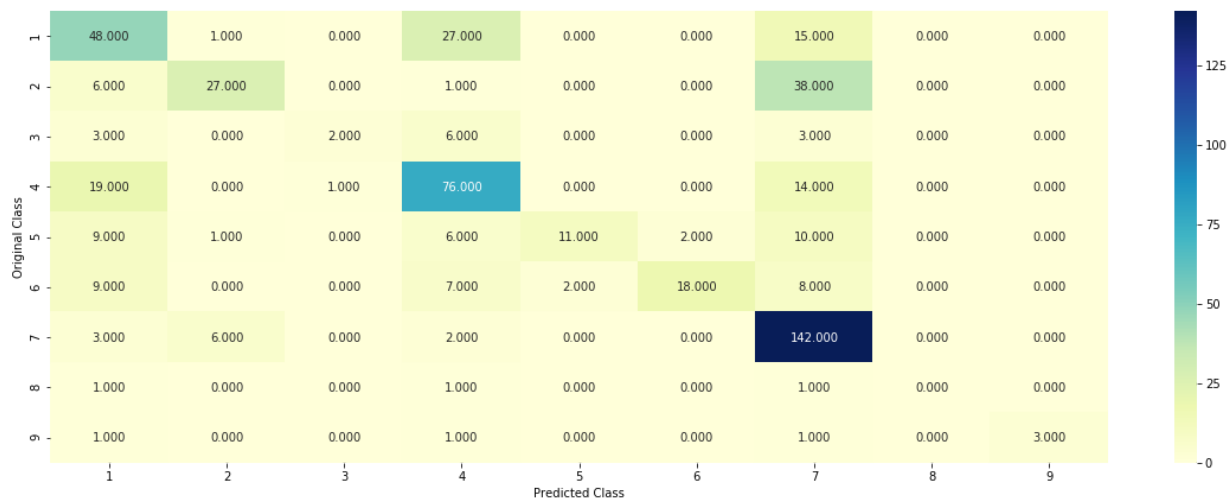
In [0]:

```
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini')
predict_and_plot_confusion_matrix(train_x_onehotCoding, train_y, cv_x_onehotCoding)
```

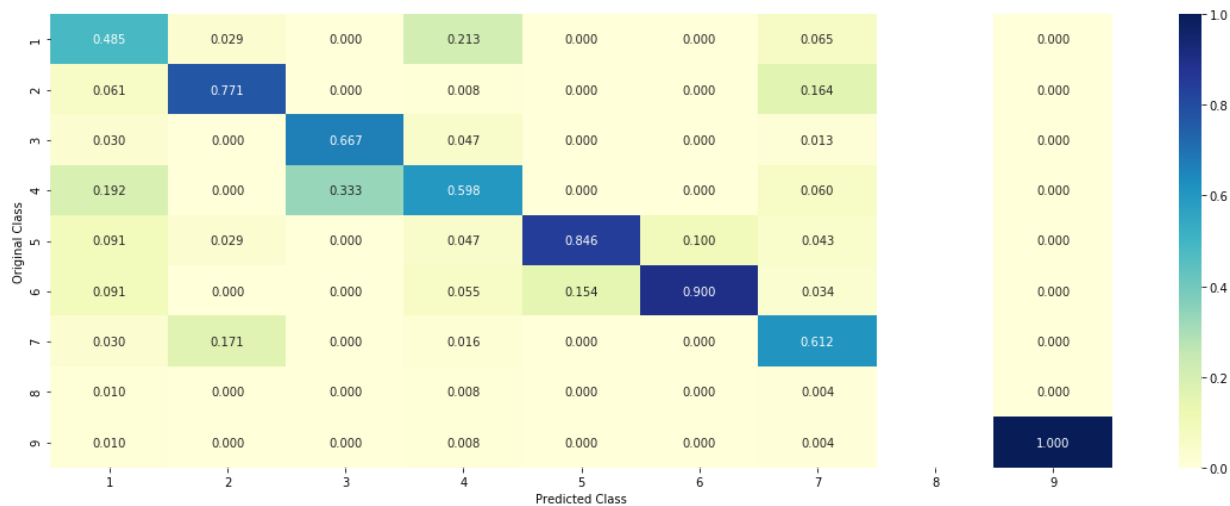
Log loss : 1.174993079576866

Number of mis-classified points : 0.38533834586466165

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



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



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



### 4.5.3. Feature Importance

#### 4.5.3.1. Correctly Classified point

```

In [0]: # test_point_index = 10
clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini')
clf.fit(train_x_onehotCoding, train_y)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(train_x_onehotCoding, train_y)

test_point_index = 1
no_feature = 100
predicted_cls = sig_clf.predict(test_x_onehotCoding[test_point_index])
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_onehotCoding[test_point_index]), 2))
print("Actual Class :", test_y[test_point_index])
indices = np.argsort(-clf.feature_importances_)
print("-"*50)
get_impfeature_names(indices[:no_feature], test_df['TEXT'].iloc[test_point_index])

```

Predicted Class : 7

Predicted Class Probabilities: [[0.0454 0.1404 0.0133 0.029 0.036 0.0294 0.6977 0.005 0.004 ]]

Actual Class : 7

```

-----
0 Text feature [inhibitors] present in test data point [True]
1 Text feature [kinase] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [tyrosine] present in test data point [True]
4 Text feature [missense] present in test data point [True]
5 Text feature [inhibitor] present in test data point [True]
7 Text feature [treatment] present in test data point [True]
8 Text feature [oncogenic] present in test data point [True]
9 Text feature [suppressor] present in test data point [True]
10 Text feature [activation] present in test data point [True]
11 Text feature [phosphorylation] present in test data point [True]
12 Text feature [kinases] present in test data point [True]
13 Text feature [nonsense] present in test data point [True]
14 Text feature [akt] present in test data point [True]
15 Text feature [function] present in test data point [True]
17 Text feature [erk] present in test data point [True]
19 Text feature [growth] present in test data point [True]
20 Text feature [variants] present in test data point [True]
22 Text feature [frameshift] present in test data point [True]
24 Text feature [therapeutic] present in test data point [True]
25 Text feature [functional] present in test data point [True]
28 Text feature [signaling] present in test data point [True]
30 Text feature [patients] present in test data point [True]
31 Text feature [cells] present in test data point [True]
32 Text feature [constitutively] present in test data point [True]
34 Text feature [trials] present in test data point [True]
35 Text feature [therapy] present in test data point [True]
37 Text feature [erk1] present in test data point [True]
38 Text feature [activate] present in test data point [True]
39 Text feature [downstream] present in test data point [True]
41 Text feature [efficacy] present in test data point [True]
42 Text feature [protein] present in test data point [True]
43 Text feature [loss] present in test data point [True]
44 Text feature [inhibited] present in test data point [True]
45 Text feature [expressing] present in test data point [True]

```

```
46 Text feature [pten] present in test data point [True]
48 Text feature [lines] present in test data point [True]
49 Text feature [treated] present in test data point [True]
50 Text feature [proliferation] present in test data point [True]
51 Text feature [drug] present in test data point [True]
57 Text feature [mek] present in test data point [True]
59 Text feature [inhibition] present in test data point [True]
61 Text feature [repair] present in test data point [True]
62 Text feature [sensitivity] present in test data point [True]
64 Text feature [receptor] present in test data point [True]
66 Text feature [assays] present in test data point [True]
68 Text feature [survival] present in test data point [True]
69 Text feature [cell] present in test data point [True]
71 Text feature [ligand] present in test data point [True]
73 Text feature [expression] present in test data point [True]
74 Text feature [variant] present in test data point [True]
75 Text feature [oncogene] present in test data point [True]
78 Text feature [extracellular] present in test data point [True]
79 Text feature [doses] present in test data point [True]
80 Text feature [mapk] present in test data point [True]
81 Text feature [hours] present in test data point [True]
84 Text feature [information] present in test data point [True]
86 Text feature [harboring] present in test data point [True]
90 Text feature [dna] present in test data point [True]
91 Text feature [concentrations] present in test data point [True]
92 Text feature [likelihood] present in test data point [True]
93 Text feature [months] present in test data point [True]
94 Text feature [binding] present in test data point [True]
96 Text feature [imatinib] present in test data point [True]
98 Text feature [preclinical] present in test data point [True]
Out of the top 100 features 65 are present in query point
```

#### 4.5.3.2. Inorrectly Classified point

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

Predicted Class : 7

Predicted Class Probabilities: [[0.1337 0.116 0.0224 0.1773 0.0674 0.0545 0.4156 0.0071 0.0059]]

Actual Class : 7

```
-----
0 Text feature [inhibitors] present in test data point [True]
1 Text feature [kinase] present in test data point [True]
2 Text feature [activating] present in test data point [True]
3 Text feature [tyrosine] present in test data point [True]
6 Text feature [activated] present in test data point [True]
8 Text feature [oncogenic] present in test data point [True]
10 Text feature [activation] present in test data point [True]
11 Text feature [phosphorylation] present in test data point [True]
12 Text feature [kinases] present in test data point [True]
14 Text feature [akt] present in test data point [True]
15 Text feature [function] present in test data point [True]
19 Text feature [growth] present in test data point [True]
21 Text feature [constitutive] present in test data point [True]
25 Text feature [functional] present in test data point [True]
28 Text feature [signaling] present in test data point [True]
31 Text feature [cells] present in test data point [True]
32 Text feature [constitutively] present in test data point [True]
38 Text feature [activate] present in test data point [True]
39 Text feature [downstream] present in test data point [True]
42 Text feature [protein] present in test data point [True]
43 Text feature [loss] present in test data point [True]
44 Text feature [inhibited] present in test data point [True]
46 Text feature [pten] present in test data point [True]
47 Text feature [transforming] present in test data point [True]
48 Text feature [lines] present in test data point [True]
50 Text feature [proliferation] present in test data point [True]
53 Text feature [neutral] present in test data point [True]
55 Text feature [transform] present in test data point [True]
56 Text feature [stability] present in test data point [True]
58 Text feature [transformation] present in test data point [True]
59 Text feature [inhibition] present in test data point [True]
62 Text feature [sensitivity] present in test data point [True]
64 Text feature [receptor] present in test data point [True]
66 Text feature [assays] present in test data point [True]
69 Text feature [cell] present in test data point [True]
75 Text feature [oncogene] present in test data point [True]
84 Text feature [information] present in test data point [True]
90 Text feature [dna] present in test data point [True]
94 Text feature [binding] present in test data point [True]
Out of the top 100 features 39 are present in query point
```

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

In [0]:

```

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)
        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_))
        print("Log Loss :",log_loss(cv_y, sig_clf_probs))
    ...

fig, ax = plt.subplots()
features = np.dot(np.array(alpha)[: ,None],np.array(max_depth)[None]).ravel()
ax.plot(features, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[int(i/4)],max_depth[int(i%4)],str(txt)), (features[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()
'''

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

predict_y = sig_clf.predict_proba(train_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The train log loss is: ", log_loss(train_y, predict_y))
predict_y = sig_clf.predict_proba(cv_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The cross validation log loss is: ", log_loss(cv_y, predict_y))
predict_y = sig_clf.predict_proba(test_x_responseCoding)
print('For values of best alpha = ', alpha[int(best_alpha/4)], "The test log loss is: ", log_loss(test_y, predict_y))

```

```

for n_estimators = 10 and max depth = 2
Log Loss : 2.2657048897349608
for n_estimators = 10 and max depth = 3
Log Loss : 1.7459205010556096
for n_estimators = 10 and max depth = 5
Log Loss : 1.4368353925512503
for n_estimators = 10 and max depth = 10
Log Loss : 1.904597809032912
for n_estimators = 50 and max depth = 2
Log Loss : 1.7221951095007484
for n_estimators = 50 and max depth = 3
Log Loss : 1.4984825877845531
for n_estimators = 50 and max depth = 5
Log Loss : 1.4593628982873716
for n_estimators = 50 and max depth = 10
Log Loss : 1.4593628982873716

```

```
Log Loss : 1.8434939703555409
for n_estimators = 100 and max depth = 2
Log Loss : 1.6182209245331227
for n_estimators = 100 and max depth = 3
Log Loss : 1.5199297988828253
for n_estimators = 100 and max depth = 5
Log Loss : 1.4177501184246677
for n_estimators = 100 and max depth = 10
Log Loss : 1.8227504417195126
for n_estimators = 200 and max depth = 2
Log Loss : 1.6622571648074496
for n_estimators = 200 and max depth = 3
Log Loss : 1.4800771339141767
for n_estimators = 200 and max depth = 5
Log Loss : 1.4412060242341358
for n_estimators = 200 and max depth = 10
Log Loss : 1.7892406351442258
for n_estimators = 500 and max depth = 2
Log Loss : 1.715950314170445
for n_estimators = 500 and max depth = 3
Log Loss : 1.5658682738699774
for n_estimators = 500 and max depth = 5
Log Loss : 1.4445360301518217
for n_estimators = 500 and max depth = 10
Log Loss : 1.8421097596928397
for n_estimators = 1000 and max depth = 2
Log Loss : 1.6834927870864949
for n_estimators = 1000 and max depth = 3
Log Loss : 1.5631973035931377
for n_estimators = 1000 and max depth = 5
Log Loss : 1.4449980792724129
for n_estimators = 1000 and max depth = 10
Log Loss : 1.85233132619749
For values of best alpha = 100 The train log loss is: 0.060702709444608406
For values of best alpha = 100 The cross validation log loss is: 1.417750118
424668
For values of best alpha = 100 The test log loss is: 1.3806278998341923
```

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



In [0]:

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

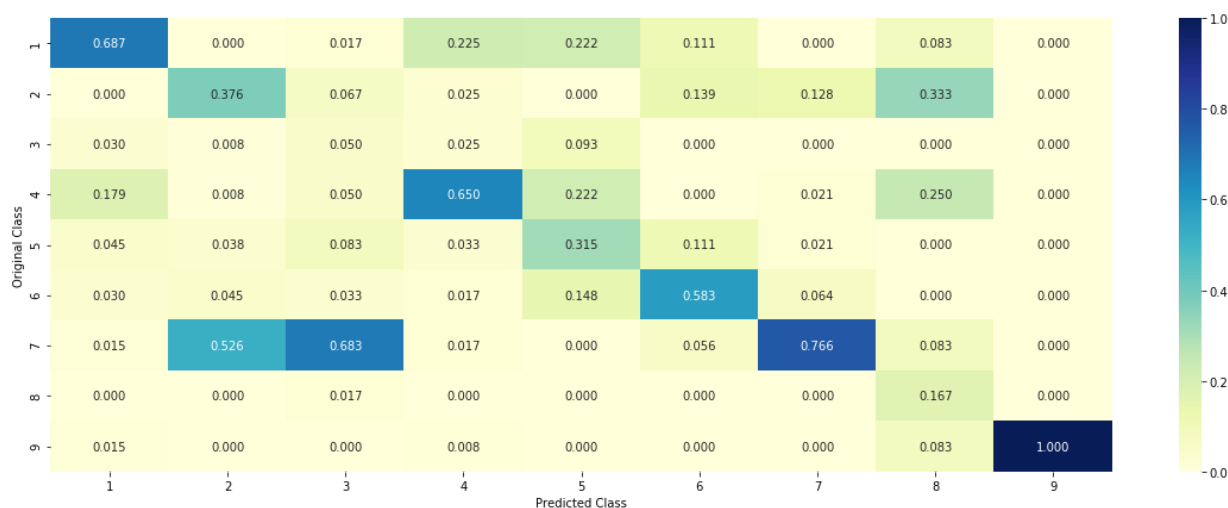
Log loss : 1.4177501184246677

Number of mis-classified points : 0.518796992481203

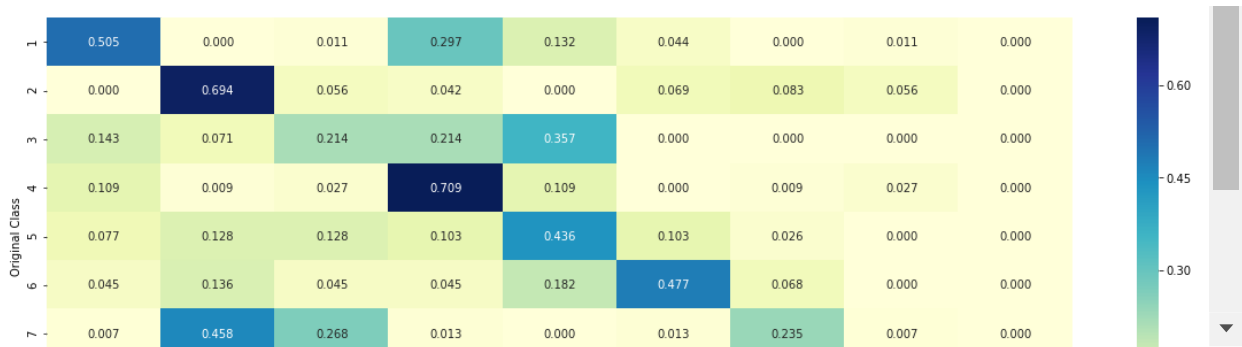
----- Confusion matrix -----



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



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



4.5.5. Feature Importance

4.5.5.1. Correctly Classified point

```

In [0]: clf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/4)], criterion='gini')
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))
print("Predicted Class :", predicted_cls[0])
print("Predicted Class Probabilities:", np.round(sig_clf.predict_proba(test_x_responseCoding[test_point_index].reshape(1)), 2))
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: [[0.0143 0.5044 0.1471 0.0191 0.0245 0.065 0.1724 0.039 0.0142]]

Actual Class : 7

```

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

```

#### 4.5.5.2. Incorrectly Classified point

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

```
Predicted Class : 7
Predicted Class Probabilities: [[0.0281 0.2006 0.203  0.0857 0.0626 0.0906 0.22
49 0.0676 0.0369]]
Actual Class : 7
```

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

## 4.7 Stack the models

### 4.7.1 testing with hyper parameter tuning

In [0]:

```

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

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

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

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

```

```

Logistic Regression : Log Loss: 1.24
Support vector machines : Log Loss: 1.72
Naive Bayes : Log Loss: 1.37

```

```

-----
Stacking Classifier : for the value of alpha: 0.000100 Log Loss: 2.179
Stacking Classifier : for the value of alpha: 0.001000 Log Loss: 2.049
Stacking Classifier : for the value of alpha: 0.010000 Log Loss: 1.577
Stacking Classifier : for the value of alpha: 0.100000 Log Loss: 1.224
Stacking Classifier : for the value of alpha: 1.000000 Log Loss: 1.366
Stacking Classifier : for the value of alpha: 10.000000 Log Loss: 1.690

```

## 4.7.2 testing the model with the best hyper parameters

```

In [0]: lr = LogisticRegression(C=0.1)
sclf = StackingClassifier(classifiers=[sig_clf1, sig_clf2, sig_clf3], meta_classifier=lr)
sclf.fit(train_x_onehotCoding, train_y)

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

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

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

print("Number of missclassified point :", np.count_nonzero((sclf.predict(test_x_onehotCoding) != test_y)))
plot_confusion_matrix(test_y=test_y, predict_y=sclf.predict(test_x_onehotCoding))

```

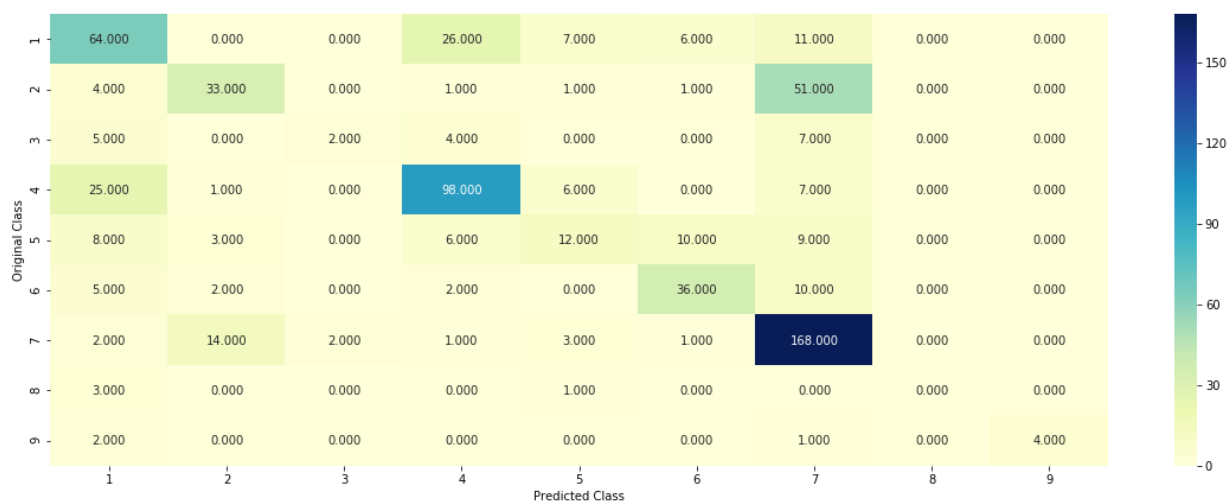
Log loss (train) on the stacking classifier : 0.6760284396805781

Log loss (CV) on the stacking classifier : 1.2243084610674686

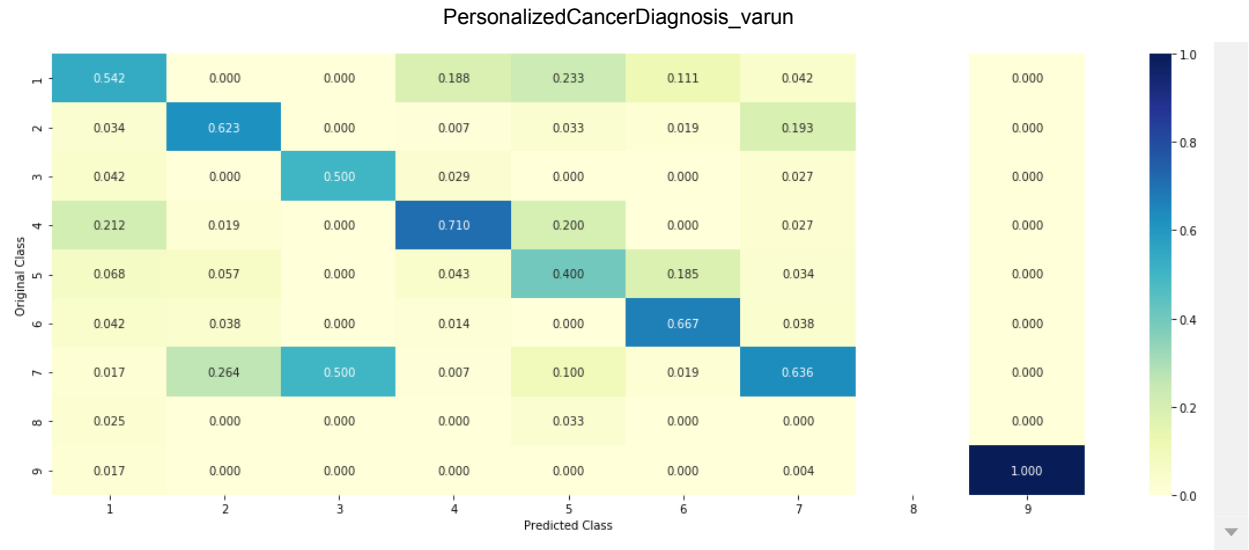
Log loss (test) on the stacking classifier : 1.1562525475350196

Number of missclassified point : 0.37293233082706767

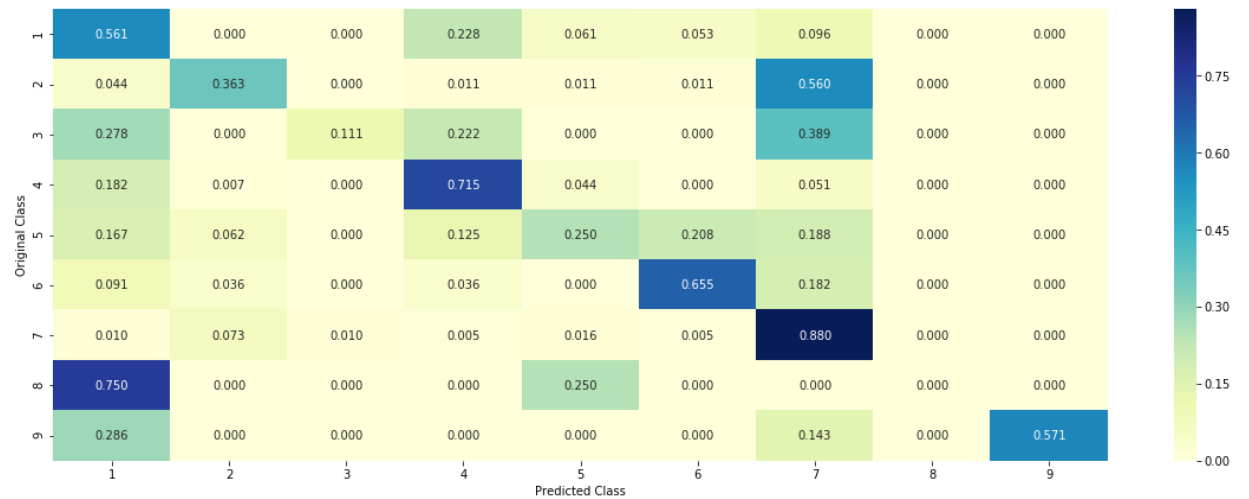
----- Confusion matrix -----



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



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



### 4.7.3 Maximum Voting classifier

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

Log loss (train) on the VotingClassifier : 0.9407598679043604

Log loss (CV) on the VotingClassifier : 1.2835402100341697

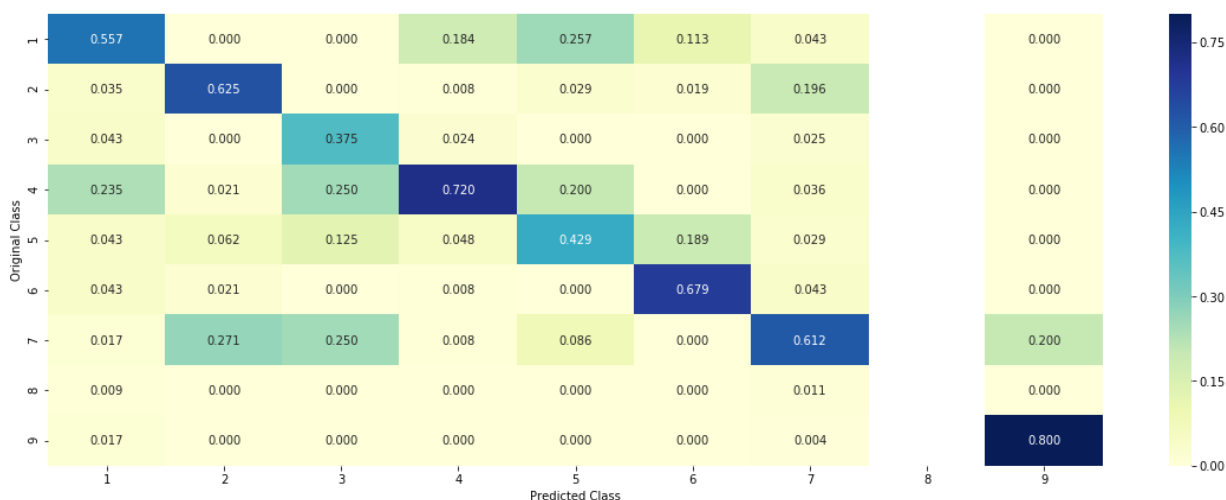
Log loss (test) on the VotingClassifier : 1.223278167176945

Number of missclassified point : 0.3819548872180451

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



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



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



