# Microsoft Malware detection

# 1.Business/Real-world Problem

## 1.1. What is Malware?

The term malware is a contraction of malicious software. Put simply, malware is any piece of software that was written with the intent of doing harm to data, devices or to people.
Source: https://www.avg.com/en/signal/what-is-malware

## 1.2. Problem Statement

In the past few years, the malware industry has grown very rapidly that, the syndicates invest heavily in technologies to evade traditional protection, forcing the anti-malware groups/communities to build more robust softwares to detect and terminate these attacks. The major part of protecting a computer system from a malware attack is to **identify whether a given piece of file/software is a malware.**

## 1.3 Source/Useful Links

Microsoft has been very active in building anti-malware products over the years and it runs it's anti-malware utilities over **150 million computers** around the world. This generates tens of millions of daily data points to be analyzed as potential malware. In order to be effective in analyzing and classifying such large amounts of data, we need to be able to group them into groups and identify their respective families.

This dataset provided by Microsoft contains about 9 classes of malware. ,

**Source:** https://www.kaggle.com/c/malware-classification

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

1. Minimize multi-class error.
2. Multi-class probability estimates.
3. Malware detection should not take hours and block the user's computer. It should fininsh in a few seconds or a minute.

# 2. Machine Learning Problem

## 2.1. Data

### 2.1.1. Data Overview

- Source : https://www.kaggle.com/c/malware-classification/data
- For every malware, we have two files
  1. .asm file (read more: https://www.reviversoft.com/file-extensions/asm)
  2. .bytes file (the raw data contains the hexadecimal representation of the file's binary content, without the PE header)

- Total train dataset consist of 200GB data out of which 50Gb of data is .bytes files and 150GB of data is .asm files:
- **Lots of Data for a single-box/computer.**
- There are total 10,868 .bytes files and 10,868 asm files total 21,736 files
- There are 9 types of malwares (9 classes) in our give data
- Types of Malware:
  1. Ramnit
  2. Lollipop
  3. Kelihos_ver3
  4. Vundo
  5. Simda
  6. Tracur
  7. Kelihos_ver1
  8. Obfuscator.ACY
  9. Gatak

### 2.1.2. Example Data Point

**.asm file**

```
.text:00401000                                         assume es:nothing, s
s:nothing, ds:_data,    fs:nothing, gs:nothing
.text:00401000 56                                      push    esi
.text:00401001 8D 44 24    08                          lea     eax,
[esp+8]
.text:00401005 50                                      push    eax
.text:00401006 8B F1                                   mov     esi, ecx
.text:00401008 E8 1C 1B    00 00                        call    ??
0exception@std@@QAE@ABQBD@Z ; std::exception::exception(char const * cons
t &)
.text:0040100D C7 06 08    BB 42 00                      mov     dwo
rd ptr [esi],    offset off_42BB08
.text:00401013 8B C6                                   mov     eax, esi
.text:00401015 5E                                      pop     esi
.text:00401016 C2 04 00                                retn    4
.text:00401016                                         ; ------------------------
--------------------------------------------------
.text:00401019 CC CC CC    CC CC CC CC                         align 10
h
.text:00401020 C7 01 08    BB 42 00                      mov     dwo
rd ptr [ecx],    offset off_42BB08
.text:00401026 E9 26 1C    00 00                         jmp     su
b_402C51
.text:00401026                                         ; ------------------------
--------------------------------------------------
.text:0040102B CC CC CC    CC CC                             align 10h
.text:00401030 56                                      push    esi
.text:00401031 8B F1                                   mov     esi, ecx
.text:00401033 C7 06 08    BB 42 00                      mov     dwo
rd ptr [esi],    offset off_42BB08
.text:00401039 E8 13 1C    00 00                         call    su
b_402C51
.text:0040103E F6 44 24    08 01                         test    by
te ptr    [esp+8], 1
.text:00401043 74 09                                   jz      short loc_
40104E
.text:00401045 56                                      push    esi
.text:00401046 E8 6C 1E    00 00                         call    ??
3@YAXPAX@Z    ; operator delete(void *)
.text:0040104B 83 C4 04                                add     esp, 4
.text:0040104E
.text:0040104E                                loc_40104E:
  ; CODE XREF: .text:00401043↓j
.text:0040104E 8B C6                                   mov     eax, esi
.text:00401050 5E                                      pop     esi
.text:00401051 C2 04 00                                retn    4
.text:00401051                                         ; ------------------------
--------------------------------------------------
```

**.bytes file**

```
00401000 00 00 80 40 40 28 00 1C 02 42 00 C4 00 20 04 20
00401010 00 00 20 09 2A 02 00 00 00 00 8E 10 41 0A 21 01
00401020 40 00 02 01 00 90 21 00 32 40 00 1C 01 40 C8 18
00401030 40 82 02 63 20 00 00 09 10 01 02 21 00 82 00 04
00401040 82 20 08 83 00 08 00 00 00 00 02 00 60 80 10 80
00401050 18 00 00 20 A9 00 00 00 00 04 04 78 01 02 70 90
00401060 00 02 00 08 20 12 00 00 00 40 10 00 80 00 40 19
00401070 00 00 00 00 11 20 80 04 80 10 00 20 00 00 25 00
00401080 00 00 01 00 00 04 00 10 02 C1 80 80 00 20 20 00
00401090 08 A0 01 01 44 28 00 00 08 10 20 00 02 08 00 00
004010A0 00 40 00 00 00 34 40 40 00 04 00 08 80 08 00 08
004010B0 10 00 40 00 68 02 40 04 E1 00 28 14 00 08 20 0A
004010C0 06 01 02 00 40 00 00 00 00 00 00 20 00 02 00 04
004010D0 80 18 90 00 00 10 A0 00 45 09 00 10 04 40 44 82
004010E0 90 00 26 10 00 00 04 00 82 00 00 00 20 40 00 00
004010F0 B4 00 00 40 00 02 20 25 08 00 00 00 00 00 00 00
00401100 08 00 00 50 00 08 40 50 00 02 06 22 08 85 30 00
00401110 00 80 00 80 60 00 09 00 04 20 00 00 00 00 00 00
00401120 00 82 40 02 00 11 46 01 4A 01 8C 01 E6 00 86 10
00401130 4C 01 22 00 64 00 AE 01 EA 01 2A 11 E8 10 26 11
00401140 4E 11 8E 11 C2 00 6C 00 0C 11 60 01 CA 00 62 10
00401150 6C 01 A0 11 CE 10 2C 11 4E 10 8C 00 CE 01 AE 01
00401160 6C 10 6C 11 A2 01 AE 00 46 11 EE 10 22 00 A8 00
00401170 EC 01 08 11 A2 01 AE 10 6C 00 6E 00 AC 11 8C 00
00401180 EC 01 2A 10 2A 01 AE 00 40 00 C8 10 48 01 4E 11
00401190 0E 00 EC 11 24 10 4A 10 04 01 C8 11 E6 01 C2 00
```

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

## 2.2.1. Type of Machine Learning Problem

```
        There are nine different classes of malware that we need to class
    ify a given a data point => Multi class classification problem
```

## 2.2.2. Performance Metric

Source: https://www.kaggle.com/c/malware-classification#evaluation
(https://www.kaggle.com/c/malware-classification#evaluation)

Metric(s):

- Multi class log-loss
- Confusion matrix

### 2.2.3. Machine Learing Objectives and Constraints

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

Constraints:
* Class probabilities are needed. * Penalize the errors in class probabilites => Metric is Log-loss. * Some Latency constraints.

## 2.3. Train and Test Dataset

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

## 2.4. Useful blogs, videos and reference papers

http://blog.kaggle.com/2015/05/26/microsoft-malware-winners-interview-1st-place-no-to-overfitting/
https://arxiv.org/pdf/1511.04317.pdf
First place solution in Kaggle competition: https://www.youtube.com/watch?v=VLQTRlLGz5Y
https://github.com/dchad/malware-detection
http://vizsec.org/files/2011/Nataraj.pdf
https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu_pIB6ua?dl=0
" Cross validation is more trustworthy than domain knowledge."

# 3. Exploratory Data Analysis

```python
import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use(u'nbAgg')
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process# this is used for multithreading
import multiprocessing
import codecs# this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```
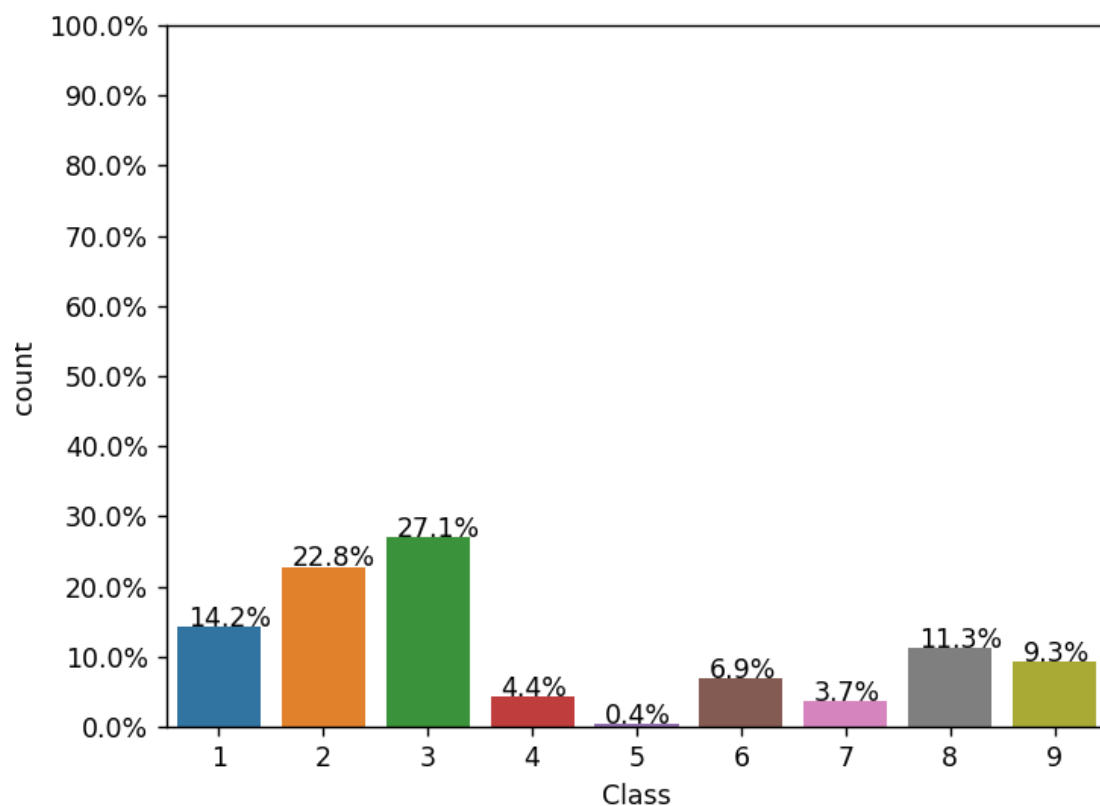
```
In [10]:   1  #separating byte files and asm files
           2
           3  source = 'train'
           4  destination = 'byteFiles'
           5
           6  # we will check if the folder 'byteFiles' exists if it not there we will creat
           7  if not os.path.isdir(destination):
           8      os.makedirs(destination)
           9
          10  # if we have folder called 'train' (train folder contains both .asm files and
          11  # for every file that we have in our 'asmFiles' directory we check if it is er
          12  # 'byteFiles' folder
          13
          14  # so by the end of this snippet we will separate all the .byte files and .asm
          15  if os.path.isdir(source):
          16      os.rename(source,'asmFiles')
          17      source='asmFiles'
          18      data_files = os.listdir(source)
          19      for file in data_files:
          20          #print(file)
          21          if (file.endswith("bytes")):
          22              #print(source+'\\'+file)
          23              shutil.move(source+'\\'+file,destination)
```

```
01azqd4InC7m9JpocGv5.asm
01azqd4InC7m9JpocGv5.bytes
asmFiles\01azqd4InC7m9JpocGv5.bytes
01IsoiSMh5gxyDYTl4CB.asm
01IsoiSMh5gxyDYTl4CB.bytes
asmFiles\01IsoiSMh5gxyDYTl4CB.bytes
01jsnpXSAlgw6aPeDxrU.asm
01jsnpXSAlgw6aPeDxrU.bytes
asmFiles\01jsnpXSAlgw6aPeDxrU.bytes
01kcPWA9K2BOxQeS5Rju.asm
01kcPWA9K2BOxQeS5Rju.bytes
asmFiles\01kcPWA9K2BOxQeS5Rju.bytes
01SuzwMJEIXsK7A8dQbl.asm
01SuzwMJEIXsK7A8dQbl.bytes
asmFiles\01SuzwMJEIXsK7A8dQbl.bytes
02IOCvYEy8mjiuAQHax3.asm
02IOCvYEy8mjiuAQHax3.bytes
asmFiles\02IOCvYEy8mjiuAQHax3.bytes
02JqQ7H3yEoD8viYWlmS.asm
```

## 3.1. Distribution of malware classes in whole data set

```
1  Y=pd.read_csv("trainLabels.csv")
2  total = len(Y)*1.
3  ax=sns.countplot(x="Class", data=Y)
4  for p in ax.patches:
5          ax.annotate('{:.1f}%'.format(100*p.get_height()/total), (p.get_x()+0.1
6
7  #put 11 ticks (therefore 10 steps), from 0 to the total number of rows in the
8  ax.yaxis.set_ticks(np.linspace(0, total, 11))
9
10 #adjust the ticklabel to the desired format, without changing the position of
11 ax.set_yticklabels(map('{:.1f}%'.format, 100*ax.yaxis.get_majorticklocs()/tota
12 plt.show()
```

<IPython.core.display.Javascript object>



## 3.2. Feature extraction

### 3.2.1 File size of byte files as a feature

```
In [7]:   1  #file sizes of byte files
          2
          3  files=os.listdir('byteFiles')
          4  filenames=Y['Id'].tolist()
          5  class_y=Y['Class'].tolist()
          6  class_bytes=[]
          7  sizebytes=[]
          8  fnames=[]
          9  for file in files:
         10      # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
         11      # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=356157170€
         12      # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=1519
         13      # read more about os.stat: here https://www.tutorialspoint.com/python/os_s
         14      statinfo=os.stat('byteFiles/'+file)
         15      # split the file name at '.' and take the first part of it i.e the file no
         16      file=file.split('.')[0]
         17      if any(file == filename for filename in filenames):
         18          i=filenames.index(file)
         19          class_bytes.append(class_y[i])
         20          # converting into Mb's
         21          sizebytes.append(statinfo.st_size/(1024.0*1024.0))
         22          fnames.append(file)
         23  data_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes}
         24  print (data_size_byte.head())
```

```
   Class                ID      size
0      9  01azqd4InC7m9JpocGv5  0.000000
1      2  01IsoiSMh5gxyDYTl4CB  5.538818
2      9  01jsnpXSAlgw6aPeDxrU  3.887939
3      1  01kcPWA9K2BOxQeS5Rju  0.574219
4      8  01SuzwMJEIXsK7A8dQbl  0.370850
```

## 3.2.2 box plots of file size (.byte files) feature

```
1  #boxplot of byte files
2  ax = sns.boxplot(x="Class", y="size", data=data_size_byte)
3  plt.title("boxplot of .bytes file sizes")
4  plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .bytes file sizes

### 3.2.3 feature extraction from byte files

```python
#removal of addres from byte files
# contents of .byte files
# ----------------
#00401000 56 8D 44 24 08 50 8B F1 E8 1C 1B 00 00 C7 06 08
#-------------------
#we remove the starting address 00401000

files = os.listdir('byteFiles')
filenames=[]
array=[]
for file in files:
    if(file.endswith("bytes")):
        tfile=file.split('.')[0]
        text_file = open('byteFiles/'+tfile+".txt", 'w+')
        with open('byteFiles/'+file,"r") as fp:
            lines=""
            for line in fp:
                a=line.rstrip().split(" ")[1:]
                b=' '.join(a)
                b=b+"\n"
                text_file.write(b)
            fp.close()
            os.remove('byteFiles/'+file)
        text_file.close()

files = os.listdir('byteFiles')
filenames2=[]
feature_matrix = np.zeros((len(files),257),dtype=int)
k=0


#program to convert into bag of words of bytefiles
#this is custom-built bag of words this is unigram bag of words
byte_feature_file=open('result.csv','w+')
byte_feature_file.write("ID,0,1,2,3,4,5,6,7,8,9,0a,0b,0c,0d,0e,0f,10,11,12,13,
for file in files:
    filenames2.append(f)
    byte_feature_file.write(file+",")
    if(file.endswith("txt")):
        with open('byteFiles/'+file,"r") as byte_flie:
            for lines in byte_flie:
                line=lines.rstrip().split(" ")
                for hex_code in line:
                    if hex_code=='??':
                        feature_matrix[k][256]+=1
                    else:
                        feature_matrix[k][int(hex_code,16)]+=1
        byte_flie.close()
    for i in feature_matrix[k]:
        byte_feature_file.write(str(i)+",")
    byte_feature_file.write("\n")

    k += 1

byte_feature_file.close()
```

```
In [9]:  1  byte_features=pd.read_csv("result.csv")
         2  print (byte_features.head())
```

```
                       ID        0      1      2      3      4      5      6      7  \
0  01azqd4InC7m9JpocGv5   601905   3905   2816   3832   3345   3242   3650   3201
1  01IsoiSMh5gxyDYTl4CB    39755   8337   7249   7186   8663   6844   8420   7589
2  01jsnpXSAlgw6aPeDxrU    93506   9542   2568   2438   8925   9330   9007   2342
3  01kcPWA9K2BOxQeS5Rju    21091   1213    726    817   1257    625    550    523
4  01SuzwMJEIXsK7A8dQbl    19764    710    302    433    559    410    262    249

         8  ...      f7     f8     f9     fa     fb     fc     fd      fe      ff      ??

0     2965  ...    2804   3687   3101   3211   3097   2758   3099    2759    5753    1824

1     9291  ...     451   6536    439    281    302   7639    518   17001   54902    8588

2     9107  ...    2325   2358   2242   2885   2863   2471   2786    2680   49144     468

3     1078  ...     478    873    485    462    516   1133    471     761    7998   13940

4      422  ...     847    947    350    209    239    653    221     242    2199    9008

[5 rows x 258 columns]
```

```
In [10]:  1  result = pd.merge(byte_features, data_size_byte,on='ID', how='left')
          2  result.head()
```

Out[10]:

| | ID | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... | f9 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 01azqd4InC7m9JpocGv5 | 601905 | 3905 | 2816 | 3832 | 3345 | 3242 | 3650 | 3201 | 2965 | ... | 3101 | 3... |
| 1 | 01IsoiSMh5gxyDYTl4CB | 39755 | 8337 | 7249 | 7186 | 8663 | 6844 | 8420 | 7589 | 9291 | ... | 439 | 2 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 93506 | 9542 | 2568 | 2438 | 8925 | 9330 | 9007 | 2342 | 9107 | ... | 2242 | 28 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 21091 | 1213 | 726 | 817 | 1257 | 625 | 550 | 523 | 1078 | ... | 485 | 4 |
| 4 | 01SuzwMJEIXsK7A8dQbl | 19764 | 710 | 302 | 433 | 559 | 410 | 262 | 249 | 422 | ... | 350 | 2 |

5 rows × 260 columns

```
In [11]:  1  result.to_csv("bytefiles_size.csv", index=False)
```

```
In [12]:   1  # https://stackoverflow.com/a/29651514
           2  def normalize(df):
           3      result1 = df.copy()
           4      for feature_name in df.columns:
           5          if (str(feature_name) != str('ID') and str(feature_name)!=str('Class')
           6              max_value = df[feature_name].max()
           7              min_value = df[feature_name].min()
           8              result1[feature_name] = (df[feature_name] - min_value) / (max_valu
           9      return result1
          10  result = normalize(result)
```

```
In [13]:  1  data_y = result['Class']
          2  result.head()
```

Out[13]:

| | ID | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|
| **0** | 01azqd4InC7m9JpocGv5 | 0.262806 | 0.005498 | 0.001567 | 0.002067 | 0.002048 | 0.001835 | 0.002058 |
| **1** | 01IsoiSMh5gxyDYTl4CB | 0.017358 | 0.011737 | 0.004033 | 0.003876 | 0.005303 | 0.003873 | 0.004747 |
| **2** | 01jsnpXSAlgw6aPeDxrU | 0.040827 | 0.013434 | 0.001429 | 0.001315 | 0.005464 | 0.005280 | 0.005078 |
| **3** | 01kcPWA9K2BOxQeS5Rju | 0.009209 | 0.001708 | 0.000404 | 0.000441 | 0.000770 | 0.000354 | 0.000310 |
| **4** | 01SuzwMJEIXsK7A8dQbl | 0.008629 | 0.001000 | 0.000168 | 0.000234 | 0.000342 | 0.000232 | 0.000148 |

5 rows × 260 columns

```
In [14]:  1  result.to_csv('bytefiles_size_normal.csv', index=False)
```

### 3.2.4 Multivariate Analysis

```
1  #multivariate analysis on byte files
2  #this is with perplexity 50
3  xtsne=TSNE(perplexity=50)
4  results=xtsne.fit_transform(result.drop(['ID','Class'], axis=1))
5  vis_x = results[:, 0]
6  vis_y = results[:, 1]
7  plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
8  plt.colorbar(ticks=range(10))
9  plt.clim(0.5, 9)
10 plt.show()
```

<IPython.core.display.Javascript object>

```
In [59]:    1  #this is with perplexity 30
            2  xtsne=TSNE(perplexity=30)
            3  results=xtsne.fit_transform(result.drop(['ID','Class'], axis=1))
            4  vis_x = results[:, 0]
            5  vis_y = results[:, 1]
            6  plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
            7  plt.colorbar(ticks=range(10))
            8  plt.clim(0.5, 9)
            9  plt.show()
```

<IPython.core.display.Javascript object>



## Train Test split
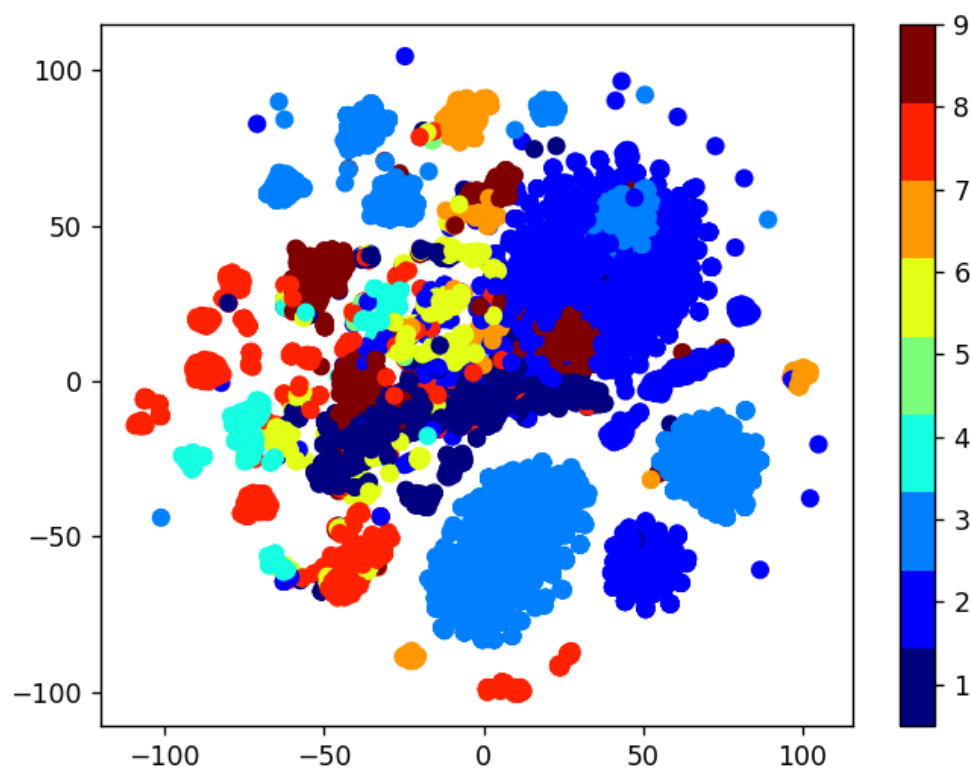
```
In [60]:    1  data_y = result['Class']
            2  # split the data into test and train by maintaining same distribution of outpu
            3  X_train, X_test, y_train, y_test = train_test_split(result.drop(['ID','Class']
            4  # split the train data into train and cross validation by maintaining same dis
            5  X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,stratify=y_tr
```

```python
print('Number of data points in train data:', X_train.shape[0])
print('Number of data points in test data:', X_test.shape[0])
print('Number of data points in cross validation data:', X_cv.shape[0])
```
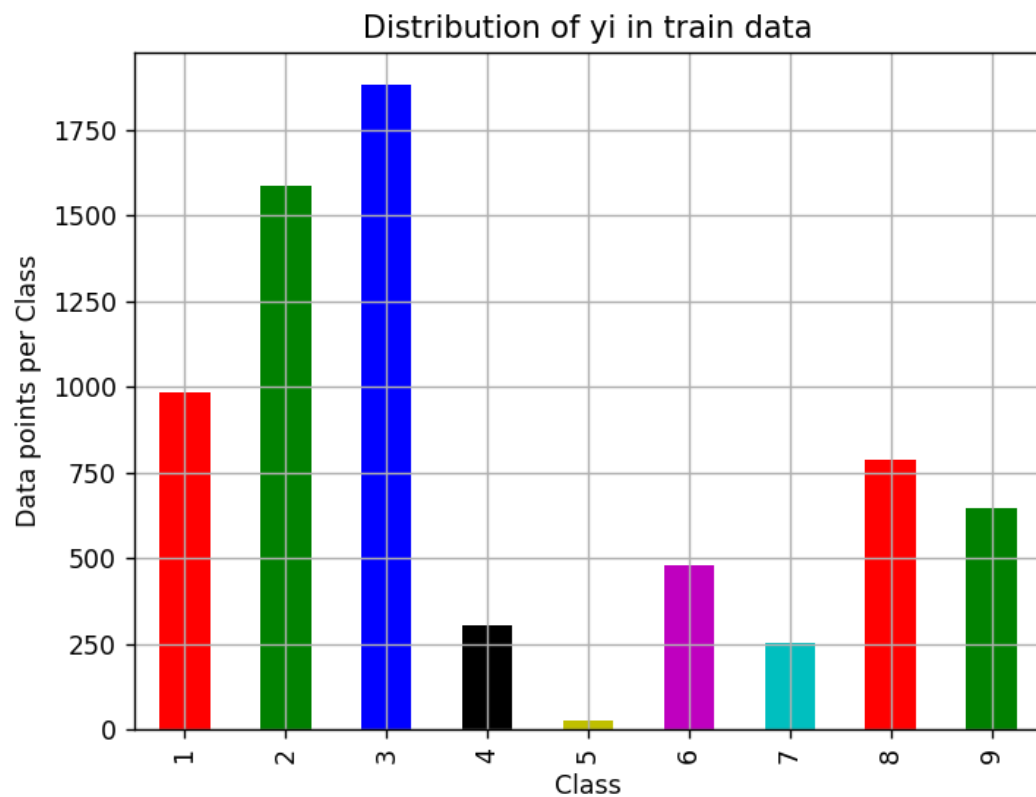
```
Number of data points in train data: 6955
Number of data points in test data: 2174
Number of data points in cross validation data: 1739
```

```python
# it returns a dict, keys as class labels and values as the number of data poi
train_class_distribution = y_train.value_counts().sortlevel()
test_class_distribution = y_test.value_counts().sortlevel()
cv_class_distribution = y_cv.value_counts().sortlevel()

my_colors = 'rgbkymc'
train_class_distribution.plot(kind='bar', color=my_colors)
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.args
# -(train_class_distribution.values): the minus sign will give us in decreasir
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.

print('-'*80)
my_colors = 'rgbkymc'
test_class_distribution.plot(kind='bar', color=my_colors)
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.args
# -(train_class_distribution.values): the minus sign will give us in decreasir
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.v

print('-'*80)
my_colors = 'rgbkymc'
cv_class_distribution.plot(kind='bar', color=my_colors)
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.args
# -(train_class_distribution.values): the minus sign will give us in decreasir
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.val
```

<IPython.core.display.Javascript object>

## Distribution of yi in train data



```
Number of data points in class 3 : 1883 ( 27.074 %)
Number of data points in class 2 : 1586 ( 22.804 %)
Number of data points in class 1 : 986 ( 14.177 %)
Number of data points in class 8 : 786 ( 11.301 %)
Number of data points in class 9 : 648 ( 9.317 %)
Number of data points in class 6 : 481 ( 6.916 %)
Number of data points in class 4 : 304 ( 4.371 %)
Number of data points in class 7 : 254 ( 3.652 %)
Number of data points in class 5 : 27 ( 0.388 %)
----------------------------------------------------------------------------
--

<IPython.core.display.Javascript object>
```

Distribution of yi in test data

Number of data points in class 3 : 588 ( 27.047 %)
Number of data points in class 2 : 496 ( 22.815 %)
Number of data points in class 1 : 308 ( 14.167 %)
Number of data points in class 8 : 246 ( 11.316 %)
Number of data points in class 9 : 203 ( 9.338 %)
Number of data points in class 6 : 150 ( 6.9 %)
Number of data points in class 4 : 95 ( 4.37 %)
Number of data points in class 7 : 80 ( 3.68 %)
Number of data points in class 5 : 8 ( 0.368 %)
-------------------------------------------------------------------------------
--

<IPython.core.display.Javascript object>

# Distribution of yi in cross validation data



Number of data points in class 3 : 471 ( 27.085 %)
Number of data points in class 2 : 396 ( 22.772 %)
Number of data points in class 1 : 247 ( 14.204 %)
Number of data points in class 8 : 196 ( 11.271 %)
Number of data points in class 9 : 162 ( 9.316 %)
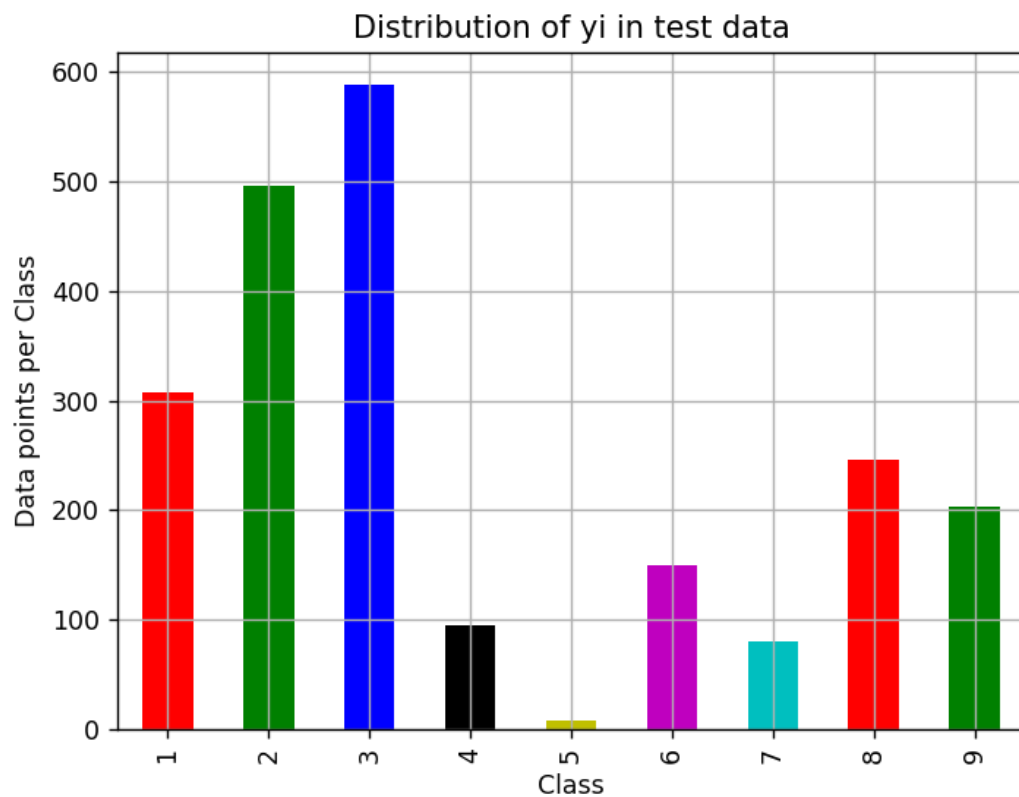Number of data points in class 6 : 120 ( 6.901 %)
Number of data points in class 4 : 76 ( 4.37 %)
Number of data points in class 7 : 64 ( 3.68 %)
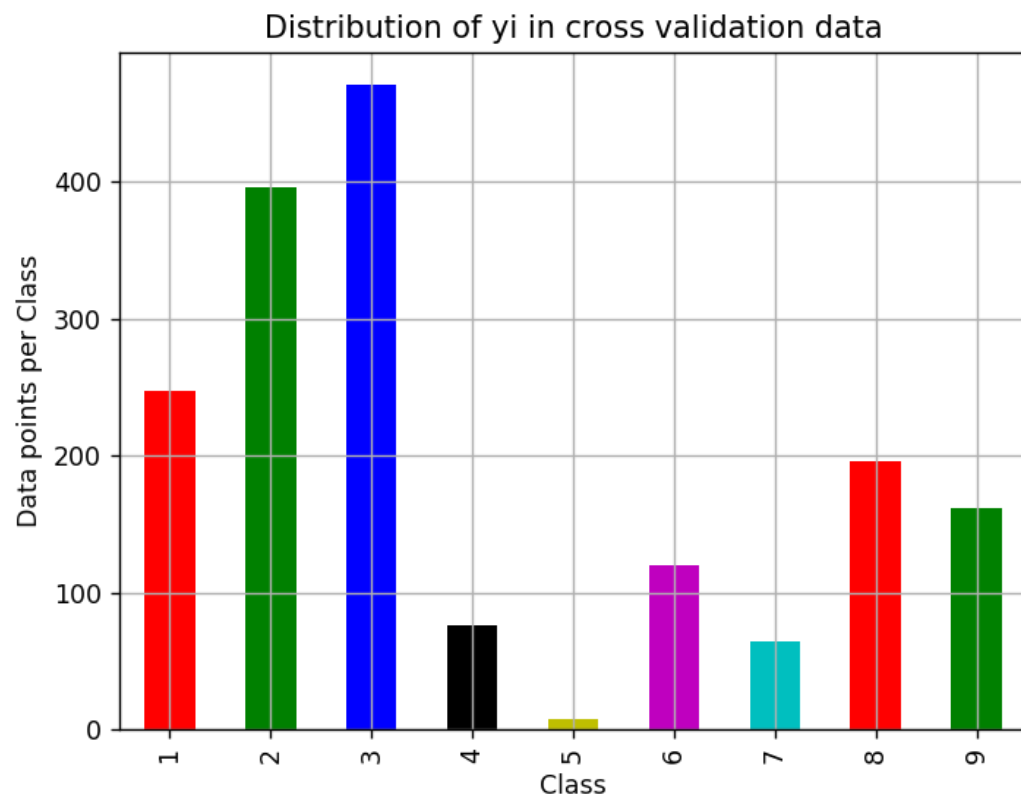Number of data points in class 5 : 7 ( 0.403 %)

```python
In [61]:   1  def plot_confusion_matrix(test_y, predict_y):
           2      C = confusion_matrix(test_y, predict_y)
           3      print("Number of misclassified points ",(len(test_y)-np.trace(C))/len(test
           4      # C = 9,9 matrix, each cell (i,j) represents number of points of class i c
           5
           6      A =(((C.T)/(C.sum(axis=1))).T)
           7      #divid each element of the confusion matrix with the sum of elements in th
           8
           9      # C = [[1, 2],
          10      #      [3, 4]]
          11      # C.T = [[1, 3],
          12      #        [2, 4]]
          13      # C.sum(axis = 1)  axis=0 corresonds to columns and axis=1 corresponds to
          14      # C.sum(axix =1) = [[3, 7]]
          15      # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
          16      #                           [2/3, 4/7]]
          17
          18      # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
          19      #                             [3/7, 4/7]]
          20      # sum of row elements = 1
          21
          22      B =(C/C.sum(axis=0))
          23      #divid each element of the confusion matrix with the sum of elements in th
          24      # C = [[1, 2],
          25      #      [3, 4]]
          26      # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to
          27      # C.sum(axix =0) = [[4, 6]]
          28      # (C/C.sum(axis=0)) = [[1/4, 2/6],
          29      #                      [3/4, 4/6]]
          30
          31      labels = [1,2,3,4,5,6,7,8,9]
          32      cmap=sns.light_palette("green")
          33      # representing A in heatmap format
          34      print("-"*50, "Confusion matrix", "-"*50)
          35      plt.figure(figsize=(10,5))
          36      sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytick
          37      plt.xlabel('Predicted Class')
          38      plt.ylabel('Original Class')
          39      plt.show()
          40
          41      print("-"*50, "Precision matrix", "-"*50)
          42      plt.figure(figsize=(10,5))
          43      sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytick
          44      plt.xlabel('Predicted Class')
          45      plt.ylabel('Original Class')
          46      plt.show()
          47      print("Sum of columns in precision matrix",B.sum(axis=0))
          48
          49      # representing B in heatmap format
          50      print("-"*50, "Recall matrix"    , "-"*50)
          51      plt.figure(figsize=(10,5))
          52      sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytick
          53      plt.xlabel('Predicted Class')
          54      plt.ylabel('Original Class')
          55      plt.show()
          56      print("Sum of rows in precision matrix",A.sum(axis=1))
```

# 4. Machine Learning Models

## 4.1. Machine Leaning Models on bytes files

### 4.1.1. Random Model

```
In [62]:   1  # we need to generate 9 numbers and the sum of numbers should be 1
           2  # one solution is to genarate 9 numbers and divide each of the numbers by thei
           3  # ref: https://stackoverflow.com/a/18662466/4084039
           4
           5  test_data_len = X_test.shape[0]
           6  cv_data_len = X_cv.shape[0]
           7
           8  # we create a output array that has exactly same size as the CV data
           9  cv_predicted_y = np.zeros((cv_data_len,9))
          10  for i in range(cv_data_len):
          11      rand_probs = np.random.rand(1,9)
          12      cv_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
          13  print("Log loss on Cross Validation Data using Random Model",log_loss(y_cv,cv_
          14
          15
          16  # Test-Set error.
          17  #we create a output array that has exactly same as the test data
          18  test_predicted_y = np.zeros((test_data_len,9))
          19  for i in range(test_data_len):
          20      rand_probs = np.random.rand(1,9)
          21      test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
          22  print("Log loss on Test Data using Random Model",log_loss(y_test,test_predicte
          23
          24  predicted_y =np.argmax(test_predicted_y, axis=1)
          25  plot_confusion_matrix(y_test, predicted_y+1)
```

```
Log loss on Cross Validation Data using Random Model 2.45615644965
Log loss on Test Data using Random Model 2.48503905509
Number of misclassified points  88.5004599816
-------------------------------------------------- Confusion matrix ----------
------------------------------------------

<IPython.core.display.Javascript object>
```



```
-------------------------------------------------- Precision matrix ----------
------------------------------------------

<IPython.core.display.Javascript object>
```

| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.156 | 0.120 | 0.131 | 0.117 | 0.154 | 0.148 | 0.172 | 0.145 | 0.129 |
| 2 | 0.197 | 0.264 | 0.236 | 0.248 | 0.211 | 0.195 | 0.243 | 0.237 | 0.218 |
| 3 | 0.279 | 0.283 | 0.266 | 0.235 | 0.276 | 0.276 | 0.243 | 0.274 | 0.302 |
| 4 | 0.061 | 0.019 | 0.061 | 0.022 | 0.041 | 0.033 | 0.056 | 0.058 | 0.040 |
| 5 | 0.000 | 0.004 | 0.000 | 0.004 | 0.000 | 0.005 | 0.007 | 0.000 | 0.012 |
| 6 | 0.066 | 0.078 | 0.074 | 0.091 | 0.065 | 0.090 | 0.060 | 0.071 | 0.032 |
| 7 | 0.020 | 0.058 | 0.035 | 0.039 | 0.024 | 0.043 | 0.030 | 0.037 | 0.044 |
| 8 | 0.102 | 0.105 | 0.118 | 0.117 | 0.130 | 0.124 | 0.112 | 0.100 | 0.113 |
| 9 | 0.119 | 0.070 | 0.079 | 0.126 | 0.098 | 0.086 | 0.078 | 0.079 | 0.109 |

```
Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
--------------------------------------------------- Recall matrix -------------
------------------------------------
```

`<IPython.core.display.Javascript object>`



| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.123 | 0.101 | 0.097 | 0.088 | 0.123 | 0.101 | 0.149 | 0.114 | 0.104 |
| 2 | 0.097 | 0.137 | 0.109 | 0.115 | 0.105 | 0.083 | 0.131 | 0.115 | 0.109 |
| 3 | 0.116 | 0.124 | 0.104 | 0.092 | 0.116 | 0.099 | 0.111 | 0.112 | 0.128 |
| 4 | 0.158 | 0.053 | 0.147 | 0.053 | 0.105 | 0.074 | 0.158 | 0.147 | 0.105 |
| 5 | 0.000 | 0.125 | 0.000 | 0.125 | 0.000 | 0.125 | 0.250 | 0.000 | 0.375 |
| 6 | 0.107 | 0.133 | 0.113 | 0.140 | 0.107 | 0.127 | 0.107 | 0.113 | 0.053 |
| 7 | 0.062 | 0.188 | 0.100 | 0.113 | 0.075 | 0.113 | 0.100 | 0.113 | 0.138 |
| 8 | 0.102 | 0.110 | 0.110 | 0.110 | 0.130 | 0.106 | 0.122 | 0.098 | 0.114 |
| 9 | 0.143 | 0.089 | 0.089 | 0.143 | 0.118 | 0.089 | 0.103 | 0.094 | 0.133 |

```
Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
```

## 4.1.2. K Nearest Neighbour Classification

```
In [68]:   1  # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/m
           2  # --------------------------
           3  # default parameter
           4  # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', lea
           5  # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)
           6
           7  # methods of
           8  # fit(X, y) : Fit the model using X as training data and y as target values
           9  # predict(X):Predict the class labels for the provided data
          10  # predict_proba(X):Return probability estimates for the test data X.
          11  #-----------------------------------
          12  # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/
          13  #-----------------------------------
          14
          15
          16  # find more about CalibratedClassifierCV here at http://scikit-learn.org/stabl
          17  # ----------------------------
          18  # default paramters
          19  # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sign
          20  #
          21  # some of the methods of CalibratedClassifierCV()
          22  # fit(X, y[, sample_weight])    Fit the calibrated model
          23  # get_params([deep])    Get parameters for this estimator.
          24  # predict(X)    Predict the target of new samples.
          25  # predict_proba(X)  Posterior probabilities of classification
          26  #-----------------------------------
          27  # video link:
          28  #-----------------------------------
          29
          30  alpha = [x for x in range(1, 15, 2)]
          31  cv_log_error_array=[]
          32  for i in alpha:
          33      k_cfl=KNeighborsClassifier(n_neighbors=i)
          34      k_cfl.fit(X_train,y_train)
          35      sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
          36      sig_clf.fit(X_train, y_train)
          37      predict_y = sig_clf.predict_proba(X_cv)
          38      cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cfl.classes_,
          39
          40  for i in range(len(cv_log_error_array)):
          41      print ('log_loss for k = ',alpha[i],'is',cv_log_error_array[i])
          42
          43  best_alpha = np.argmin(cv_log_error_array)
          44
          45  fig, ax = plt.subplots()
          46  ax.plot(alpha, cv_log_error_array,c='g')
          47  for i, txt in enumerate(np.round(cv_log_error_array,3)):
          48      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
          49  plt.grid()
          50  plt.title("Cross Validation Error for each alpha")
          51  plt.xlabel("Alpha i's")
          52  plt.ylabel("Error measure")
          53  plt.show()
          54
          55  k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
          56  k_cfl.fit(X_train,y_train)
```

```
57  sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
58  sig_clf.fit(X_train, y_train)
59
60  predict_y = sig_clf.predict_proba(X_train)
61  print ('For values of best alpha = ', alpha[best_alpha], "The train log loss i
62  predict_y = sig_clf.predict_proba(X_cv)
63  print('For values of best alpha = ', alpha[best_alpha], "The cross validation
64  predict_y = sig_clf.predict_proba(X_test)
65  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:
66  plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```
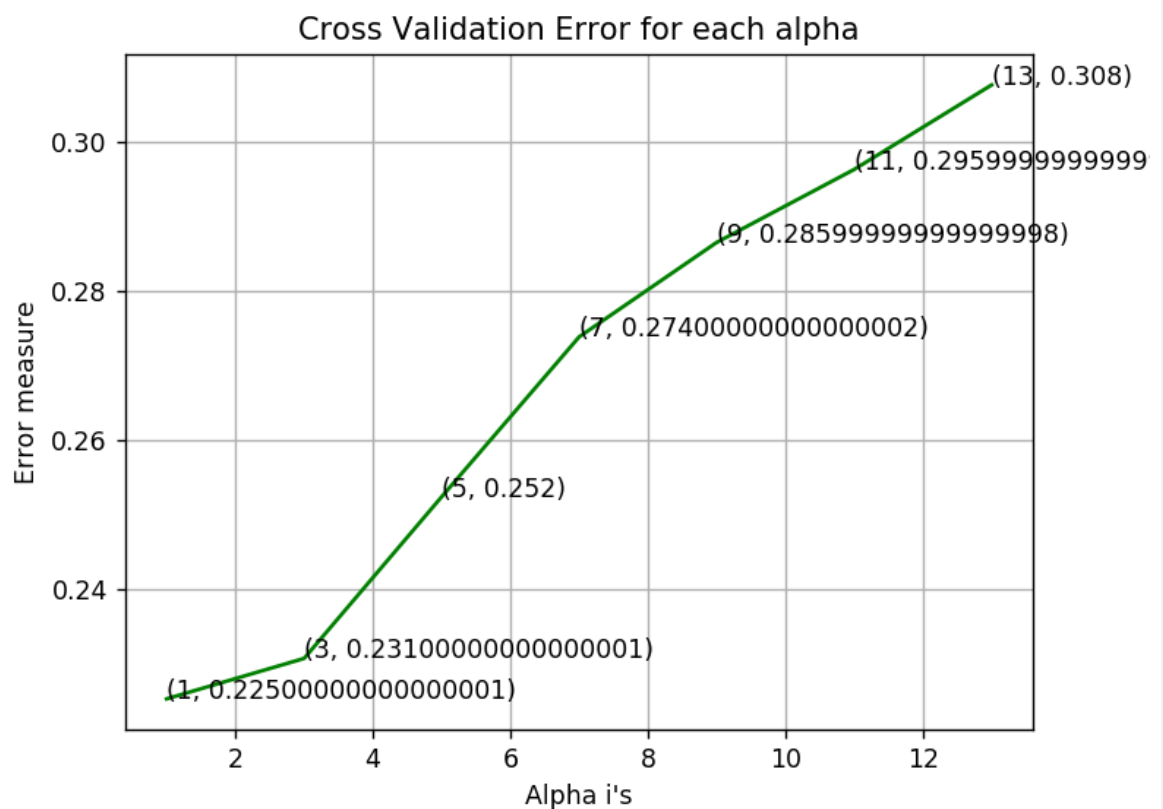
```
log_loss for k =   1 is 0.225386237304
log_loss for k =   3 is 0.230795229168
log_loss for k =   5 is 0.252421408646
log_loss for k =   7 is 0.273827486888
log_loss for k =   9 is 0.286469181555
log_loss for k =  11 is 0.29623391147
log_loss for k =  13 is 0.307551203154

<IPython.core.display.Javascript object>
```
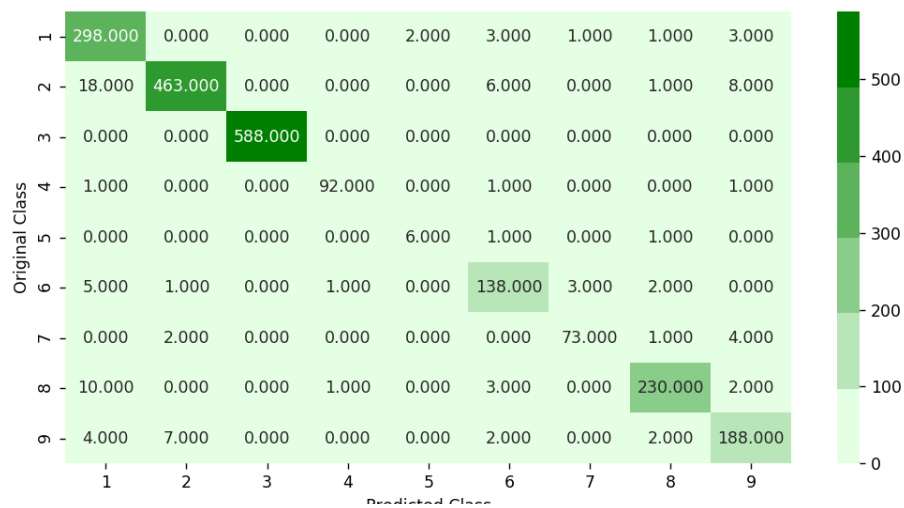


Cross Validation Error for each alpha

```
For values of best alpha =   1 The train log loss is: 0.0782947669247
For values of best alpha =   1 The cross validation log loss is: 0.225386237304
For values of best alpha =   1 The test log loss is: 0.241508604195
Number of misclassified points  4.50781968721
---------------------------------------------------- Confusion matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>
```
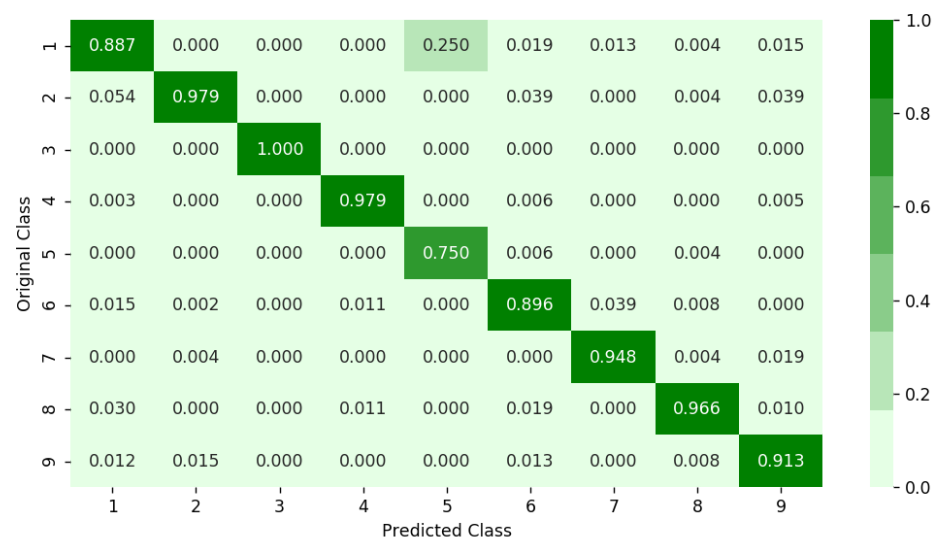
-------------------------------------------------- Precision matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>



Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
-------------------------------------------------- Recall matrix -------------
-------------------------------------

<IPython.core.display.Javascript object>

Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]

## 4.1.3. Logistic Regression

```
In [71]:   1  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/ge
           2  # ------------------------------
           3  # default parameters
           4  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_i
           5  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_
           6  # class_weight=None, warm_start=False, average=False, n_iter=None)
           7
           8  # some of methods
           9  # fit(X, y[, coef_init, intercept_init, …]) Fit linear model with Stochastic (
          10  # predict(X)    Predict class labels for samples in X.
          11
          12  #------------------------------
          13  # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/
          14  #------------------------------
          15
          16  alpha = [10 ** x for x in range(-5, 4)]
          17  cv_log_error_array=[]
          18  for i in alpha:
          19      logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
          20      logisticR.fit(X_train,y_train)
          21      sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
          22      sig_clf.fit(X_train, y_train)
          23      predict_y = sig_clf.predict_proba(X_cv)
          24      cv_log_error_array.append(log_loss(y_cv, predict_y, labels=logisticR.class
          25
          26  for i in range(len(cv_log_error_array)):
          27      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
          28
          29  best_alpha = np.argmin(cv_log_error_array)
          30
          31  fig, ax = plt.subplots()
          32  ax.plot(alpha, cv_log_error_array,c='g')
          33  for i, txt in enumerate(np.round(cv_log_error_array,3)):
          34      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
          35  plt.grid()
          36  plt.title("Cross Validation Error for each alpha")
          37  plt.xlabel("Alpha i's")
          38  plt.ylabel("Error measure")
          39  plt.show()
          40
          41  logisticR=LogisticRegression(penalty='l2',C=alpha[best_alpha],class_weight='ba
          42  logisticR.fit(X_train,y_train)
          43  sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
          44  sig_clf.fit(X_train, y_train)
          45  pred_y=sig_clf.predict(X_test)
          46
          47  predict_y = sig_clf.predict_proba(X_train)
          48  print ('log loss for train data',log_loss(y_train, predict_y, labels=logisticR
          49  predict_y = sig_clf.predict_proba(X_cv)
          50  print ('log loss for cv data',log_loss(y_cv, predict_y, labels=logisticR.class
          51  predict_y = sig_clf.predict_proba(X_test)
          52  print ('log loss for test data',log_loss(y_test, predict_y, labels=logisticR.c
          53  plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```

```
log_loss for c =  1e-05 is 1.56916911178
log_loss for c =  0.0001 is 1.57336384417
```

```
log_loss for c =  0.001 is 1.53598598273
log_loss for c =  0.01 is 1.01720972418
log_loss for c =  0.1 is 0.857766083873
log_loss for c =  1 is 0.711154393309
log_loss for c =  10 is 0.583929522635
log_loss for c =  100 is 0.549929846589
log_loss for c =  1000 is 0.624746769121
```

<IPython.core.display.Javascript object>



Cross Validation Error for each alpha

log loss for train data 0.498923428696
log loss for cv data 0.549929846589
log loss for test data 0.528347316704
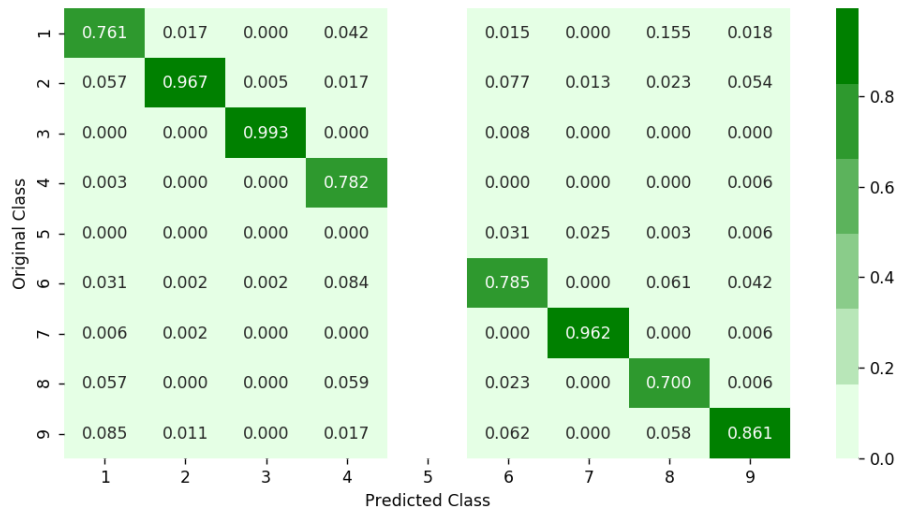Number of misclassified points  12.3275068997
------------------------------------------------ Confusion matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 242.000 | 8.000 | 0.000 | 5.000 | 0.000 | 2.000 | 0.000 | 48.000 | 3.000 |
| 2 | 18.000 | 446.000 | 3.000 | 2.000 | 0.000 | 10.000 | 1.000 | 7.000 | 9.000 |
| 3 | 0.000 | 0.000 | 587.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 |
| 4 | 1.000 | 0.000 | 0.000 | 93.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 4.000 | 2.000 | 1.000 | 1.000 |
| 6 | 10.000 | 1.000 | 1.000 | 10.000 | 0.000 | 102.000 | 0.000 | 19.000 | 7.000 |
| 7 | 2.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 76.000 | 0.000 | 1.000 |
| 8 | 18.000 | 0.000 | 0.000 | 7.000 | 0.000 | 3.000 | 0.000 | 217.000 | 1.000 |
| 9 | 27.000 | 5.000 | 0.000 | 2.000 | 0.000 | 8.000 | 0.000 | 18.000 | 143.000 |

---------------------------------------------- Precision matrix -----------------------------------------------

<IPython.core.display.Javascript object>



|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.761 | 0.017 | 0.000 | 0.042 | 0.015 | 0.000 | 0.155 | 0.018 |
| 2 | 0.057 | 0.967 | 0.005 | 0.017 | 0.077 | 0.013 | 0.023 | 0.054 |
| 3 | 0.000 | 0.000 | 0.993 | 0.000 | 0.008 | 0.000 | 0.000 | 0.000 |
| 4 | 0.003 | 0.000 | 0.000 | 0.782 | 0.000 | 0.000 | 0.000 | 0.006 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.031 | 0.025 | 0.003 | 0.006 |
| 6 | 0.031 | 0.002 | 0.002 | 0.084 | 0.785 | 0.000 | 0.061 | 0.042 |
| 7 | 0.006 | 0.002 | 0.000 | 0.000 | 0.000 | 0.962 | 0.000 | 0.006 |
| 8 | 0.057 | 0.000 | 0.000 | 0.059 | 0.023 | 0.000 | 0.700 | 0.006 |
| 9 | 0.085 | 0.011 | 0.000 | 0.017 | 0.062 | 0.000 | 0.058 | 0.861 |

Sum of columns in precision matrix [  1.   1.   1.   1.   nan   1.   1.   1.   1.]
---------------------------------------------- Recall matrix -------------------------------------------------

<IPython.core.display.Javascript object>

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.786 | 0.026 | 0.000 | 0.016 | 0.000 | 0.006 | 0.000 | 0.156 | 0.010 |
| 2 | 0.036 | 0.899 | 0.006 | 0.004 | 0.000 | 0.020 | 0.002 | 0.014 | 0.018 |
| 3 | 0.000 | 0.000 | 0.998 | 0.000 | 0.000 | 0.002 | 0.000 | 0.000 | 0.000 |
| 4 | 0.011 | 0.000 | 0.000 | 0.979 | 0.000 | 0.000 | 0.000 | 0.000 | 0.011 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.500 | 0.250 | 0.125 | 0.125 |
| 6 | 0.067 | 0.007 | 0.007 | 0.067 | 0.000 | 0.680 | 0.000 | 0.127 | 0.047 |
| 7 | 0.025 | 0.013 | 0.000 | 0.000 | 0.000 | 0.000 | 0.950 | 0.000 | 0.013 |
| 8 | 0.073 | 0.000 | 0.000 | 0.028 | 0.000 | 0.012 | 0.000 | 0.882 | 0.004 |
| 9 | 0.133 | 0.025 | 0.000 | 0.010 | 0.000 | 0.039 | 0.000 | 0.089 | 0.704 |

Original Class (rows) / Predicted Class (columns)

```
Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
```

### 4.1.4. Random Forest Classifier

```
In [72]:  1  # --------------------------------
          2  # default parameters
          3  # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', m
          4  # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_l
          5  # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_s
          6  # class_weight=None)
          7
          8  # Some of methods of RandomForestClassifier()
          9  # fit(X, y, [sample_weight])    Fit the SVM model according to the given trair
         10  # predict(X)    Perform classification on samples in X.
         11  # predict_proba (X) Perform classification on samples in X.
         12
         13  # some of attributes of  RandomForestClassifier()
         14  # feature_importances_  : array of shape = [n_features]
         15  # The feature importances (the higher, the more important the feature).
         16
         17  # --------------------------------
         18  # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/
         19  # --------------------------------
         20
         21  alpha=[10,50,100,500,1000,2000,3000]
         22  cv_log_error_array=[]
         23  train_log_error_array=[]
         24  from sklearn.ensemble import RandomForestClassifier
         25  for i in alpha:
         26      r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
         27      r_cfl.fit(X_train,y_train)
         28      sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
         29      sig_clf.fit(X_train, y_train)
         30      predict_y = sig_clf.predict_proba(X_cv)
         31      cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_,
         32
         33  for i in range(len(cv_log_error_array)):
         34      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
         35
         36
         37  best_alpha = np.argmin(cv_log_error_array)
         38
         39  fig, ax = plt.subplots()
         40  ax.plot(alpha, cv_log_error_array,c='g')
         41  for i, txt in enumerate(np.round(cv_log_error_array,3)):
         42      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
         43  plt.grid()
         44  plt.title("Cross Validation Error for each alpha")
         45  plt.xlabel("Alpha i's")
         46  plt.ylabel("Error measure")
         47  plt.show()
         48
         49
         50  r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_
         51  r_cfl.fit(X_train,y_train)
         52  sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
         53  sig_clf.fit(X_train, y_train)
         54
         55  predict_y = sig_clf.predict_proba(X_train)
         56  print('For values of best alpha = ', alpha[best_alpha], "The train log loss is
```

```
57  predict_y = sig_clf.predict_proba(X_cv)
58  print('For values of best alpha = ', alpha[best_alpha], "The cross validation
59  predict_y = sig_clf.predict_proba(X_test)
60  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:
61  plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```
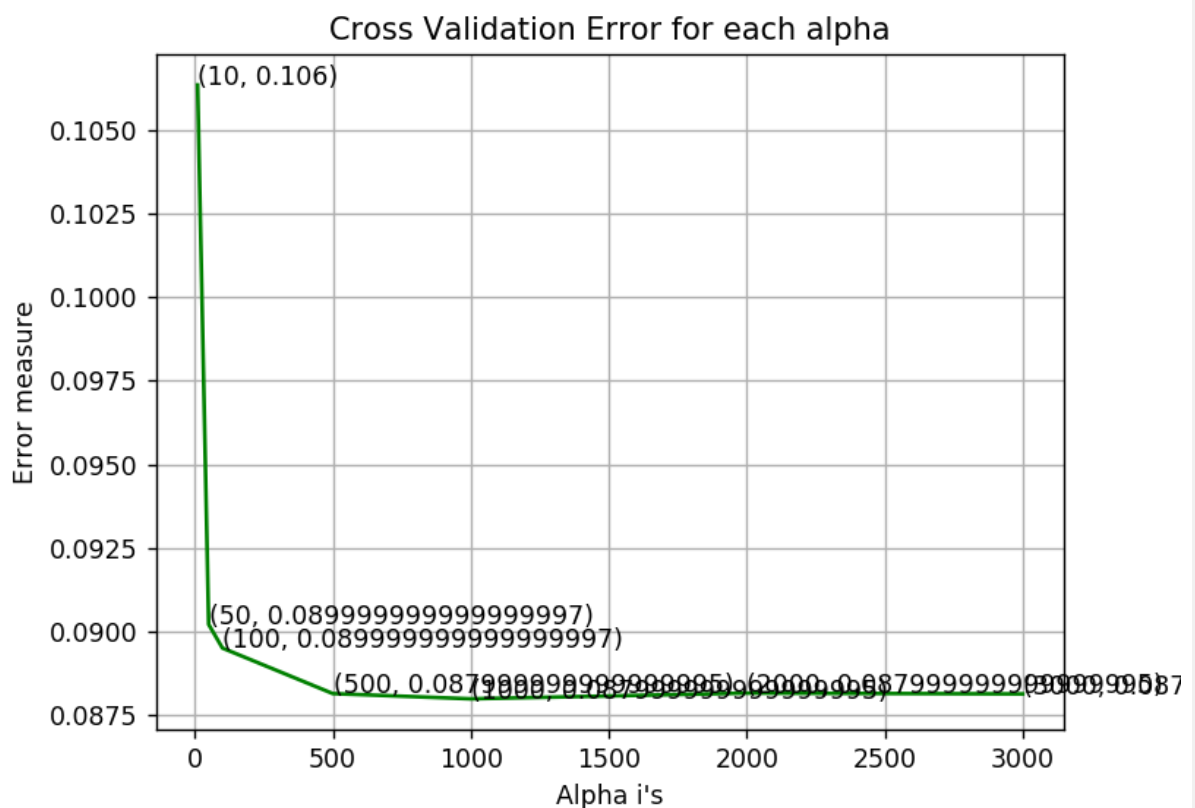
```
log_loss for c =   10 is 0.106357709164
log_loss for c =   50 is 0.0902124124145
log_loss for c =  100 is 0.0895043339776
log_loss for c =  500 is 0.0881420869288
log_loss for c = 1000 is 0.0879849524621
log_loss for c = 2000 is 0.0881566647295
log_loss for c = 3000 is 0.0881318948443
```

<IPython.core.display.Javascript object>

Cross Validation Error for each alpha

(10, 0.106)

(50, 0.08999999999999997)
(100, 0.08999999999999997)

(500, 0.0879...)(1000,...)...(2000,...0879999999...)(3000, ...)

Error measure

Alpha i's

For values of best alpha =  1000 The train log loss is: 0.0266476291801
For values of best alpha =  1000 The cross validation log loss is: 0.087984952
4621
For values of best alpha =  1000 The test log loss is: 0.0858346961407
Number of misclassified points   2.02391904324
------------------------------------------------ Confusion matrix ----------
-----------------------------------------

<IPython.core.display.Javascript object>

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 303.000 | 0.000 | 0.000 | 1.000 | 0.000 | 2.000 | 0.000 | 1.000 | 1.000 |
| 2 | 1.000 | 495.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 588.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 0.000 | 0.000 | 0.000 | 92.000 | 0.000 | 2.000 | 0.000 | 1.000 | 0.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 7.000 | 0.000 | 0.000 | 1.000 | 0.000 |
| 6 | 2.000 | 0.000 | 0.000 | 1.000 | 0.000 | 143.000 | 0.000 | 4.000 | 0.000 |
| 7 | 2.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 76.000 | 2.000 | 0.000 |
| 8 | 9.000 | 2.000 | 0.000 | 2.000 | 0.000 | 2.000 | 0.000 | 229.000 | 2.000 |
| 9 | 3.000 | 0.000 | 1.000 | 0.000 | 0.000 | 1.000 | 0.000 | 1.000 | 197.000 |

Original Class / Predicted Class

---------------------------------------------------- Precision matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>



|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.947 | 0.000 | 0.000 | 0.010 | 0.000 | 0.013 | 0.000 | 0.004 | 0.005 |
| 2 | 0.003 | 0.996 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 0.998 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 0.000 | 0.000 | 0.000 | 0.958 | 0.000 | 0.013 | 0.000 | 0.004 | 0.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.004 | 0.000 |
| 6 | 0.006 | 0.000 | 0.000 | 0.010 | 0.000 | 0.953 | 0.000 | 0.017 | 0.000 |
| 7 | 0.006 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.008 | 0.000 |
| 8 | 0.028 | 0.004 | 0.000 | 0.021 | 0.000 | 0.013 | 0.000 | 0.958 | 0.010 |
| 9 | 0.009 | 0.000 | 0.002 | 0.000 | 0.000 | 0.007 | 0.000 | 0.004 | 0.985 |

Original Class / Predicted Class

Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
---------------------------------------------------- Recall matrix ------------
-------------------------------------

<IPython.core.display.Javascript object>

Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]

## 4.1.5. XgBoost Classification

```
In [74]:   1  # Training a hyper-parameter tuned Xg-Boost regressor on our train data
           2
           3  # find more about XGBClassifier function here http://xgboost.readthedocs.io/er
           4  # ------------------------
           5  # default paramters
           6  # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100
           7  # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma
           8  # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_
           9  # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None,
          10
          11  # some of methods of RandomForestRegressor()
          12  # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppin
          13  # get_params([deep])    Get parameters for this estimator.
          14  # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE:
          15  # get_score(importance_type='weight') -> get the feature importance
          16  # ----------------------
          17  # video link1: https://www.appliedaicourse.com/course/applied-ai-course-online
          18  # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online
          19  # ----------------------
          20
          21  alpha=[10,50,100,500,1000,2000]
          22  cv_log_error_array=[]
          23  for i in alpha:
          24      x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
          25      x_cfl.fit(X_train,y_train)
          26      sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
          27      sig_clf.fit(X_train, y_train)
          28      predict_y = sig_clf.predict_proba(X_cv)
          29      cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_,
          30
          31  for i in range(len(cv_log_error_array)):
          32      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
          33
          34
          35  best_alpha = np.argmin(cv_log_error_array)
          36
          37  fig, ax = plt.subplots()
          38  ax.plot(alpha, cv_log_error_array,c='g')
          39  for i, txt in enumerate(np.round(cv_log_error_array,3)):
          40      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
          41  plt.grid()
          42  plt.title("Cross Validation Error for each alpha")
          43  plt.xlabel("Alpha i's")
          44  plt.ylabel("Error measure")
          45  plt.show()
          46
          47  x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
          48  x_cfl.fit(X_train,y_train)
          49  sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
          50  sig_clf.fit(X_train, y_train)
          51
          52  predict_y = sig_clf.predict_proba(X_train)
          53  print ('For values of best alpha = ', alpha[best_alpha], "The train log loss i
          54  predict_y = sig_clf.predict_proba(X_cv)
          55  print('For values of best alpha = ', alpha[best_alpha], "The cross validation
          56  predict_y = sig_clf.predict_proba(X_test)
```

```
57  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:
58  plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```

```
log_loss for c =   10 is 0.20615980494
log_loss for c =   50 is 0.123888382365
log_loss for c =  100 is 0.099919437112
log_loss for c =  500 is 0.0931035681289
log_loss for c = 1000 is 0.0933084876012
log_loss for c = 2000 is 0.0938395690309
```

&lt;IPython.core.display.Javascript object&gt;



Cross Validation Error for each alpha

```
For values of best alpha =  500 The train log loss is: 0.0225231805824
For values of best alpha =  500 The cross validation log loss is: 0.0931035681
289
For values of best alpha =  500 The test log loss is: 0.0792067651731
Number of misclassified points  1.24195032199
------------------------------------------------- Confusion matrix ----------
-----------------------------------------
```

&lt;IPython.core.display.Javascript object&gt;

| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 306.000 | 0.000 | 0.000 | 1.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 |
| 2 | 1.000 | 495.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 588.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 0.000 | 0.000 | 0.000 | 94.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 6.000 | 1.000 | 0.000 | 1.000 | 0.000 |
| 6 | 1.000 | 1.000 | 0.000 | 1.000 | 0.000 | 147.000 | 0.000 | 0.000 | 0.000 |
| 7 | 2.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 77.000 | 0.000 | 0.000 |
| 8 | 7.000 | 1.000 | 0.000 | 2.000 | 0.000 | 1.000 | 0.000 | 234.000 | 1.000 |
| 9 | 3.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 200.000 |

--------------------------------------------------- Precision matrix ----------------------------------------------------

<IPython.core.display.Javascript object>



| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.956 | 0.000 | 0.000 | 0.010 | 0.000 | 0.007 | 0.000 | 0.000 | 0.000 |
| 2 | 0.003 | 0.996 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 0.000 | 0.000 | 0.000 | 0.959 | 0.000 | 0.000 | 0.000 | 0.000 | 0.005 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.007 | 0.000 | 0.004 | 0.000 |
| 6 | 0.003 | 0.002 | 0.000 | 0.010 | 0.000 | 0.974 | 0.000 | 0.000 | 0.000 |
| 7 | 0.006 | 0.000 | 0.000 | 0.000 | 0.000 | 0.007 | 1.000 | 0.000 | 0.000 |
| 8 | 0.022 | 0.002 | 0.000 | 0.020 | 0.000 | 0.007 | 0.000 | 0.996 | 0.005 |
| 9 | 0.009 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.990 |

Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
--------------------------------------------------- Recall matrix ----------------------------------------------------

<IPython.core.display.Javascript object>

```
Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
```

## 4.1.5. XgBoost Classification with best hyper parameters using RandomSearch

```
In [75]:   1  # https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning
           2  x_cfl=XGBClassifier()
           3
           4  prams={
           5      'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
           6       'n_estimators':[100,200,500,1000,2000],
           7       'max_depth':[3,5,10],
           8      'colsample_bytree':[0.1,0.3,0.5,1],
           9      'subsample':[0.1,0.3,0.5,1]
          10  }
          11  random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jo
          12  random_cfl1.fit(X_train,y_train)
```

```
Fitting 3 folds for each of 10 candidates, totalling 30 fits

[Parallel(n_jobs=-1)]: Done    2 tasks       | elapsed:    26.5s
[Parallel(n_jobs=-1)]: Done    9 tasks       | elapsed:  5.8min
[Parallel(n_jobs=-1)]: Done   19 out of   30 | elapsed:  9.3min remaining:  5.4m
in
[Parallel(n_jobs=-1)]: Done   23 out of   30 | elapsed: 10.1min remaining:  3.1m
in
[Parallel(n_jobs=-1)]: Done   27 out of   30 | elapsed: 14.0min remaining:  1.6m
in
[Parallel(n_jobs=-1)]: Done   30 out of   30 | elapsed: 14.2min finished
```

```
Out[75]: RandomizedSearchCV(cv=None, error_score='raise',
              estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, colsampl
         e_bytree=1,
              gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
              min_child_weight=1, missing=None, n_estimators=100, nthread=-1,
              objective='binary:logistic', reg_alpha=0, reg_lambda=1,
              scale_pos_weight=1, seed=0, silent=True, subsample=1),
              fit_params=None, iid=True, n_iter=10, n_jobs=-1,
              param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15,
         0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'co
         lsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
              pre_dispatch='2*n_jobs', random_state=None, refit=True,
              return_train_score=True, scoring=None, verbose=10)
```

```
In [76]:   1  print (random_cfl1.best_params_)
```

```
{'subsample': 1, 'n_estimators': 500, 'max_depth': 5, 'learning_rate': 0.05,
 'colsample_bytree': 0.5}
```

```
In [80]:    1   # Training a hyper-parameter tuned Xg-Boost regressor on our train data
            2
            3   # find more about XGBClassifier function here http://xgboost.readthedocs.io/er
            4   # -----------------------
            5   # default paramters
            6   # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100
            7   # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gammc
            8   # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_
            9   # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None,
           10
           11   # some of methods of RandomForestRegressor()
           12   # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppir
           13   # get_params([deep])     Get parameters for this estimator.
           14   # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE:
           15   # get_score(importance_type='weight') -> get the feature importance
           16   # ----------------------
           17   # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online
           18   # ----------------------
           19
           20   x_cfl=XGBClassifier(n_estimators=2000, learning_rate=0.05, colsample_bytree=1,
           21   x_cfl.fit(X_train,y_train)
           22   c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
           23   c_cfl.fit(X_train,y_train)
           24
           25   predict_y = c_cfl.predict_proba(X_train)
           26   print ('train loss',log_loss(y_train, predict_y))
           27   predict_y = c_cfl.predict_proba(X_cv)
           28   print ('cv loss',log_loss(y_cv, predict_y))
           29   predict_y = c_cfl.predict_proba(X_test)
           30   print ('test loss',log_loss(y_test, predict_y))
```

```
train loss 0.022540976086
cv loss 0.0928710624158
test loss 0.0782688587098
```

## 4.2 Modeling with .asm files

There are 10868 files of asm
All the files make up about 150 GB
The asm files contains :
1. Address
2. Segments
3. Opcodes
4. Registers
5. function calls
6. APIs
With the help of parallel processing we extracted all the features. In parallel we can use all the cores that are present in our computer.

Here we extracted 52 features from all the asm files which are important.

We read the top solutions and handpicked the features from those papers/v
ideos/blogs.
 Refer:https://www.kaggle.com/c/malware-classification/discussion

## 4.2.1 Feature extraction from asm files

- To extract the unigram features from the .asm files we need to process ~150GB of data
- **Note: Below two cells will take lot of time (over 48 hours to complete)**
- We will provide you the output file of these two cells, which you can directly use it

```
In [ ]:
1   #intially create five folders
2   #first
3   #second
4   #thrid
5   #fourth
6   #fifth
7   #this code tells us about random split of files into five folders
8   folder_1 ='first'
9   folder_2 ='second'
10  folder_3 ='third'
11  folder_4 ='fourth'
12  folder_5 ='fifth'
13  folder_6 = 'output'
14  for i in [folder_1,folder_2,folder_3,folder_4,folder_5,folder_6]:
15      if not os.path.isdir(i):
16          os.makedirs(i)
17
18  source='train/'
19  files = os.listdir('train')
20  ID=df['Id'].tolist()
21  data=range(0,10868)
22  r.shuffle(data)
23  count=0
24  for i in range(0,10868):
25      if i % 5==0:
26          shutil.move(source+files[data[i]],'first')
27      elif i%5==1:
28          shutil.move(source+files[data[i]],'second')
29      elif i%5 ==2:
30          shutil.move(source+files[data[i]],'thrid')
31      elif i%5 ==3:
32          shutil.move(source+files[data[i]],'fourth')
33      elif i%5==4:
34          shutil.move(source+files[data[i]],'fifth')
```

```python
#http://flint.cs.yale.edu/cs421/papers/x86-asm/asm.html

def firstprocess():
    #The prefixes tells about the segments that are present in the asm files
    #There are 450 segments(approx) present in all asm files.
    #this prefixes are best segments that gives us best values.
    #https://en.wikipedia.org/wiki/Data_segment

    prefixes = ['HEADER:','.text:','.Pav:','.idata:','.data:','.bss:','.rdata
    #this are opcodes that are used to get best results
    #https://en.wikipedia.org/wiki/X86_instruction_listings

    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
    #best keywords that are taken from different blogs
    keywords = ['.dll','std::',':dword']
    #Below taken registers are general purpose registers and special register
    #All the registers which are taken are best
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\asmsmallfile.txt","w+")
    files = os.listdir('first')
    for f in files:
        #filling the values with zeros into the arrays
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        # https://docs.python.org/3/library/codecs.html#codecs.ignore_errors
        # https://docs.python.org/3/library/codecs.html#codecs.Codec.encode
        with codecs.open('first/'+f,encoding='cp1252',errors ='replace') as f
            for lines in fli:
                # https://www.tutorialspoint.com/python3/string_rstrip.htm
                line=lines.rstrip().split()
                l=line[0]
                #counting the prefixs in each and every line
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                #counting the opcodes in each and every line
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                #counting registers in the line
                for i in range(len(registers)):
                    for li in line:
                        # we will use registers only in 'text' and 'CODE' segm
                        if registers[i] in li and ('text' in l or 'CODE' in l
                            registerscount[i]+=1
                #counting keywords in the line
                for i in range(len(keywords)):
```

```python
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
            #pushing the values into the file after reading whole file
            for prefix in prefixescount:
                file1.write(str(prefix)+",")
            for opcode in opcodescount:
                file1.write(str(opcode)+",")
            for register in registerscount:
                file1.write(str(register)+",")
            for key in keywordcount:
                file1.write(str(key)+",")
            file1.write("\n")
    file1.close()


#same as above
def secondprocess():
    prefixes = ['HEADER:','.text:','.Pav:','.idata:','.data:','.bss:','.rdata
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
    keywords = ['.dll','std::',':dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\mediumasmfile.txt","w+")
    files = os.listdir('second')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('second/'+f,encoding='cp1252',errors ='replace') as
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                for i in range(len(registers)):
                    for li in line:
                        if registers[i] in li and ('text' in l or 'CODE' in l
                            registerscount[i]+=1
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
        for prefix in prefixescount:
            file1.write(str(prefix)+",")
        for opcode in opcodescount:
            file1.write(str(opcode)+",")
```

```python
114            for register in registerscount:
115                file1.write(str(register)+",")
116            for key in keywordcount:
117                file1.write(str(key)+",")
118            file1.write("\n")
119        file1.close()
120
121    # same as smallprocess() functions
122    def thirdprocess():
123        prefixes = ['HEADER:','.text:','.Pav:','.idata:','.data:','.bss:','.rdata
124        opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
125        keywords = ['.dll','std::',':dword']
126        registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
127        file1=open("output\largeasmfile.txt","w+")
128        files = os.listdir('thrid')
129        for f in files:
130            prefixescount=np.zeros(len(prefixes),dtype=int)
131            opcodescount=np.zeros(len(opcodes),dtype=int)
132            keywordcount=np.zeros(len(keywords),dtype=int)
133            registerscount=np.zeros(len(registers),dtype=int)
134            features=[]
135            f2=f.split('.')[0]
136            file1.write(f2+",")
137            opcodefile.write(f2+" ")
138            with codecs.open('thrid/'+f,encoding='cp1252',errors ='replace') as f
139                for lines in fli:
140                    line=lines.rstrip().split()
141                    l=line[0]
142                    for i in range(len(prefixes)):
143                        if prefixes[i] in line[0]:
144                            prefixescount[i]+=1
145                    line=line[1:]
146                    for i in range(len(opcodes)):
147                        if any(opcodes[i]==li for li in line):
148                            features.append(opcodes[i])
149                            opcodescount[i]+=1
150                    for i in range(len(registers)):
151                        for li in line:
152                            if registers[i] in li and ('text' in l or 'CODE' in l
153                                registerscount[i]+=1
154                    for i in range(len(keywords)):
155                        for li in line:
156                            if keywords[i] in li:
157                                keywordcount[i]+=1
158            for prefix in prefixescount:
159                file1.write(str(prefix)+",")
160            for opcode in opcodescount:
161                file1.write(str(opcode)+",")
162            for register in registerscount:
163                file1.write(str(register)+",")
164            for key in keywordcount:
165                file1.write(str(key)+",")
166            file1.write("\n")
167        file1.close()
168
169
170    def fourthprocess():
```

```python
        prefixes = ['HEADER:','.text:','.Pav:','.idata:','.data:','.bss:','.rdata
        opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
        keywords = ['.dll','std::',':dword']
        registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
        file1=open("output\hugeasmfile.txt","w+")
        files = os.listdir('fourth/')
        for f in files:
            prefixescount=np.zeros(len(prefixes),dtype=int)
            opcodescount=np.zeros(len(opcodes),dtype=int)
            keywordcount=np.zeros(len(keywords),dtype=int)
            registerscount=np.zeros(len(registers),dtype=int)
            features=[]
            f2=f.split('.')[0]
            file1.write(f2+",")
            opcodefile.write(f2+" ")
            with codecs.open('fourth/'+f,encoding='cp1252',errors ='replace') as 
                for lines in fli:
                    line=lines.rstrip().split()
                    l=line[0]
                    for i in range(len(prefixes)):
                        if prefixes[i] in line[0]:
                            prefixescount[i]+=1
                    line=line[1:]
                    for i in range(len(opcodes)):
                        if any(opcodes[i]==li for li in line):
                            features.append(opcodes[i])
                            opcodescount[i]+=1
                    for i in range(len(registers)):
                        for li in line:
                            if registers[i] in li and ('text' in l or 'CODE' in l
                                registerscount[i]+=1
                    for i in range(len(keywords)):
                        for li in line:
                            if keywords[i] in li:
                                keywordcount[i]+=1
            for prefix in prefixescount:
                file1.write(str(prefix)+",")
            for opcode in opcodescount:
                file1.write(str(opcode)+",")
            for register in registerscount:
                file1.write(str(register)+",")
            for key in keywordcount:
                file1.write(str(key)+",")
            file1.write("\n")
        file1.close()


def fifthprocess():
    prefixes = ['HEADER:','.text:','.Pav:','.idata:','.data:','.bss:','.rdata
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
    keywords = ['.dll','std::',':dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\trainasmfile.txt","w+")
    files = os.listdir('fifth/')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
```

```python
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('fifth/'+f,encoding='cp1252',errors ='replace') as f
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                for i in range(len(registers)):
                    for li in line:
                        if registers[i] in li and ('text' in l or 'CODE' in l
                            registerscount[i]+=1
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
        for prefix in prefixescount:
            file1.write(str(prefix)+",")
        for opcode in opcodescount:
            file1.write(str(opcode)+",")
        for register in registerscount:
            file1.write(str(register)+",")
        for key in keywordcount:
            file1.write(str(key)+",")
        file1.write("\n")
    file1.close()


def main():
    #the below code is used for multiprogramming
    #the number of process depends upon the number of cores present System
    #process is used to call multiprogramming
    manager=multiprocessing.Manager()
    p1=Process(target=firstprocess)
    p2=Process(target=secondprocess)
    p3=Process(target=thirdprocess)
    p4=Process(target=fourthprocess)
    p5=Process(target=fifthprocess)
    #p1.start() is used to start the thread execution
    p1.start()
    p2.start()
    p3.start()
    p4.start()
    p5.start()
    #After completion all the threads are joined
    p1.join()
    p2.join()
```

```
285        p3.join()
286        p4.join()
287        p5.join()
288
289    if __name__=="__main__":
290        main()
```

In [137]:
```
1   # asmoutputfile.csv(output genarated from the above two cells) will contain a
2   # this file will be uploaded in the drive, you can directly use this
3   dfasm=pd.read_csv("asmoutputfile.csv")
4   Y.columns = ['ID', 'Class']
5   result_asm = pd.merge(dfasm, Y,on='ID', how='left')
6   result_asm.head()
```

Out[137]:

| | ID | HEADER: | .text: | .Pav: | .idata: | .data: | .bss: | .rdata: | .edata: | .rsrc: | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 01kcPWA9K2BOxQeS5Rju | 19 | 744 | 0 | 127 | 57 | 0 | 323 | 0 | 3 | ... |
| **1** | 1E93CpP60RHFNiT5Qfvn | 17 | 838 | 0 | 103 | 49 | 0 | 0 | 0 | 3 | ... |
| **2** | 3ekVow2ajZHbTnBcsDfX | 17 | 427 | 0 | 50 | 43 | 0 | 145 | 0 | 3 | ... |
| **3** | 3X2nY7iQaPBIWDrAZqJe | 17 | 227 | 0 | 43 | 19 | 0 | 0 | 0 | 3 | ... |
| **4** | 46OZzdsSKDCFV8h7XWxf | 17 | 402 | 0 | 59 | 170 | 0 | 0 | 0 | 3 | ... |

5 rows × 53 columns

**4.2.1.1 Files sizes of each .asm file**

```
In [138]:    1   #file sizes of byte files
             2
             3   files=os.listdir('asmFiles')
             4   filenames=Y['ID'].tolist()
             5   class_y=Y['Class'].tolist()
             6   class_bytes=[]
             7   sizebytes=[]
             8   fnames=[]
             9   for file in files:
            10       # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
            11       # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=356157170
            12       # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=151
            13       # read more about os.stat: here https://www.tutorialspoint.com/python/os_
            14       statinfo=os.stat('asmFiles/'+file)
            15       # split the file name at '.' and take the first part of it i.e the file n
            16       file=file.split('.')[0]
            17       if any(file == filename for filename in filenames):
            18           i=filenames.index(file)
            19           class_bytes.append(class_y[i])
            20           # converting into Mb's
            21           sizebytes.append(statinfo.st_size/(1024.0*1024.0))
            22           fnames.append(file)
            23   asm_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes}
            24   print (asm_size_byte.head())
```

```
     Class                  ID       size
0        9  01azqd4InC7m9JpocGv5  56.229886
1        2  01IsoiSMh5gxyDYTl4CB  13.999378
2        9  01jsnpXSAlgw6aPeDxrU   8.507785
3        1  01kcPWA9K2BOxQeS5Rju   0.078190
4        8  01SuzwMJEIXsK7A8dQbl   0.996723
```

**4.2.1.2 Distribution of .asm file sizes**

```
In [139]:    1  #boxplot of asm files
             2  ax = sns.boxplot(x="Class", y="size", data=asm_size_byte)
             3  plt.title("boxplot of .bytes file sizes")
             4  plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .bytes file sizes

```
In [140]:    1  # add the file size feature to previous extracted features
             2  print(result_asm.shape)
             3  print(asm_size_byte.shape)
             4  result_asm = pd.merge(result_asm, asm_size_byte.drop(['Class'], axis=1),on='I
             5  result_asm.head()
```

(10868, 53)
(10868, 3)

Out[140]:

| | ID | HEADER: | .text: | .Pav: | .idata: | .data: | .bss: | .rdata: | .edata: | .rsrc: | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 01kcPWA9K2BOxQeS5Rju | 19 | 744 | 0 | 127 | 57 | 0 | 323 | 0 | 3 | ... |
| 1 | 1E93CpP60RHFNiT5Qfvn | 17 | 838 | 0 | 103 | 49 | 0 | 0 | 0 | 3 | ... |
| 2 | 3ekVow2ajZHbTnBcsDfX | 17 | 427 | 0 | 50 | 43 | 0 | 145 | 0 | 3 | ... |
| 3 | 3X2nY7iQaPBIWDrAZqJe | 17 | 227 | 0 | 43 | 19 | 0 | 0 | 0 | 3 | ... |
| 4 | 46OZzdsSKDCFV8h7XWxf | 17 | 402 | 0 | 59 | 170 | 0 | 0 | 0 | 3 | ... |

5 rows × 54 columns

```
In [145]:    1  # we normalize the data each column
             2  result_asm = normalize(result_asm)
             3  result_asm.head()
```

Out[145]:

| | ID | HEADER: | .text: | .Pav: | .idata: | .data: | .bss: | .rdata: | .edata |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 01kcPWA9K2BOxQeS5Rju | 0.107345 | 0.001092 | 0.0 | 0.000761 | 0.000023 | 0.0 | 0.000084 | 0.0 |
| 1 | 1E93CpP60RHFNiT5Qfvn | 0.096045 | 0.001230 | 0.0 | 0.000617 | 0.000019 | 0.0 | 0.000000 | 0.0 |
| 2 | 3ekVow2ajZHbTnBcsDfX | 0.096045 | 0.000627 | 0.0 | 0.000300 | 0.000017 | 0.0 | 0.000038 | 0.0 |
| 3 | 3X2nY7iQaPBIWDrAZqJe | 0.096045 | 0.000333 | 0.0 | 0.000258 | 0.000008 | 0.0 | 0.000000 | 0.0 |
| 4 | 46OZzdsSKDCFV8h7XWxf | 0.096045 | 0.000590 | 0.0 | 0.000353 | 0.000068 | 0.0 | 0.000000 | 0.0 |

5 rows × 54 columns

## 4.2.2 Univariate analysis on asm file features

```
In [146]:    1  ax = sns.boxplot(x="Class", y=".text:", data=result_asm)
             2  plt.title("boxplot of .asm text segment")
             3  plt.show()
```
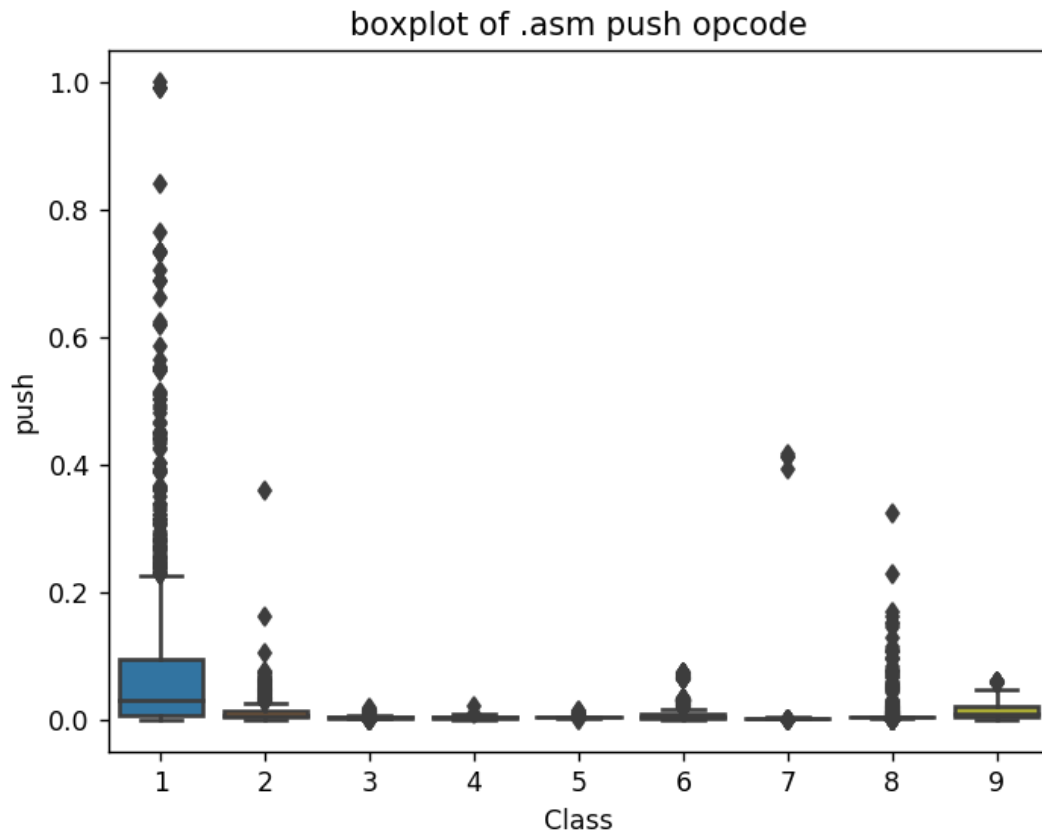
<IPython.core.display.Javascript object>

The plot is between Text and class

Class 1,2 and 9 can be easly separated

In [115]:
```python
ax = sns.boxplot(x="Class", y=".Pav:", data=result_asm)
plt.title("boxplot of .asm pav segment")
plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .asm pav segment

```
In [116]:   1  ax = sns.boxplot(x="Class", y=".data:", data=result_asm)
            2  plt.title("boxplot of .asm data segment")
            3  plt.show()
```
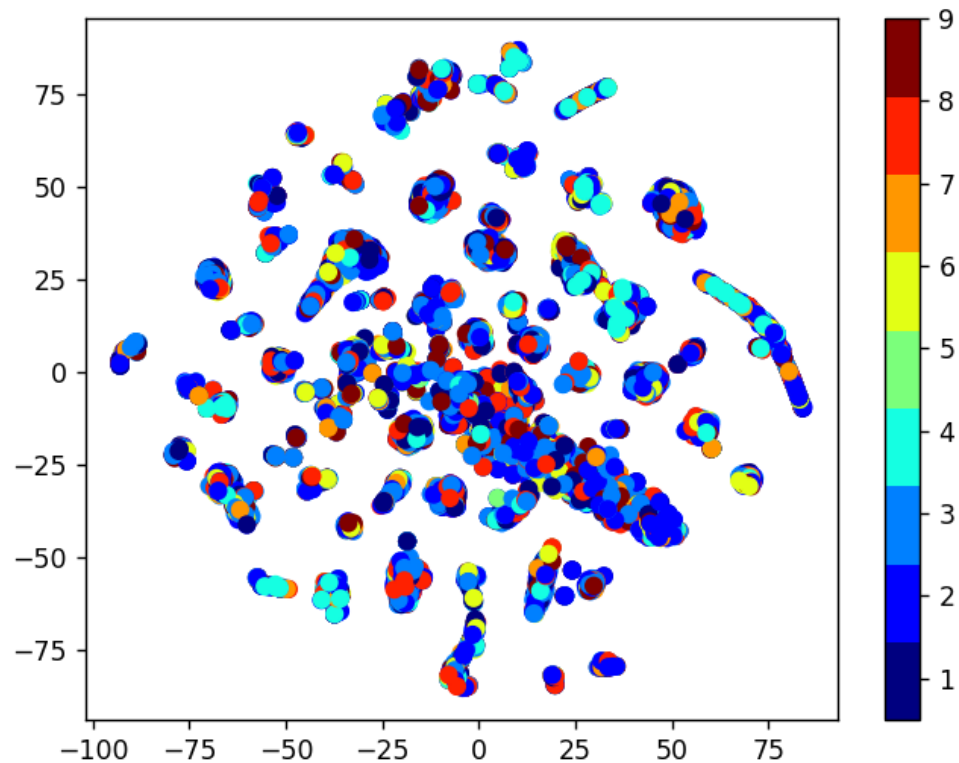
<IPython.core.display.Javascript object>

boxplot of .asm data segment



The plot is between data segment and class label
class 6 and class 9 can be easily separated from given points

```
1  ax = sns.boxplot(x="Class", y=".bss:", data=result_asm)
2  plt.title("boxplot of .asm bss segment")
3  plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .asm bss segment

plot between bss segment and class label
very less number of files are having bss segment

```
In [118]:    1  ax = sns.boxplot(x="Class", y=".rdata:", data=result_asm)
             2  plt.title("boxplot of .asm rdata segment")
             3  plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .asm rdata segment

Plot between rdata segment and Class segment
Class 2 can be easily separated 75 pecentile files are having 1M rdata li
nes

```
1  ax = sns.boxplot(x="Class", y="jmp", data=result_asm)
2  plt.title("boxplot of .asm jmp opcode")
3  plt.show()
```

<IPython.core.display.Javascript object>

boxplot of .asm jmp opcode



plot between jmp and Class label
Class 1 is having frequency of 2000 approx in 75 perentile of files

```
1  ax = sns.boxplot(x="Class", y="mov", data=result_asm)
2  plt.title("boxplot of .asm mov opcode")
3  plt.show()
```

<IPython.core.display.Javascript object>

boxplot of .asm mov opcode



plot between Class label and mov opcode
Class 1 is having frequency of 2000 approx in 75 perentile of files

```
In [121]:  1  ax = sns.boxplot(x="Class", y="retf", data=result_asm)
           2  plt.title("boxplot of .asm retf opcode")
           3  plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .asm retf opcode

plot between Class label and retf
Class 6 can be easily separated with opcode retf
The frequency of retf is approx of 250.

```
1  ax = sns.boxplot(x="Class", y="push", data=result_asm)
2  plt.title("boxplot of .asm push opcode")
3  plt.show()
```

<IPython.core.display.Javascript object>



boxplot of .asm push opcode

plot between push opcode and Class label
Class 1 is having 75 precentile files with push opcodes of frequency 1000

### 4.2.2 Multivariate Analysis on .asm file features

```
In [129]:   1   # check out the course content for more explantion on tsne algorithm
            2   # https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/t-d
            3
            4   #multivariate analysis on byte files
            5   #this is with perplexity 50
            6   xtsne=TSNE(perplexity=50)
            7   results=xtsne.fit_transform(result_asm.drop(['ID','Class'], axis=1).fillna(0)
            8   vis_x = results[:, 0]
            9   vis_y = results[:, 1   ]
           10   plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
           11   plt.colorbar(ticks=range(10))
           12   plt.clim(0.5, 9)
           13   plt.show()
```

<IPython.core.display.Javascript object>

```
1  # by univariate analysis on the .asm file features we are getting very neglig
2  # 'rtn', '.BSS:' '.CODE' features, so heare we are trying multivariate analys
3  # the plot looks very messy
4
5  xtsne=TSNE(perplexity=30)
6  results=xtsne.fit_transform(result_asm.drop(['ID','Class', 'rtn', '.BSS:', '.
7  vis_x = results[:, 0]
8  vis_y = results[:, 1]
9  plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
10 plt.colorbar(ticks=range(10))
11 plt.clim(0.5, 9)
12 plt.show()
```

<IPython.core.display.Javascript object>



TSNE for asm data with perplexity 50

## 4.2.3 Conclusion on EDA

- We have taken only 52 features from asm files (after reading through many blogs and research papers)
- The univariate analysis was done only on few important features.
- Take-aways

- 1. Class 3 can be easily separated because of the frequency of segments,opcodes and keywords being less
- 2. Each feature has its unique importance in separating the Class labels.

## 4.3 Train and test split

```
In [149]:   1  asm_y = result_asm['Class']
            2  asm_x = result_asm.drop(['ID','Class','.BSS:','rtn','.CODE'], axis=1)
```

```
In [150]:   1  X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x,asm
            2  X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm,
```

```
In [153]:  1  print( X_cv_asm.isnull().all())
```

```
HEADER:    False
.text:     False
.Pav:      False
.idata:    False
.data:     False
.bss:      False
.rdata:    False
.edata:    False
.rsrc:     False
.tls:      False
.reloc:    False
jmp        False
mov        False
retf       False
push       False
pop        False
xor        False
retn       False
nop        False
sub        False
inc        False
dec        False
add        False
imul       False
xchg       False
or         False
shr        False
cmp        False
call       False
shl        False
ror        False
rol        False
jnb        False
jz         False
lea        False
movzx      False
.dll       False
std::      False
:dword     False
edx        False
esi        False
eax        False
ebx        False
ecx        False
edi        False
ebp        False
esp        False
eip        False
size       False
dtype: bool
```

## 4.4. Machine Learning models on features of .asm files

### 4.4.1 K-Nearest Neigbors

```python
# find more about KNeighborsClassifier() here http://scikit-learn.org/stable/
# ---------------------------
# default parameter
# KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', le
# metric='minkowski', metric_params=None, n_jobs=1, **kwargs)

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


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

alpha = [x for x in range(1, 21,2)]
cv_log_error_array=[]
for i in alpha:
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=k_cfl.clas

for i in range(len(cv_log_error_array)):
    print ('log_loss for k = ',alpha[i],'is',cv_log_error_array[i])

best_alpha = np.argmin(cv_log_error_array)

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()

k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
k_cfl.fit(X_train_asm,y_train_asm)
```

```
57  sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
58  sig_clf.fit(X_train_asm, y_train_asm)
59  pred_y=sig_clf.predict(X_test_asm)
60
61
62  predict_y = sig_clf.predict_proba(X_train_asm)
63  print ('log loss for train data',log_loss(y_train_asm, predict_y))
64  predict_y = sig_clf.predict_proba(X_cv_asm)
65  print ('log loss for cv data',log_loss(y_cv_asm, predict_y))
66  predict_y = sig_clf.predict_proba(X_test_asm)
67  print ('log loss for test data',log_loss(y_test_asm, predict_y))
68  plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```

```
log_loss for k =  1 is 0.104531321344
log_loss for k =  3 is 0.0958800580948
log_loss for k =  5 is 0.0995466557335
log_loss for k =  7 is 0.107227274345
log_loss for k =  9 is 0.119239543547
log_loss for k =  11 is 0.133926642781
log_loss for k =  13 is 0.147643793967
log_loss for k =  15 is 0.159439699615
log_loss for k =  17 is 0.16878376444
log_loss for k =  19 is 0.178020728839

<IPython.core.display.Javascript object>
```

Cross Validation Error for each alpha



```
log loss for train data 0.0476773462198
log loss for cv data 0.0958800580948
log loss for test data 0.0894810720832
```

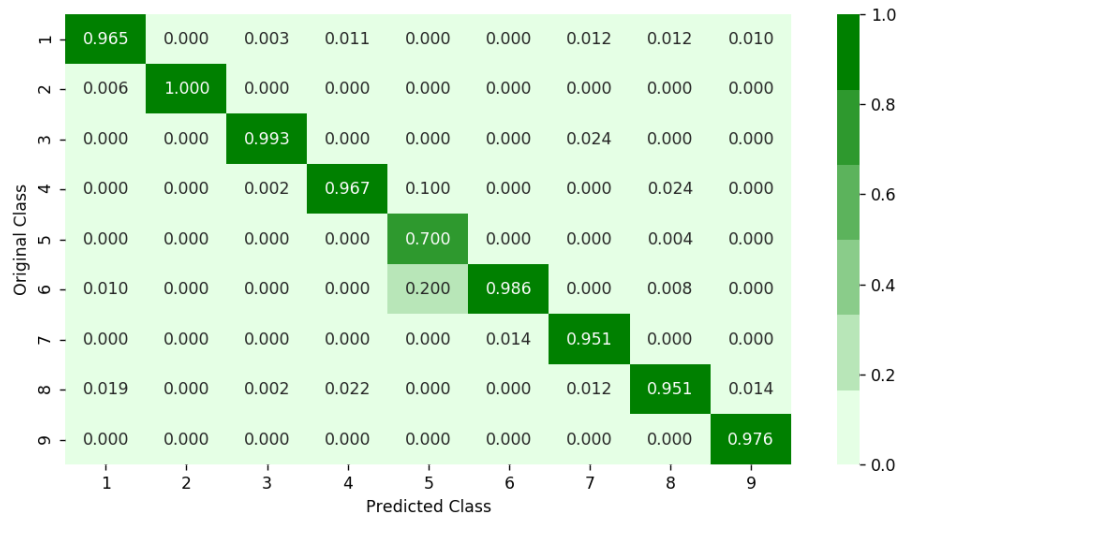Number of misclassified points  2.02391904324
---------------------------------------------------- Confusion matrix --------
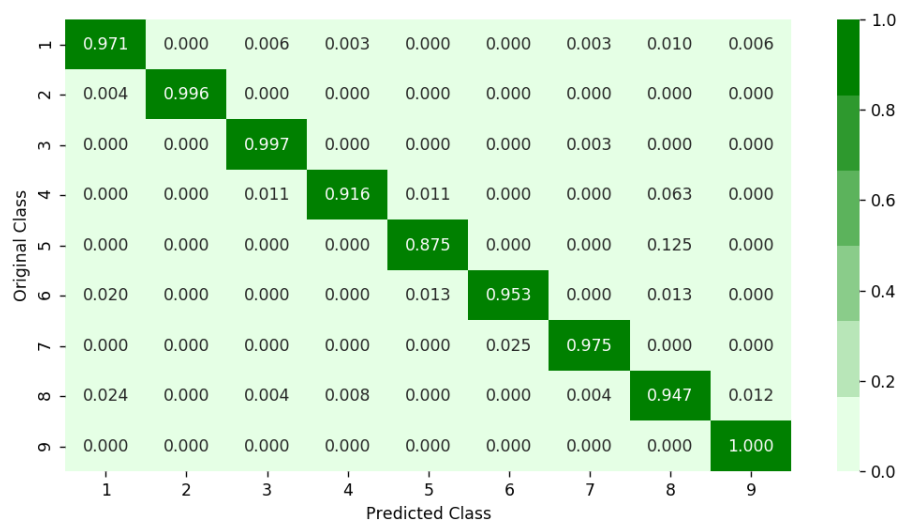-------------------------------------------

<IPython.core.display.Javascript object>



---------------------------------------------------- Precision matrix ----------
---------------------------------------

<IPython.core.display.Javascript object>



Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
---------------------------------------------------- Recall matrix -------------
-------------------------------------

<IPython.core.display.Javascript object>

Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]

## 4.4.2 Logistic Regression

```python
In [160]:
 1   # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/g
 2   # ------------------------------
 3   # default parameters
 4   # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_
 5   # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning
 6   # class_weight=None, warm_start=False, average=False, n_iter=None)
 7
 8   # some of methods
 9   # fit(X, y[, coef_init, intercept_init, …]) Fit linear model with Stochastic
10   # predict(X)    Predict class labels for samples in X.
11
12   #------------------------------
13   # video link: https://www.appliedaicourse.com/course/applied-ai-course-online
14   #------------------------------
15
16
17   alpha = [10 ** x for x in range(-5, 4)]
18   cv_log_error_array=[]
19   for i in alpha:
20       logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
21       logisticR.fit(X_train_asm,y_train_asm)
22       sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
23       sig_clf.fit(X_train_asm, y_train_asm)
24       predict_y = sig_clf.predict_proba(X_cv_asm)
25       cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.
26
27   for i in range(len(cv_log_error_array)):
28       print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
29
30   best_alpha = np.argmin(cv_log_error_array)
31
32   fig, ax = plt.subplots()
33   ax.plot(alpha, cv_log_error_array,c='g')
34   for i, txt in enumerate(np.round(cv_log_error_array,3)):
35       ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
36   plt.grid()
37   plt.title("Cross Validation Error for each alpha")
38   plt.xlabel("Alpha i's")
39   plt.ylabel("Error measure")
40   plt.show()
41
42   logisticR=LogisticRegression(penalty='l2',C=alpha[best_alpha],class_weight='b
43   logisticR.fit(X_train_asm,y_train_asm)
44   sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
45   sig_clf.fit(X_train_asm, y_train_asm)
46
47   predict_y = sig_clf.predict_proba(X_train_asm)
48   print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=log
49   predict_y = sig_clf.predict_proba(X_cv_asm)
50   print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=logisticR
51   predict_y = sig_clf.predict_proba(X_test_asm)
52   print ('log loss for test data',(log_loss(y_test_asm, predict_y, labels=logis
53   plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```
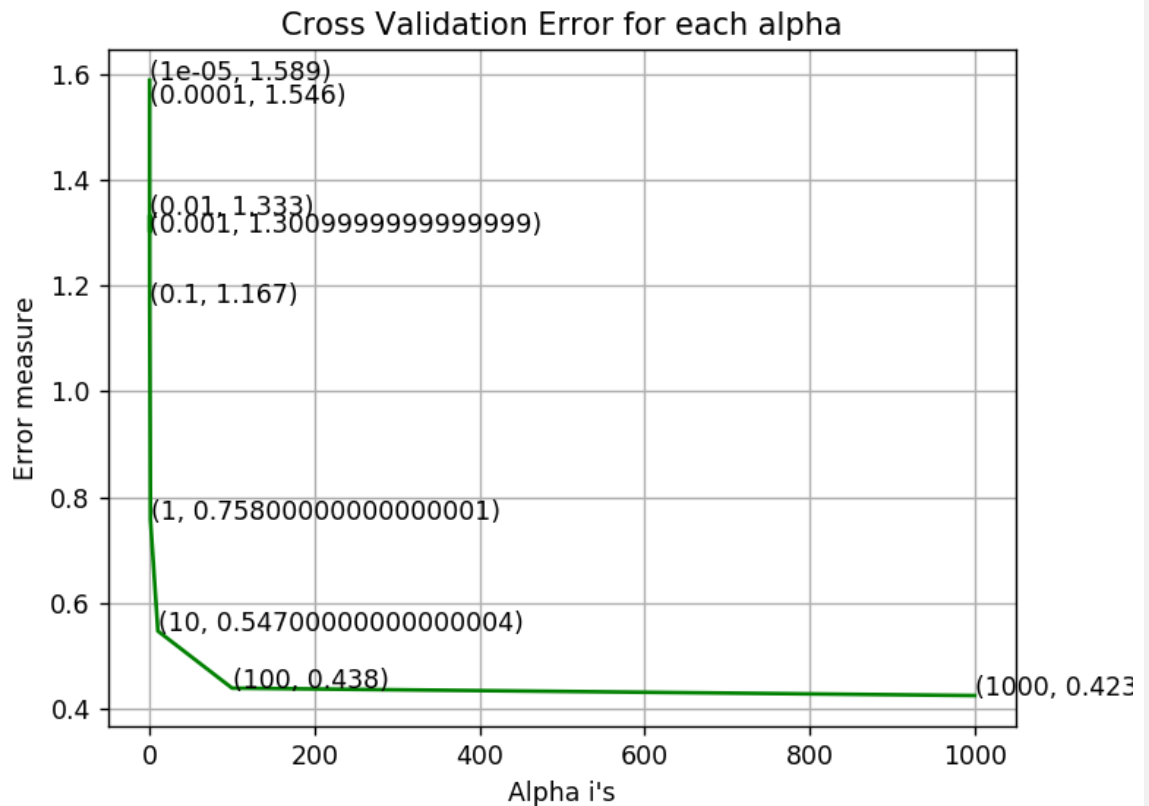
```
log_loss for c =  1e-05 is 1.58867274165
```

```
log_loss for c =   0.0001 is 1.54560797884
log_loss for c =   0.001 is 1.30137786807
log_loss for c =   0.01 is 1.33317456931
log_loss for c =   0.1 is 1.16705751378
log_loss for c =   1 is 0.757667807779
log_loss for c =   10 is 0.546533939819
log_loss for c =   100 is 0.438414998062
log_loss for c =   1000 is 0.424423536526
```
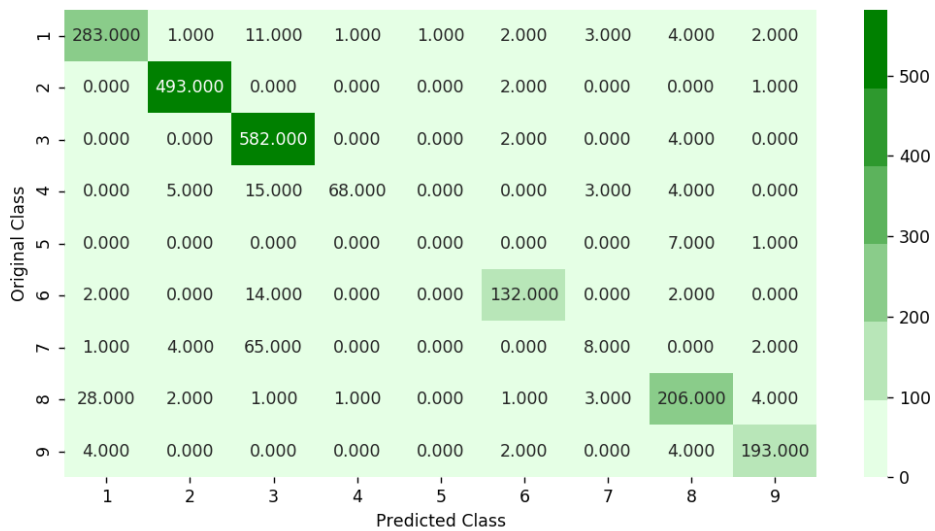
&lt;IPython.core.display.Javascript object&gt;

**Cross Validation Error for each alpha**

(1e-05, 1.589)
(0.0001, 1.546)

(0.01, 1.333)
(0.001, 1.3009999999999999)

(0.1, 1.167)

(1, 0.75800000000000001)

(10, 0.54700000000000004)

(100, 0.438)

(1000, 0.423

Error measure (y-axis)
Alpha i's (x-axis): 0, 200, 400, 600, 800, 1000

```
log loss for train data 0.396219394701
log loss for cv data 0.424423536526
log loss for test data 0.415685592517
Number of misclassified points  9.61361545538
------------------------------------------------ Confusion matrix ----------
----------------------------------------
```
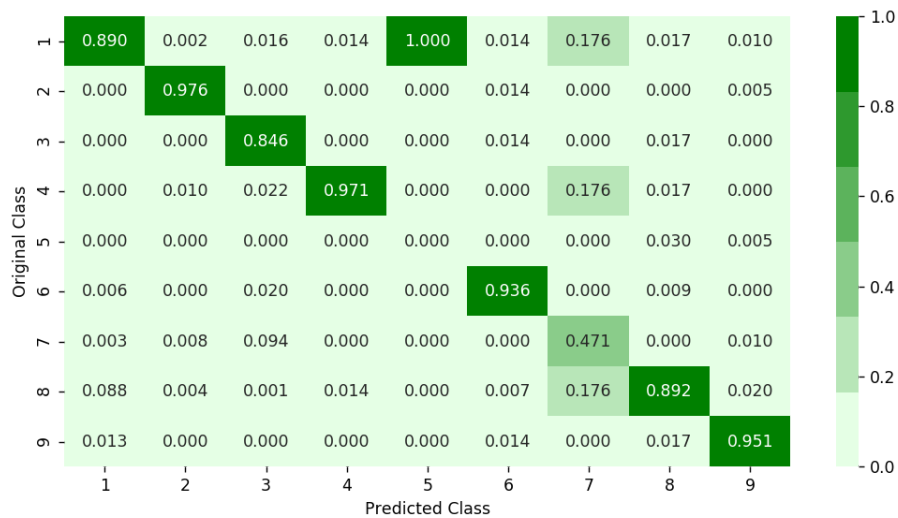
&lt;IPython.core.display.Javascript object&gt;

---------------------------------------------------- Precision matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>



Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
---------------------------------------------------- Recall matrix ------------
--------------------------------------

<IPython.core.display.Javascript object>

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.919 | 0.003 | 0.036 | 0.003 | 0.003 | 0.006 | 0.010 | 0.013 | 0.006 |
| 2 | 0.000 | 0.994 | 0.000 | 0.000 | 0.000 | 0.004 | 0.000 | 0.000 | 0.002 |
| 3 | 0.000 | 0.000 | 0.990 | 0.000 | 0.000 | 0.003 | 0.000 | 0.007 | 0.000 |
| 4 | 0.000 | 0.053 | 0.158 | 0.716 | 0.000 | 0.000 | 0.032 | 0.042 | 0.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.875 | 0.125 |
| 6 | 0.013 | 0.000 | 0.093 | 0.000 | 0.000 | 0.880 | 0.000 | 0.013 | 0.000 |
| 7 | 0.013 | 0.050 | 0.812 | 0.000 | 0.000 | 0.000 | 0.100 | 0.000 | 0.025 |
| 8 | 0.114 | 0.008 | 0.004 | 0.004 | 0.000 | 0.004 | 0.012 | 0.837 | 0.016 |
| 9 | 0.020 | 0.000 | 0.000 | 0.000 | 0.000 | 0.010 | 0.000 | 0.020 | 0.951 |

Original Class / Predicted Class

```
Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
```

### 4.4.3 Random Forest Classifier

```python
# --------------------------------
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini',
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_
# class_weight=None)

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

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

# --------------------------------
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online
# --------------------------------

alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.clas

for i in range(len(cv_log_error_array)):
    print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])


best_alpha = np.argmin(cv_log_error_array)

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()

r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=sig
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=sig_clf.c
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',(log_loss(y_test_asm, predict_y, labels=sig_c
```

```
57  plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
◀ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬                                          ▶
```

```
log_loss for c =   10 is 0.0581657906023
log_loss for c =   50 is 0.0515443148419
log_loss for c =  100 is 0.0513084973231
log_loss for c =  500 is 0.0499021761479
log_loss for c = 1000 is 0.0497972474298
log_loss for c = 2000 is 0.0497091690815
log_loss for c = 3000 is 0.0496706817633
```
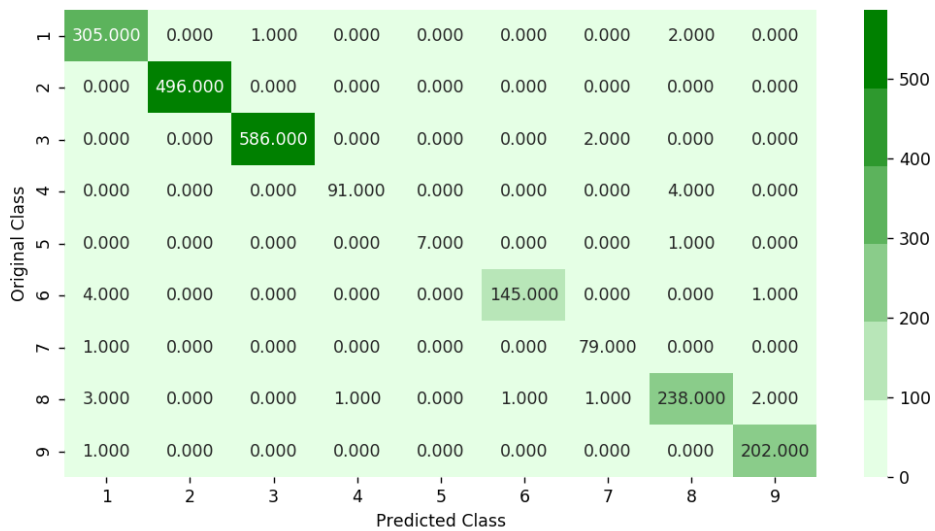
<IPython.core.display.Javascript object>



Cross Validation Error for each alpha

```
log loss for train data 0.0116517052676
log loss for cv data 0.0496706817633
log loss for test data 0.0571239496453
Number of misclassified points  1.14995400184
------------------------------------------------ Confusion matrix ----------
--------------------------------------
```
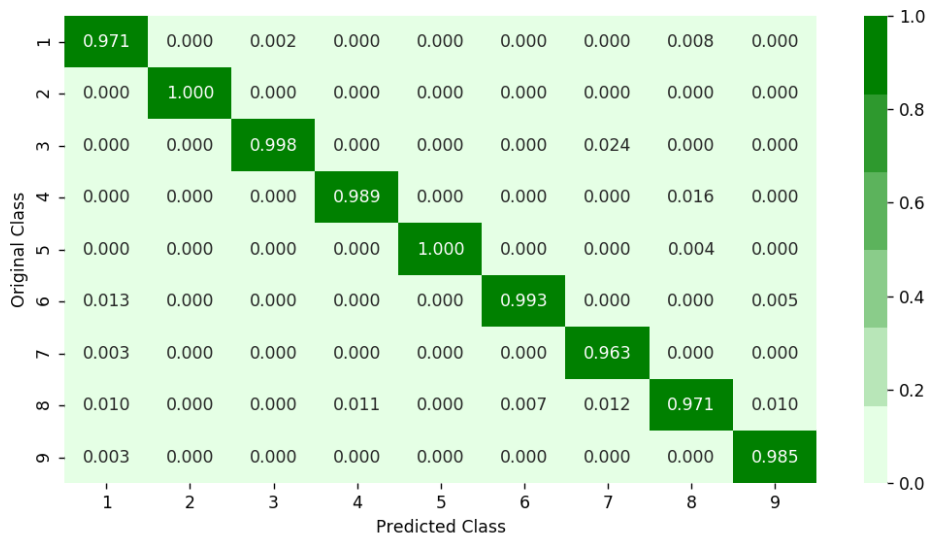
<IPython.core.display.Javascript object>

---------------------------------------------------- Precision matrix ----------
---------------------------------------

<IPython.core.display.Javascript object>



Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
---------------------------------------------------- Recall matrix ------------
---------------------------------------

<IPython.core.display.Javascript object>

Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]

## 4.4.4 XgBoost Classifier

```
In [162]:  1   # Training a hyper-parameter tuned Xg-Boost regressor on our train data
           2
           3   # find more about XGBClassifier function here http://xgboost.readthedocs.io/e
           4   # ------------------------
           5   # default paramters
           6   # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=10
           7   # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamm
           8   # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg
           9   # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None
          10
          11   # some of methods of RandomForestRegressor()
          12   # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppi
          13   # get_params([deep])    Get parameters for this estimator.
          14   # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE
          15   # get_score(importance_type='weight') -> get the feature importance
          16   # -----------------------
          17   # video link2: https://www.appliedaicourse.com/course/applied-ai-course-onlin
          18   # ----------------------
          19
          20   alpha=[10,50,100,500,1000,2000,3000]
          21   cv_log_error_array=[]
          22   for i in alpha:
          23       x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
          24       x_cfl.fit(X_train_asm,y_train_asm)
          25       sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
          26       sig_clf.fit(X_train_asm, y_train_asm)
          27       predict_y = sig_clf.predict_proba(X_cv_asm)
          28       cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=x_cfl.clas
          29
          30   for i in range(len(cv_log_error_array)):
          31       print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
          32
          33
          34   best_alpha = np.argmin(cv_log_error_array)
          35
          36   fig, ax = plt.subplots()
          37   ax.plot(alpha, cv_log_error_array,c='g')
          38   for i, txt in enumerate(np.round(cv_log_error_array,3)):
          39       ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
          40   plt.grid()
          41   plt.title("Cross Validation Error for each alpha")
          42   plt.xlabel("Alpha i's")
          43   plt.ylabel("Error measure")
          44   plt.show()
          45
          46   x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
          47   x_cfl.fit(X_train_asm,y_train_asm)
          48   sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
          49   sig_clf.fit(X_train_asm, y_train_asm)
          50
          51   predict_y = sig_clf.predict_proba(X_train_asm)
          52
          53   print ('For values of best alpha = ', alpha[best_alpha], "The train log loss
          54   predict_y = sig_clf.predict_proba(X_cv_asm)
          55   print('For values of best alpha = ', alpha[best_alpha], "The cross validation
          56   predict_y = sig_clf.predict_proba(X_test_asm)
```

```
57  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
58  plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```
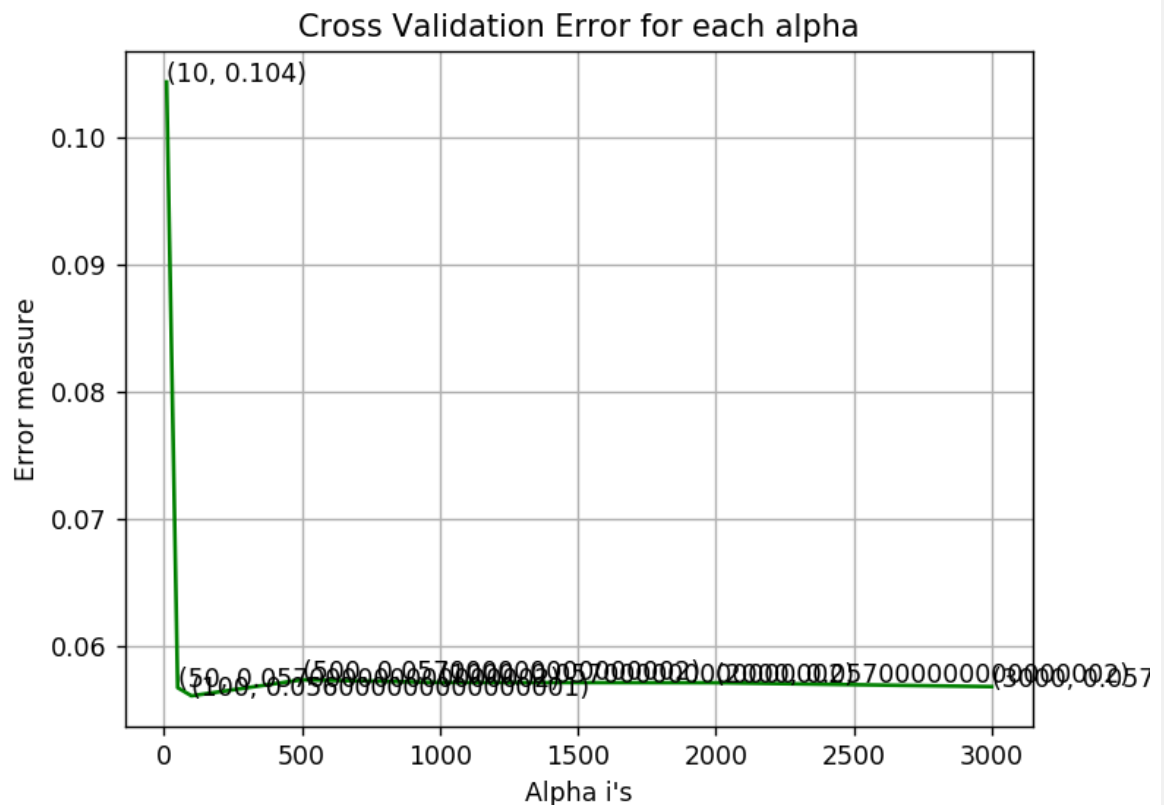
```
log_loss for c =   10 is 0.104344888454
log_loss for c =   50 is 0.0567190635611
log_loss for c =  100 is 0.056075038646
log_loss for c =  500 is 0.057336051683
log_loss for c =  1000 is 0.0571265109903
log_loss for c =  2000 is 0.057103406781
log_loss for c =  3000 is 0.0567993215778
```

<IPython.core.display.Javascript object>



Cross Validation Error for each alpha

For values of best alpha =   100 The train log loss is: 0.0117883742574
For values of best alpha =   100 The cross validation log loss is: 0.0560750386
46
For values of best alpha =   100 The test log loss is: 0.0491647763845
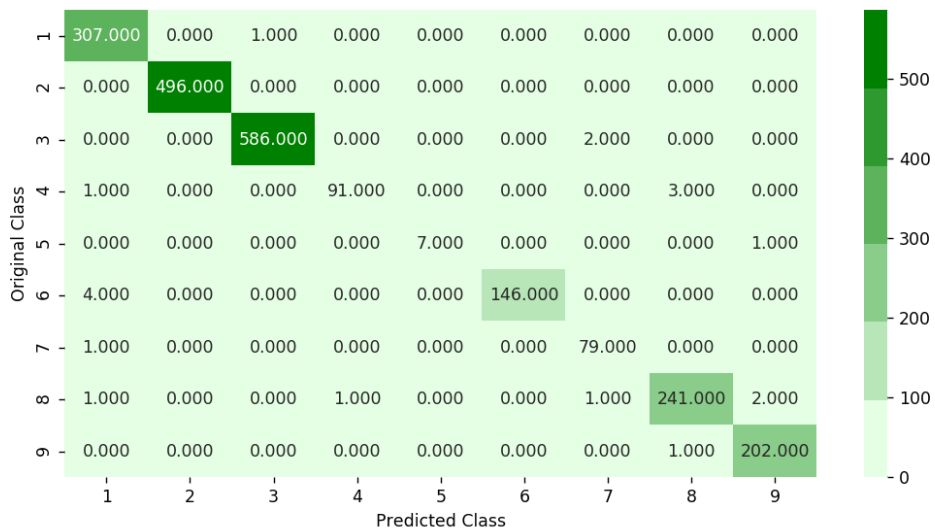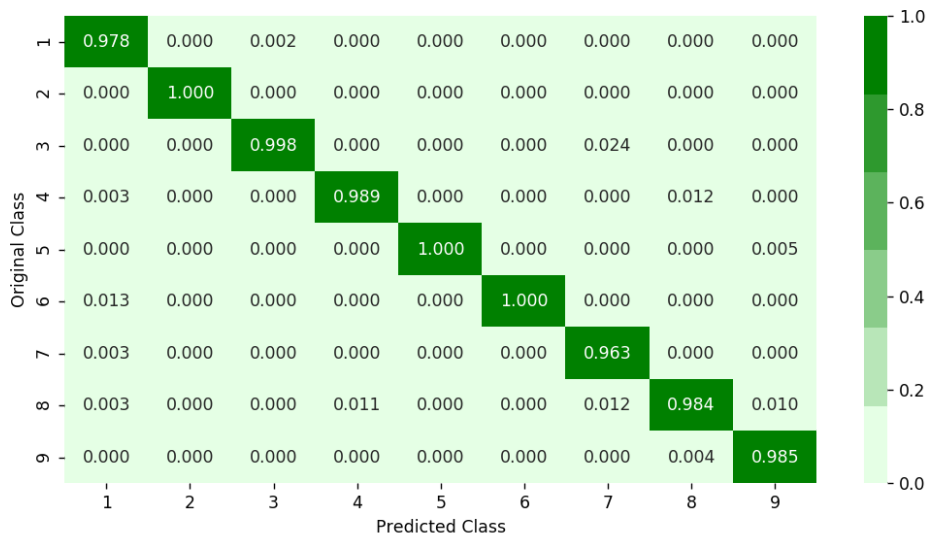Number of misclassified points  0.873965041398
------------------------------------------------ Confusion matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>

---------------------------------------------------- Precision matrix ----------
----------------------------------------

<IPython.core.display.Javascript object>



Sum of columns in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]
---------------------------------------------------- Recall matrix ------------
-----------------------------------

<IPython.core.display.Javascript object>

Sum of rows in precision matrix [ 1.  1.  1.  1.  1.  1.  1.  1.  1.]

## 4.4.5 Xgboost Classifier with best hyperparameters

```
In [163]:  1  x_cfl=XGBClassifier()
           2
           3  prams={
           4      'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
           5      'n_estimators':[100,200,500,1000,2000],
           6      'max_depth':[3,5,10],
           7      'colsample_bytree':[0.1,0.3,0.5,1],
           8      'subsample':[0.1,0.3,0.5,1]
           9  }
          10  random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jo
          11  random_cfl.fit(X_train_asm,y_train_asm)
```

Fitting 3 folds for each of 10 candidates, totalling 30 fits

```
[Parallel(n_jobs=-1)]: Done    2 tasks       | elapsed:    8.1s
[Parallel(n_jobs=-1)]: Done    9 tasks       | elapsed:   32.8s
[Parallel(n_jobs=-1)]: Done   19 out of  30 | elapsed:  1.1min remaining:   39.
3s
[Parallel(n_jobs=-1)]: Done   23 out of  30 | elapsed:  1.3min remaining:   23.
0s
[Parallel(n_jobs=-1)]: Done   27 out of  30 | elapsed:  1.4min remaining:    9.
2s
[Parallel(n_jobs=-1)]: Done   30 out of  30 | elapsed:  2.3min finished
```

Out[163]: RandomizedSearchCV(cv=None, error_score='raise',
          estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, colsampl
          e_bytree=1,
             gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
             min_child_weight=1, missing=None, n_estimators=100, nthread=-1,
             objective='binary:logistic', reg_alpha=0, reg_lambda=1,
             scale_pos_weight=1, seed=0, silent=True, subsample=1),
          fit_params=None, iid=True, n_iter=10, n_jobs=-1,
          param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15,
          0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'co
          lsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
          pre_dispatch='2*n_jobs', random_state=None, refit=True,
          return_train_score=True, scoring=None, verbose=10)

```
In [164]:  1  print (random_cfl.best_params_)
```

{'subsample': 1, 'n_estimators': 200, 'max_depth': 5, 'learning_rate': 0.15,
'colsample_bytree': 0.5}

```
In [170]:  1  # Training a hyper-parameter tuned Xg-Boost regressor on our train data
           2
           3  # find more about XGBClassifier function here http://xgboost.readthedocs.io/e
           4  # ------------------------
           5  # default paramters
           6  # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=10
           7  # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamm
           8  # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg
           9  # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None
          10
          11  # some of methods of RandomForestRegressor()
          12  # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppi
          13  # get_params([deep])    Get parameters for this estimator.
          14  # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE
          15  # get_score(importance_type='weight') -> get the feature importance
          16  # ----------------------
          17  # video link2: https://www.appliedaicourse.com/course/applied-ai-course-onlin
          18  # ----------------------
          19
          20  x_cfl=XGBClassifier(n_estimators=200,subsample=0.5,learning_rate=0.15,colsamp
          21  x_cfl.fit(X_train_asm,y_train_asm)
          22  c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
          23  c_cfl.fit(X_train_asm,y_train_asm)
          24
          25  predict_y = c_cfl.predict_proba(X_train_asm)
          26  print ('train loss',log_loss(y_train_asm, predict_y))
          27  predict_y = c_cfl.predict_proba(X_cv_asm)
          28  print ('cv loss',log_loss(y_cv_asm, predict_y))
          29  predict_y = c_cfl.predict_proba(X_test_asm)
          30  print ('test loss',log_loss(y_test_asm, predict_y))
```

```
train loss 0.0102661325822
cv loss 0.0501201796687
test loss 0.0483908764397
```

## 4.5. Machine Learning models on features of both .asm and .bytes files
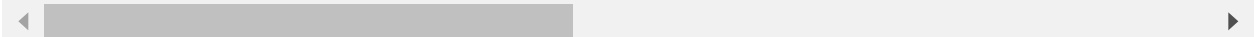
### 4.5.1. Merging both asm and byte file features

```
In [171]:   1  result.head()
```

Out[171]:

|   | ID | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|----|---|---|---|---|---|---|---|
| 0 | 01azqd4InC7m9JpocGv5 | 0.262806 | 0.005498 | 0.001567 | 0.002067 | 0.002048 | 0.001835 | 0.002058 |
| 1 | 01IsoiSMh5gxyDYTl4CB | 0.017358 | 0.011737 | 0.004033 | 0.003876 | 0.005303 | 0.003873 | 0.004747 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 0.040827 | 0.013434 | 0.001429 | 0.001315 | 0.005464 | 0.005280 | 0.005078 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 0.009209 | 0.001708 | 0.000404 | 0.000441 | 0.000770 | 0.000354 | 0.000310 |
| 4 | 01SuzwMJEIXsK7A8dQbI | 0.008629 | 0.001000 | 0.000168 | 0.000234 | 0.000342 | 0.000232 | 0.000148 |

5 rows × 260 columns

```
In [174]:   1  result_asm.head()
```

Out[174]:

|   | ID | HEADER: | .text: | .Pav: | .idata: | .data: | .bss: | .rdata: | .edata |
|---|----|---------|--------|-------|---------|--------|-------|---------|--------|
| 0 | 01kcPWA9K2BOxQeS5Rju | 0.107345 | 0.001092 | 0.0 | 0.000761 | 0.000023 | 0.0 | 0.000084 | 0.0 |
| 1 | 1E93CpP60RHFNiT5Qfvn | 0.096045 | 0.001230 | 0.0 | 0.000617 | 0.000019 | 0.0 | 0.000000 | 0.0 |
| 2 | 3ekVow2ajZHbTnBcsDfX | 0.096045 | 0.000627 | 0.0 | 0.000300 | 0.000017 | 0.0 | 0.000038 | 0.0 |
| 3 | 3X2nY7iQaPBIWDrAZqJe | 0.096045 | 0.000333 | 0.0 | 0.000258 | 0.000008 | 0.0 | 0.000000 | 0.0 |
| 4 | 46OZzdsSKDCFV8h7XWxf | 0.096045 | 0.000590 | 0.0 | 0.000353 | 0.000068 | 0.0 | 0.000000 | 0.0 |

5 rows × 54 columns

```
In [173]:   1  print(result.shape)
            2  print(result_asm.shape)
```

```
(10868, 260)
(10868, 54)
```

```
In [182]:   1  result_x = pd.merge(result,result_asm.drop(['Class'], axis=1),on='ID', how='l
            2  result_y = result_x['Class']
            3  result_x = result_x.drop(['ID','rtn','.BSS:','.CODE','Class'], axis=1)
            4  result_x.head()
```

Out[182]:

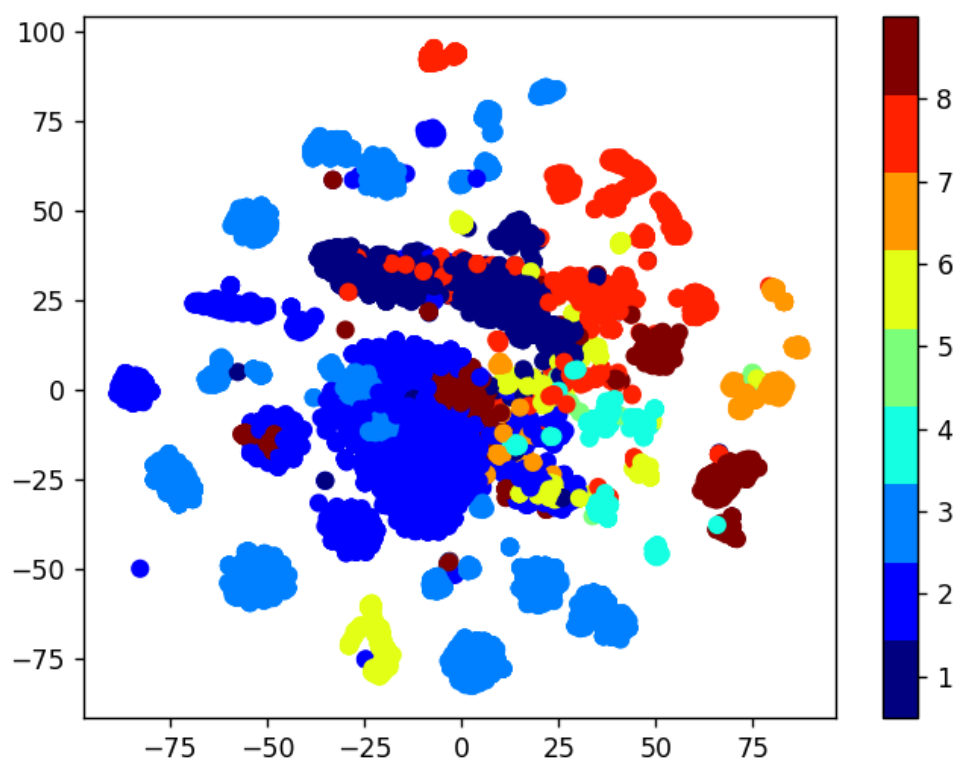|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.262806 | 0.005498 | 0.001567 | 0.002067 | 0.002048 | 0.001835 | 0.002058 | 0.002946 | 0.002638 | 0.003 |
| 1 | 0.017358 | 0.011737 | 0.004033 | 0.003876 | 0.005303 | 0.003873 | 0.004747 | 0.006984 | 0.008267 | 0.000 |
| 2 | 0.040827 | 0.013434 | 0.001429 | 0.001315 | 0.005464 | 0.005280 | 0.005078 | 0.002155 | 0.008104 | 0.002 |
| 3 | 0.009209 | 0.001708 | 0.000404 | 0.000441 | 0.000770 | 0.000354 | 0.000310 | 0.000481 | 0.000959 | 0.000 |
| 4 | 0.008629 | 0.001000 | 0.000168 | 0.000234 | 0.000342 | 0.000232 | 0.000148 | 0.000229 | 0.000376 | 0.000 |

5 rows × 307 columns

### 4.5.2. Multivariate Analysis on final fearures

```
In [181]:   1  xtsne=TSNE(perplexity=50)
            2  results=xtsne.fit_transform(result_x, axis=1))
            3  vis_x = results[:, 0]
            4  vis_y = results[:, 1]
            5  plt.scatter(vis_x, vis_y, c=result_y, cmap=plt.cm.get_cmap("jet", 9))
            6  plt.colorbar(ticks=range(9))
            7  plt.clim(0.5, 9)
            8  plt.show()
```

<IPython.core.display.Javascript object>



### 4.5.3. Train and Test split

```
In [183]:   1  X_train, X_test_merge, y_train, y_test_merge = train_test_split(result_x, res
            2  X_train_merge, X_cv_merge, y_train_merge, y_cv_merge = train_test_split(X_tra
```

### 4.5.4. Random Forest Classifier on final features

```python
In [185]:
     1  # --------------------------------
     2  # default parameters
     3  # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini',
     4  # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_
     5  # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_
     6  # class_weight=None)
     7
     8  # Some of methods of RandomForestClassifier()
     9  # fit(X, y, [sample_weight])     Fit the SVM model according to the given trai
    10  # predict(X)    Perform classification on samples in X.
    11  # predict_proba (X) Perform classification on samples in X.
    12
    13  # some of attributes of  RandomForestClassifier()
    14  # feature_importances_ : array of shape = [n_features]
    15  # The feature importances (the higher, the more important the feature).
    16
    17  # --------------------------------
    18  # video link: https://www.appliedaicourse.com/course/applied-ai-course-online
    19  # --------------------------------
    20
    21  alpha=[10,50,100,500,1000,2000,3000]
    22  cv_log_error_array=[]
    23  from sklearn.ensemble import RandomForestClassifier
    24  for i in alpha:
    25      r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    26      r_cfl.fit(X_train_merge,y_train_merge)
    27      sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    28      sig_clf.fit(X_train_merge, y_train_merge)
    29      predict_y = sig_clf.predict_proba(X_cv_merge)
    30      cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=r_cfl.cl
    31
    32  for i in range(len(cv_log_error_array)):
    33      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
    34
    35
    36  best_alpha = np.argmin(cv_log_error_array)
    37
    38  fig, ax = plt.subplots()
    39  ax.plot(alpha, cv_log_error_array,c='g')
    40  for i, txt in enumerate(np.round(cv_log_error_array,3)):
    41      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
    42  plt.grid()
    43  plt.title("Cross Validation Error for each alpha")
    44  plt.xlabel("Alpha i's")
    45  plt.ylabel("Error measure")
    46  plt.show()
    47
    48
    49  r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n
    50  r_cfl.fit(X_train_merge,y_train_merge)
    51  sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    52  sig_clf.fit(X_train_merge, y_train_merge)
    53
    54  predict_y = sig_clf.predict_proba(X_train_merge)
    55  print ('For values of best alpha = ', alpha[best_alpha], "The train log loss
    56  predict_y = sig_clf.predict_proba(X_cv_merge)
```

```
57  print('For values of best alpha = ', alpha[best_alpha], "The cross validation
58  predict_y = sig_clf.predict_proba(X_test_merge)
59  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
```
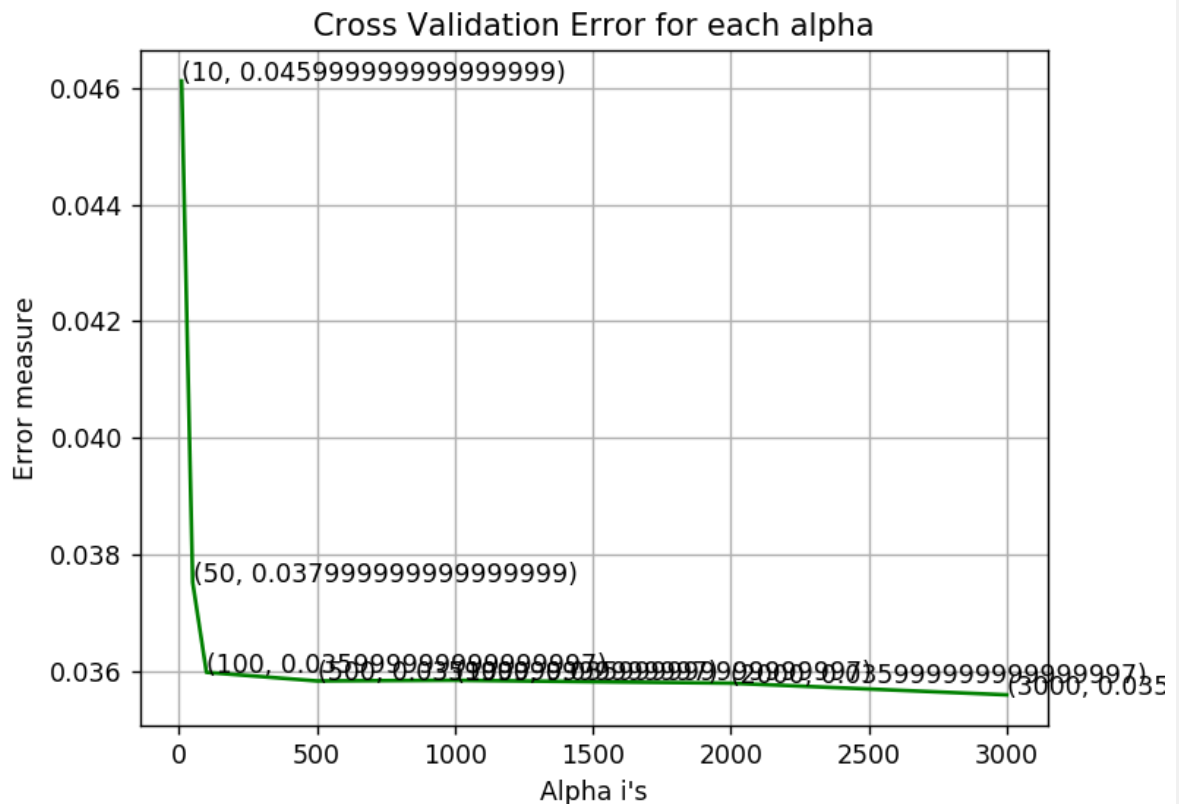
```
log_loss for c =   10 is 0.0461221662017
log_loss for c =   50 is 0.0375229563452
log_loss for c =  100 is 0.0359765822455
log_loss for c =  500 is 0.0358291883873
log_loss for c =  1000 is 0.0358403093496
log_loss for c =  2000 is 0.0357908022178
log_loss for c =  3000 is 0.0355909487962
```

```
<IPython.core.display.Javascript object>
```



Cross Validation Error for each alpha

```
For values of best alpha =   3000 The train log loss is: 0.0166267614753
For values of best alpha =   3000 The cross validation log loss is: 0.035590948
7962
For values of best alpha =   3000 The test log loss is: 0.0401141303589
```

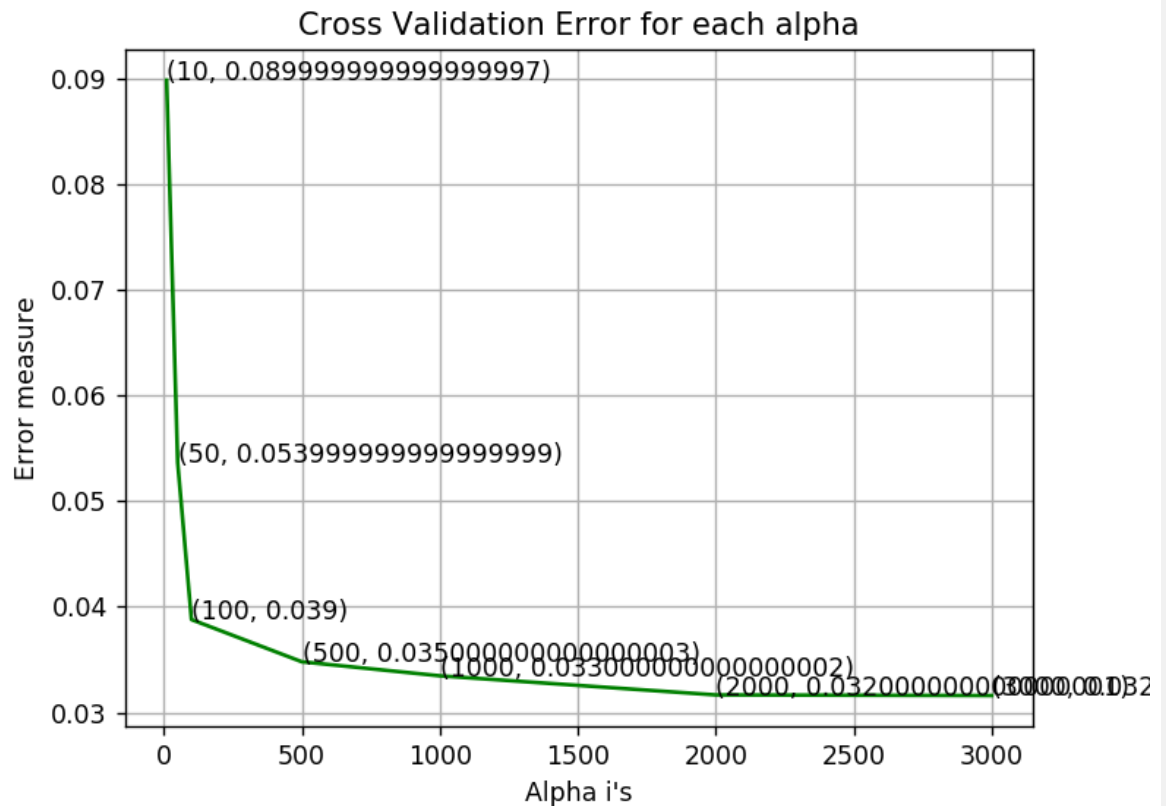## 4.5.5. XgBoost Classifier on final features

```python
In [186]:   1  # Training a hyper-parameter tuned Xg-Boost regressor on our train data
            2
            3  # find more about XGBClassifier function here http://xgboost.readthedocs.io/e
            4  # ------------------------
            5  # default paramters
            6  # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=10
            7  # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamm
            8  # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg
            9  # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None
           10
           11  # some of methods of RandomForestRegressor()
           12  # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppi
           13  # get_params([deep])    Get parameters for this estimator.
           14  # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE
           15  # get_score(importance_type='weight') -> get the feature importance
           16  # ----------------------
           17  # video link2: https://www.appliedaicourse.com/course/applied-ai-course-onlin
           18  # ----------------------
           19
           20  alpha=[10,50,100,500,1000,2000,3000]
           21  cv_log_error_array=[]
           22  for i in alpha:
           23      x_cfl=XGBClassifier(n_estimators=i)
           24      x_cfl.fit(X_train_merge,y_train_merge)
           25      sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
           26      sig_clf.fit(X_train_merge, y_train_merge)
           27      predict_y = sig_clf.predict_proba(X_cv_merge)
           28      cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=x_cfl.cl
           29
           30  for i in range(len(cv_log_error_array)):
           31      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
           32
           33
           34  best_alpha = np.argmin(cv_log_error_array)
           35
           36  fig, ax = plt.subplots()
           37  ax.plot(alpha, cv_log_error_array,c='g')
           38  for i, txt in enumerate(np.round(cv_log_error_array,3)):
           39      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
           40  plt.grid()
           41  plt.title("Cross Validation Error for each alpha")
           42  plt.xlabel("Alpha i's")
           43  plt.ylabel("Error measure")
           44  plt.show()
           45
           46  x_cfl=XGBClassifier(n_estimators=3000,nthread=-1)
           47  x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
           48  sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
           49  sig_clf.fit(X_train_merge, y_train_merge)
           50
           51  predict_y = sig_clf.predict_proba(X_train_merge)
           52  print ('For values of best alpha = ', alpha[best_alpha], "The train log loss
           53  predict_y = sig_clf.predict_proba(X_cv_merge)
           54  print('For values of best alpha = ', alpha[best_alpha], "The cross validation
           55  predict_y = sig_clf.predict_proba(X_test_merge)
           56  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
```

```
log_loss for c =   10 is 0.0898979446265
log_loss for c =   50 is 0.0536946658041
log_loss for c =   100 is 0.0387968186177
log_loss for c =   500 is 0.0347960327293
log_loss for c =   1000 is 0.0334668083237
log_loss for c =   2000 is 0.0316569078846
log_loss for c =   3000 is 0.0315972694477

<IPython.core.display.Javascript object>
```



Cross Validation Error for each alpha

```
For values of best alpha =   3000 The train log loss is: 0.0111918809342
For values of best alpha =   3000 The cross validation log loss is: 0.031597269
4477
For values of best alpha =   3000 The test log loss is: 0.0323978515915
```

## 4.5.5. XgBoost Classifier on final features with best hyper parameters using Random search

```
In [187]:   1  x_cfl=XGBClassifier()
            2
            3  prams={
            4      'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
            5      'n_estimators':[100,200,500,1000,2000],
            6      'max_depth':[3,5,10],
            7      'colsample_bytree':[0.1,0.3,0.5,1],
            8      'subsample':[0.1,0.3,0.5,1]
            9  }
           10  random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jo
           11  random_cfl.fit(X_train_merge, y_train_merge)
```

Fitting 3 folds for each of 10 candidates, totalling 30 fits

```
[Parallel(n_jobs=-1)]: Done    2 tasks       | elapsed:  1.1min
[Parallel(n_jobs=-1)]: Done    9 tasks       | elapsed:  2.2min
[Parallel(n_jobs=-1)]: Done   19 out of  30 | elapsed:  4.5min remaining:  2.6m
in
[Parallel(n_jobs=-1)]: Done   23 out of  30 | elapsed:  5.8min remaining:  1.8m
in
[Parallel(n_jobs=-1)]: Done   27 out of  30 | elapsed:  6.7min remaining:   44.
5s
[Parallel(n_jobs=-1)]: Done   30 out of  30 | elapsed:  7.4min finished
```

```
Out[187]: RandomizedSearchCV(cv=None, error_score='raise',
              estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, colsampl
          e_bytree=1,
              gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
              min_child_weight=1, missing=None, n_estimators=100, nthread=-1,
              objective='binary:logistic', reg_alpha=0, reg_lambda=1,
              scale_pos_weight=1, seed=0, silent=True, subsample=1),
              fit_params=None, iid=True, n_iter=10, n_jobs=-1,
              param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15,
          0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'co
          lsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
              pre_dispatch='2*n_jobs', random_state=None, refit=True,
              return_train_score=True, scoring=None, verbose=10)
```

```
In [188]:   1  print (random_cfl.best_params_)
```

```
{'subsample': 1, 'n_estimators': 1000, 'max_depth': 10, 'learning_rate': 0.15,
'colsample_bytree': 0.3}
```

```python
# find more about XGBClassifier function here http://xgboost.readthedocs.io/e
# --------------------------
# default paramters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=10
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamm
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppi
# get_params([deep])    Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE
# get_score(importance_type='weight') -> get the feature importance
# ----------------------
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-onlin
# ----------------------

x_cfl=XGBClassifier(n_estimators=1000,max_depth=10,learning_rate=0.15,colsamp
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

predict_y = sig_clf.predict_proba(X_train_merge)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss
predict_y = sig_clf.predict_proba(X_cv_merge)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation
predict_y = sig_clf.predict_proba(X_test_merge)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_merge))
```

```
For values of best alpha =  3000 The train log loss is: 0.0121922832297
For values of best alpha =  3000 The cross validation log loss is: 0.034495548
7471
For values of best alpha =  3000 The test log loss is: 0.0317041132442
```