

# 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 finish 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, ss:nothin
g, 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     ??0exceptio
n@std@@@QAE@ABQBD@Z ; std::exception::exception(char const * const &)
.text:0040100D C7 06 08    BB 42 00                                mov      dword ptr [e
si],      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 10h
.text:00401020 C7 01 08    BB 42 00                                mov      dword ptr [e
cx],      offset off_42BB08
.text:00401026 E9 26 1C    00 00                                jmp      sub_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      dword ptr [e
si],      offset off_42BB08
.text:00401039 E8 13 1C    00 00                                call     sub_402C51
.text:0040103E F6 44 24    08 01                                test     byte 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     ???@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 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
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
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 classify 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 Learning 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 probabilities => 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=VLQTRILGz5Y>

<https://github.com/dchad/malware-detection>

<http://vizsec.org/files/2011/Nataraj.pdf>

[https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu\\_plB6ua?dl=0](https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu_plB6ua?dl=0)

" Cross validation is more trustworthy than domain knowledge."

## 3. Exploratory Data Analysis

In [1]:

```
1 import warnings
2 warnings.filterwarnings("ignore")
3
4 # File operations:
5 import shutil
6 import os
7 import codecs# this is used for file operations
8
9 # Data structures:
10 import pandas as pd
11 import pandas_profiling as pp
12 import numpy as np
13 import random as r
14
15 # Visualizations:
16 import matplotlib
17 matplotlib.use(u'nbAgg')
18 import matplotlib.pyplot as plt
19 %matplotlib inline
20 import seaborn as sns
21 from sklearn.manifold import TSNE
22
23 # Code speed improvement:
24 from multiprocessing import Process# this is used for multithreading
25 import multiprocessing as mp
26
27 # Transformations:
28 from sklearn import preprocessing
29
30 # Models:
31 from xgboost import XGBClassifier
32 from sklearn.tree import DecisionTreeClassifier
33 from sklearn.calibration import CalibratedClassifierCV
34 from sklearn.neighbors import KNeighborsClassifier
35 from sklearn.linear_model import LogisticRegression
36 from sklearn.ensemble import RandomForestClassifier
37
38 # CV:
39 from sklearn.model_selection import train_test_split, RandomizedSearchCV
40
41 # Performance Metrics:
42 from sklearn.metrics import confusion_matrix, log_loss, accuracy_score
43
44 # Object persistance:
45 import pickle
46 from sklearn.externals import joblib
47
48 from tqdm import tqdm
49
50 # date and time:
51 from datetime import datetime
52 import time
53
54 from prettytable import PrettyTable
```

c:\users\byron\applications\pythonmaster\lib\site-packages\sklearn\ensemble

```
\weight_boosting.py:29: DeprecationWarning: numpy.core.umath_tests is an internal NumPy module and should not be imported. It will be removed in a future NumPy release.
```

```
from numpy.core.umath_tests import inner1d
```

In [2]:

```
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 create a folder
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 .bytes files)
11 # for every file that we have in our 'asmFiles' directory we check if it is ending with .asm
12 # 'byteFiles' folder
13
14 # so by the end of this snippet we will separate all the .byte files and .asm files
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         if (file.endswith("bytes")):
21             shutil.move(source+'\\'+file, destination)
```

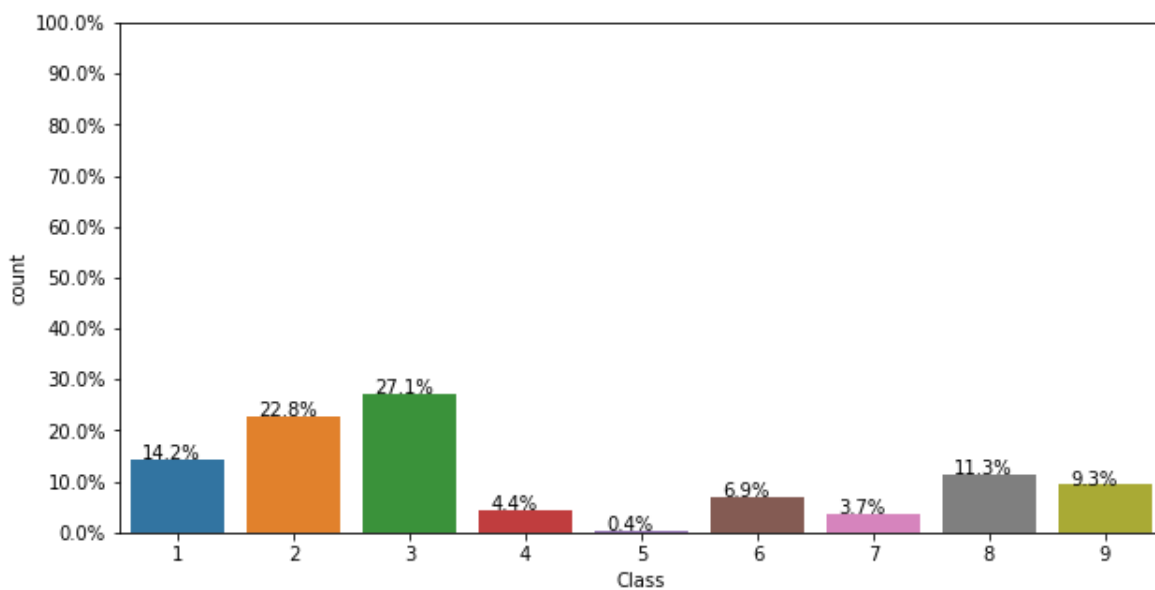
### 3.1. Distribution of malware classes in whole data set

In [3]:

```

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

```



## 3.2. Feature extraction

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



In [4]:

```

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  num_lines=[]
9  fnames=[]
10 for file in tqdm(files):
11     # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
12     # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, st_nli
13     # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=1519638522)
14     # read more about os.stat: here https://www.tutorialspoint.com/python/os_stat.htm
15     statinfo=os.stat('byteFiles/'+file)
16     with open('byteFiles/'+file,'r') as content:
17         num_of_lines = len(content.readlines())
18     # split the file name at '.' and take the first part of it i.e the file name
19     file=file.split('.')[0]
20     if any(file == filename for filename in filenames):
21         i=filenames.index(file)
22         class_bytes.append(class_y[i])
23         # converting into Mb's
24         sizebytes.append(statinfo.st_size/(1024.0*1024.0))
25         num_lines.append(num_of_lines)
26         fnames.append(file)
27 data_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes, 'filelines':num_lines, 'Cla
28 print (data_size_byte.head())

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 10868/10868 [07:39<00:00, 23.65it/s]

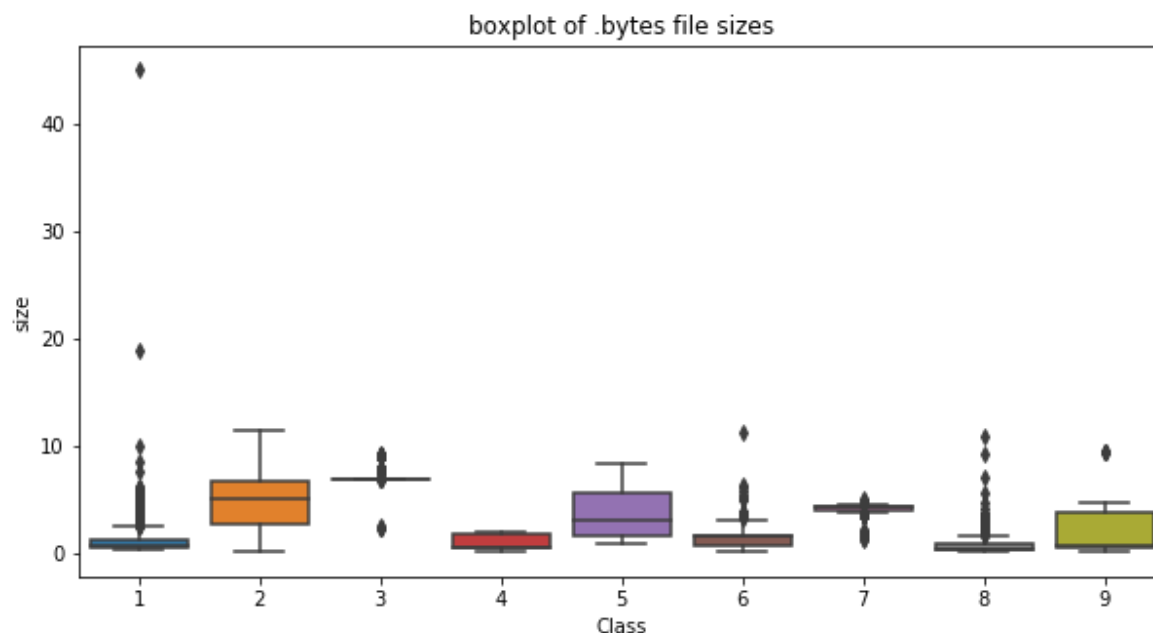
```

	ID	size	filelines	Class
0	01azqd4InC7m9JpocGv5	4.234863	90624	9
1	01IsoiSMh5gxyDYTl4CB	5.538818	118528	2
2	01jsnpXSAlgW6aPeDxrU	3.887939	83200	9
3	01kcPWA9K2B0xQeS5Rju	0.574219	12288	1
4	01SuzwMJEIXsK7A8dQbl	0.370850	7936	8

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

In [5]:

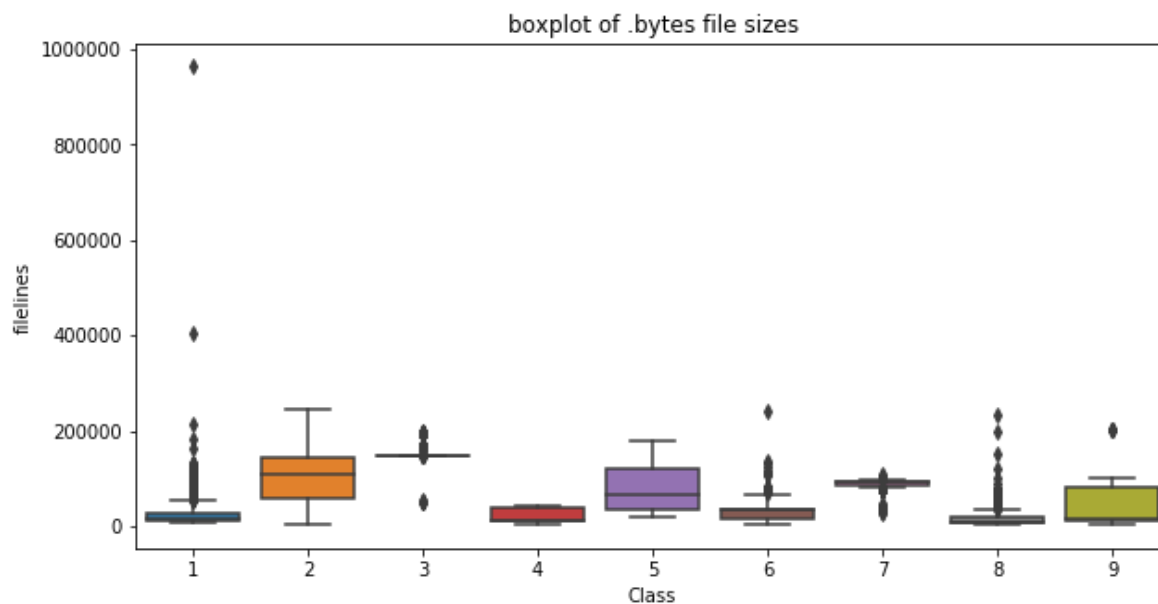
```
1 #boxplot of byte files
2 plt.figure(figsize = (10,5))
3 ax = sns.boxplot(x="Class", y="size", data=data_size_byte)
4 plt.title("boxplot of .bytes file sizes")
5 plt.show()
```



### 3.2.2.2 box plots of file row numbers (.byte files) feature

In [6]:

```
1 #boxplot of byte files
2 plt.figure(figsize = (10,5))
3 ax = sns.boxplot(x="Class", y="filelines", data=data_size_byte)
4 plt.title("boxplot of .bytes file sizes")
5 plt.show()
```



### 3.2.3 feature extraction from byte files

In [7]:

```

1  #removal of address from byte files
2  # contents of .byte files
3  # -----
4  #00401000 56 8D 44 24 08 50 8B F1 E8 1C 1B 00 00 C7 06 08
5  #-----
6  #we remove the starting address 00401000
7  files = os.listdir('byteFiles')
8  filenames=[]
9  array=[]
10 for file in tqdm(files):
11     if(file.endswith("bytes")):
12         filename=file.split('.')[0]
13         text_file = open('byteFiles\\'+filename+'.txt', 'w+')
14         with open('byteFiles\\'+file,"r") as fp:
15             lines=""
16             for line in fp:
17                 a=line.rstrip().split(" ")[1:]
18                 b=' '.join(a)
19                 b=b+"\n"
20                 text_file.write(b)
21             fp.close()
22             os.remove('byteFiles\\'+file)
23         text_file.close()

```

```
100%|███████████| 
| 10868/10868 [00:00<00:00, 1815794.13it/s]
```

In [8]:

```
1  #program to convert into bag of words of bytefiles
2  #this is custom-built bag of words this is unigram bag of words
3  #this runs for about 3 hours so be patient on this section
4  if not os.path.exists('result.csv'):
5      files = os.listdir('byteFiles')
6      filenames2=[]
7      feature_matrix = np.zeros((len(files),257),dtype=int)
8      k=0
9      byte_feature_file=open('result.csv','w+')
10     byte_feature_file.write("ID,0,1,2,3,4,5,6,7,8,9,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,17,18,19,1a,1b,1c,1d,1e,1f,20,21,22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,4a,4b,4c,4d,4e,4f,50,51,52,53,54,55,56,57,58,59,5a,5b,5c,5d,5e,5f,60,61,62,63,64,65,66,67,68,69,6a,6b,6c,6d,6e,6f,70,71,72,73,74,75,76,77,78,79,7a,7b,7c,7d,7e,7f,80,81,82,83,84,85,86,87,88,89,8a,8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,a0,a1,a2,a3,a4,a5,a6,a7,a8,a9,aa,ab,ac,ad,ae,af,b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,ba,bb,bc,bd,be,bf,c0,c1,c2,c3,c4,c5,c6,c7,c8,c9,ca,cb,cc,cd,ce,cf,d0,d1,d2,d3,d4,d5,d6,d7,d8,d9,da,db,dc,dd,de,df,e0,e1,e2,e3,e4,e5,e6,e7,e8,e9,ea,eb,ec,ed,ee,ef,f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,ff")
11     byte_feature_file.write("\n")
12     for file in tqdm(files):
13         filenames2.append(file.split('.')[0])
14         byte_feature_file.write(file.split('.')[0]+",")
15         if(file.endswith("txt")):
16             with open('byteFiles\\'+file,"r") as byte_flie:
17                 for lines in byte_flie:
18                     line=line.rstrip().split(" ")
19                     for hex_code in line:
20                         if hex_code=='?':
21                             feature_matrix[k][256]+=1
22                         else:
23                             feature_matrix[k][int(hex_code,16)]+=1
24             byte_flie.close()
25         for i in feature_matrix[k]:
26             byte_feature_file.write(str(i)+",")
27         byte_feature_file.write("\n")
28
29     k += 1
30     byte_feature_file.close()
```

In [9]:

```
1 byte_features=pd.read_csv("result.csv",index_col=False)
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	182
4												
1	9291	...	451	6536	439	281	302	7639	518	17001	54902	858
8												
2	9107	...	2325	2358	2242	2885	2863	2471	2786	2680	49144	46
8												
3	1078	...	478	873	485	462	516	1133	471	761	7998	1394
0												
4	422	...	847	947	350	209	239	653	221	242	2199	900
8												

[5 rows x 258 columns]

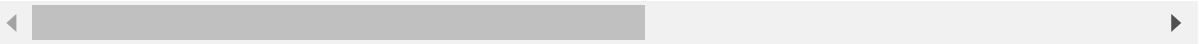
In [10]:

```
1 result = pd.merge(byte_features, data_size_byte,on='ID', how='left')
2 result['bytezeroweighted'] = result['0']*2
3 result.head()
```

Out[10]:

	ID	0	1	2	3	4	5	6	7	8	...	fb
0	01azqd4InC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...	3097
1	01IsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...	302
2	01jsnpXSAlgW6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	9107	...	2863
3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	1078	...	516
4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	422	...	239

5 rows × 262 columns



In [11]:

```

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_value - min_value)
9     return result1
10 result = normalize(result)

```

In [12]:

```

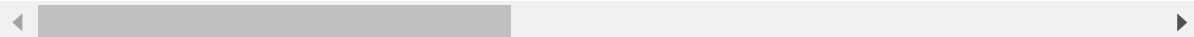
1 data_y = result['Class']
2 result.head()

```

Out[12]:

	ID	0	1	2	3	4	5	
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.0020
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.0047
2	01jsnpXSAIgw6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.0050
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.0003
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.0001

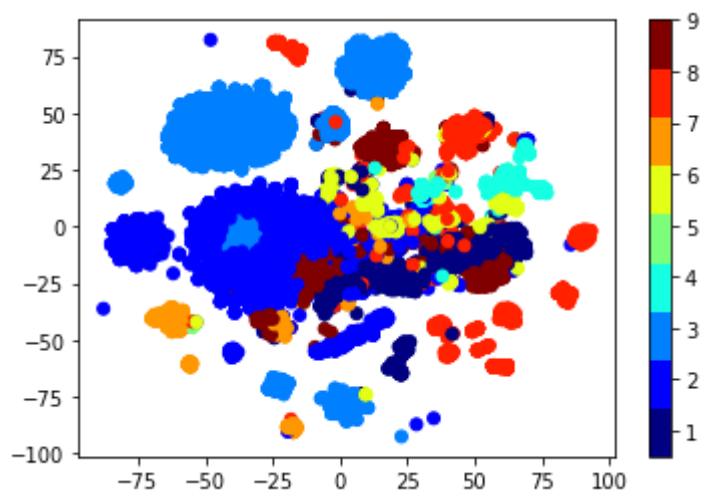
5 rows × 262 columns



### 3.2.4 Multivariate Analysis

In [13]:

```
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()
```



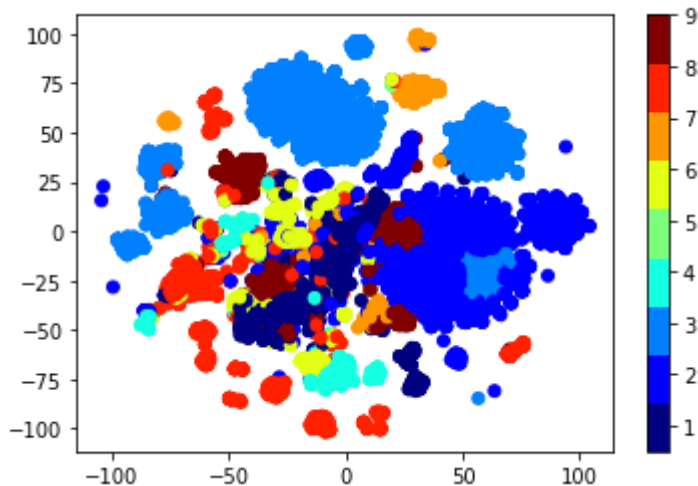


In [14]:

```

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()

```



## dChad Features

In [15]:

```

1 dchad = pd.DataFrame()
2 for file in os.listdir('dchad'):
3     if file[-3:] == 'csv':
4         filelocation = 'dchad/'+file
5         dchad = pd.concat(objs=[dchad,pd.read_csv(filelocation)],axis=0)
6 dchad.rename(columns={'filename':'ID'},inplace=True)
7 dchad.head()

```

Out[15]:

	ID	vertex_count	edge_count	delta_max	density
0	4jKA1GUDv6TMNpPulxER	3243	5440	85	0.005500
1	4ZBJzEqnW52fFUw0PG3v	40	39	36	13.000000
2	6m8NxLfg2MR0nwXFuEq5	95	94	61	0.029012
3	28U1hRkQ6YI57493ZdXD	120	119	10	0.200000
4	45Wy3TxE98HfiXreOCSu	339	337	50	0.047999

In [16]:

```
1 dchad1 = normalize(dchad)
```

In [17]:

```
1 dchad1.head()
```

Out[17]:

	ID	vertex_count	edge_count	delta_max	density
0	4jKA1GUDv6TMNpPulxER	0.113854	0.139936	0.010612	0.000015
1	4ZBJzEqnW52fFUw0PG3v	0.001370	0.001003	0.004494	0.034685
2	6m8NxLfg2MR0nwXFuEq5	0.003301	0.002418	0.007615	0.000077
3	28U1hRkQ6YI57493ZdXD	0.004179	0.003061	0.001248	0.000534
4	45Wy3TxE98HfiXreOCSu	0.011870	0.008669	0.006242	0.000128

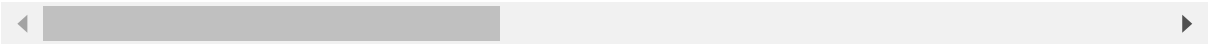
In [18]:

```
1 result = pd.merge(result,dchad1, on = 'ID', how='inner')
2 result.head()
```

Out[18]:

	ID	0	1	2	3	4	5	
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.0020
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.0047
2	01jsnpXSAIgw6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.0050
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.0003
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.0001

5 rows × 266 columns



# Train Test split

In [19]:

```
1 data_y = result['Class']
2 # split the data into test and train by maintaining same distribution of output variable
3 X_train, X_test, y_train, y_test = train_test_split(result.drop(['ID', 'Class'], axis=1), data_y, test_size=0.1, random_state=42)
4 # split the train data into train and cross validation by maintaining same distribution
5 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.1, random_state=42)
```

In [20]:

```
1 print('Number of data points in train data:', X_train.shape[0])
2 print('Number of data points in test data:', X_test.shape[0])
3 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

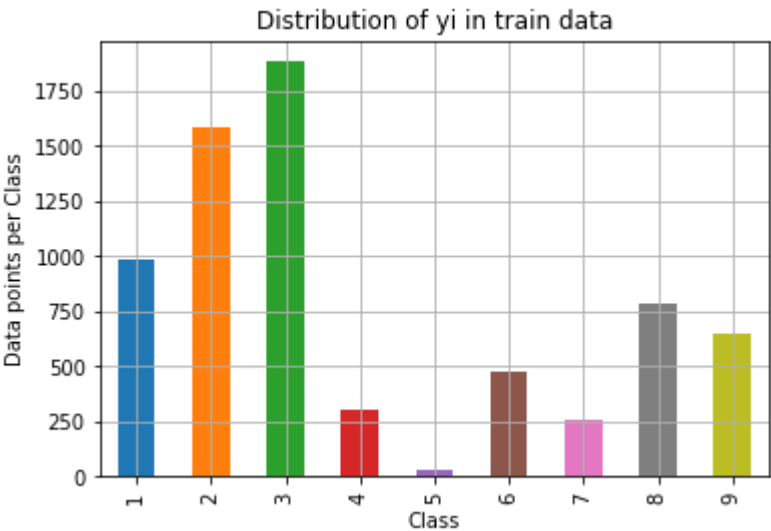
Number of data points in cross validation data: 1739

In [21]:

```

1  # it returns a dict, keys as class labels and values as the number of data points in t
2  train_class_distribution = y_train.value_counts().sortlevel()
3  test_class_distribution = y_test.value_counts().sortlevel()
4  cv_class_distribution = y_cv.value_counts().sortlevel()
5
6  train_class_distribution.plot(kind='bar')
7  plt.xlabel('Class')
8  plt.ylabel('Data points per Class')
9  plt.title('Distribution of yi in train data')
10 plt.grid()
11 plt.show()
12
13 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
14 # -(train_class_distribution.values): the minus sign will give us in decreasing order
15 sorted_yi = np.argsort(-train_class_distribution.values)
16 for i in sorted_yi:
17     print('Number of data points in class', i+1, ':', train_class_distribution.values[i]
18
19
20 print('-'*80)
21 test_class_distribution.plot(kind='bar')
22 plt.xlabel('Class')
23 plt.ylabel('Data points per Class')
24 plt.title('Distribution of yi in test data')
25 plt.grid()
26 plt.show()
27
28 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
29 # -(train_class_distribution.values): the minus sign will give us in decreasing order
30 sorted_yi = np.argsort(-test_class_distribution.values)
31 for i in sorted_yi:
32     print('Number of data points in class', i+1, ':', test_class_distribution.values[i]
33
34 print('-'*80)
35 cv_class_distribution.plot(kind='bar')
36 plt.xlabel('Class')
37 plt.ylabel('Data points per Class')
38 plt.title('Distribution of yi in cross validation data')
39 plt.grid()
40 plt.show()
41
42 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
43 # -(train_class_distribution.values): the minus sign will give us in decreasing order
44 sorted_yi = np.argsort(-train_class_distribution.values)
45 for i in sorted_yi:
46     print('Number of data points in class', i+1, ':', cv_class_distribution.values[i],
47

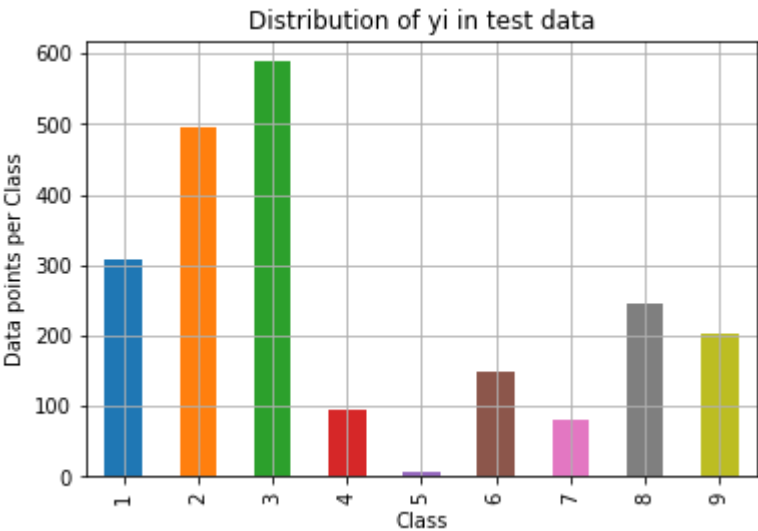
```



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 %)

-----

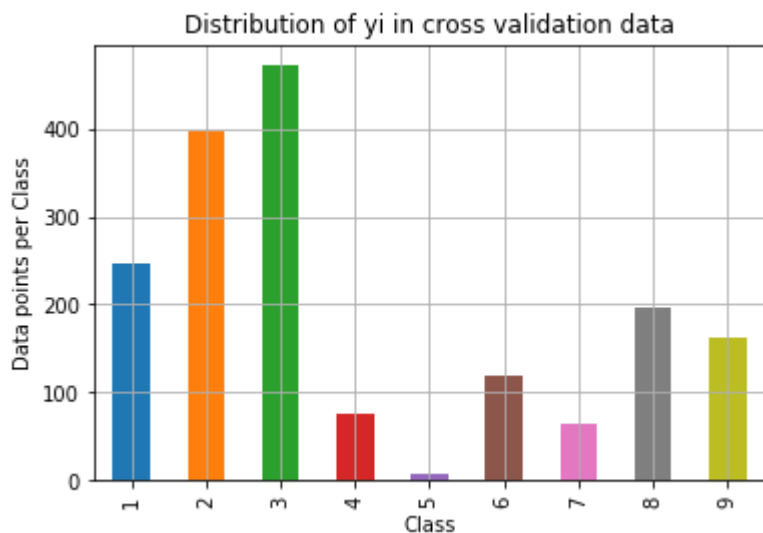
----



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 %)

-----

----



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 %)

In [22]:

```

1  def plot_confusion_matrix(test_y, predict_y, model):
2
3      C = confusion_matrix(test_y, predict_y)
4      # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted as class j
5
6      A = ((C.T)/(C.sum(axis=1))).T
7      #divid each element of the confusion matrix with the sum of elements in that column
8
9      # C = [[1, 2],
10     #      [3, 4]]
11     # C.T = [[1, 3],
12     #        [2, 4]]
13     # C.sum(axis = 1) axis=0 corresponds to columns and axis=1 corresponds to rows in C
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 that row
24     # C = [[1, 2],
25     #      [3, 4]]
26     # C.sum(axis = 0) axis=0 corresponds to columns and axis=1 corresponds to rows in C
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, yticklabels=labels)
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, yticklabels=labels)
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, yticklabels=labels)
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))
57

```

```
58 miss_classified_rate = round((1 - accuracy_score(test_y, predict_y))*100,2)
59 print("Percentage of misclassified points using {} Model: {}".format(model, str(miss_
```

## 4. Machine Learning Models

### 4.1. Machine Learning Models on bytes files

#### 4.1.1. Random Model



In [23]:

```

1  # we need to generate 9 numbers and the sum of numbers should be 1
2  # one solution is to generate 9 numbers and divide each of the numbers by their sum
3  # ref: https://stackoverflow.com/a/18662466/4084039
4
5  train_data_len = X_train.shape[0]
6  test_data_len = X_test.shape[0]
7  cv_data_len = X_cv.shape[0]
8
9  # we create a output array that has exactly same size as the train data
10 random_pred_train = np.random.rand(train_data_len,9)
11 for key, value in dict(zip(np.sum(random_pred_train,axis=1),random_pred_train)).items():
12     dict(zip(np.sum(random_pred_train,axis=1),random_pred_train))[key] = value/key
13 print("Log loss on Training Data using Random Model: ",log_loss(y_train,random_pred_train))
14
15 # we create a output array that has exactly same size as the CV data
16 random_pred_cv = np.random.rand(cv_data_len,9)
17 for key, value in dict(zip(np.sum(random_pred_cv,axis=1),random_pred_cv)).items():
18     dict(zip(np.sum(random_pred_cv,axis=1),random_pred_cv))[key] = value/key
19 print("Log loss on Cross Validation Data using Random Model: ",log_loss(y_cv,random_pred_cv))
20
21
22 # we create a output array that has exactly same as the test data
23 random_pred_test = np.random.rand(test_data_len,9)
24 for key, value in dict(zip(np.sum(random_pred_test,axis=1),random_pred_test)).items():
25     dict(zip(np.sum(random_pred_test,axis=1),random_pred_test))[key] = value/key
26 print("Log loss on Test Data using Random Model: ",log_loss(y_test,random_pred_test, epsilon=1e-15))
27
28 predicted_y = np.argmax(random_pred_test, axis=1) + 1 # here we return the index(class)
29 plot_confusion_matrix(y_test, predicted_y, 'Random')

```

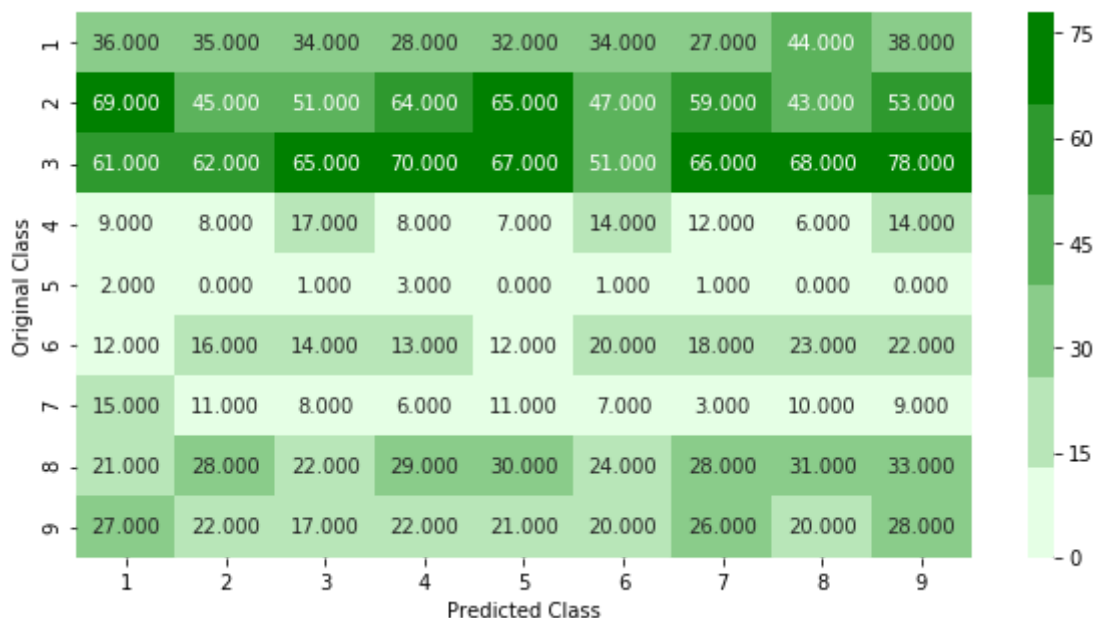
Log loss on Training Data using Random Model: 2.485348708199024

Log loss on Cross Validation Data using Random Model: 2.4860901717051105

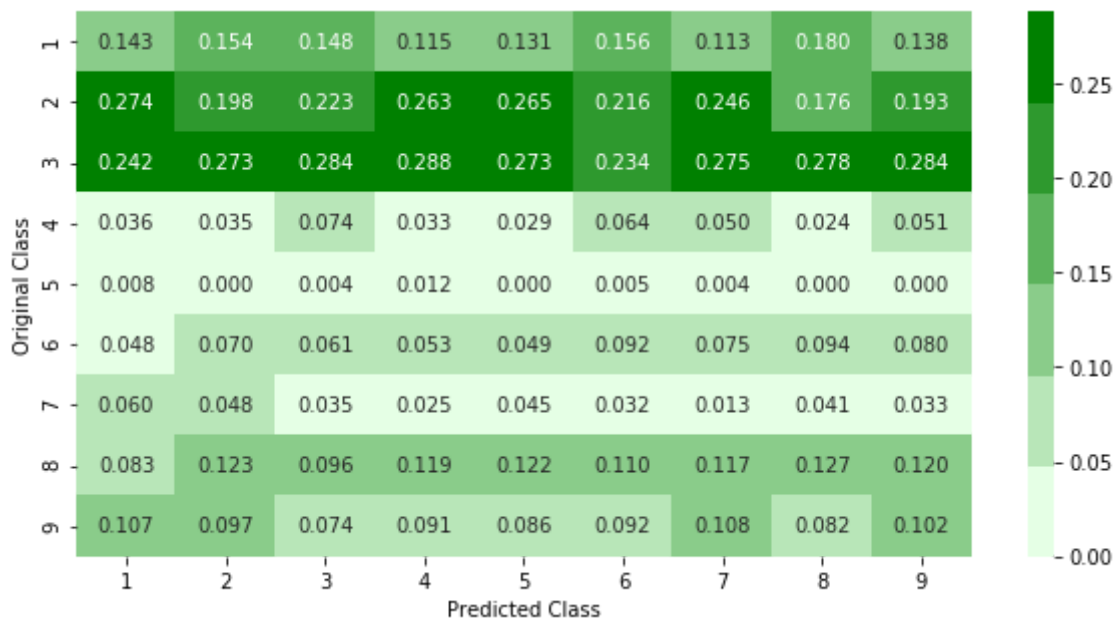
Log loss on Test Data using Random Model: 2.4631112739428467

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

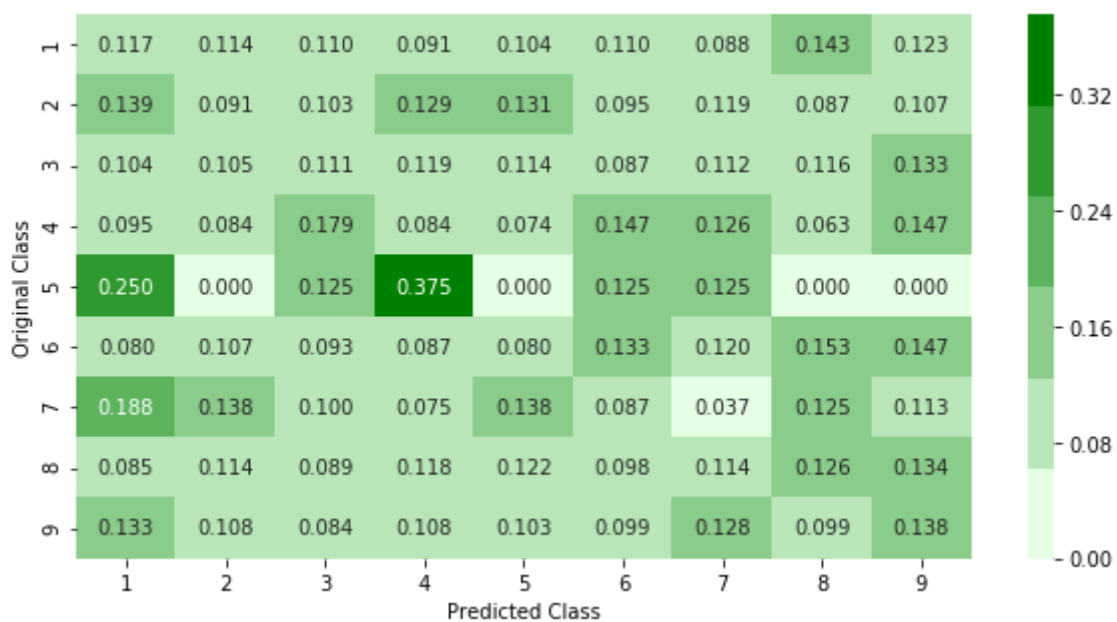
-----



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



----- Recall matrix -----



Percentage of misclassified points using Random Model: 89.14%

#### 4.1.2. K Nearest Neighbour Classification

In [24]:

```

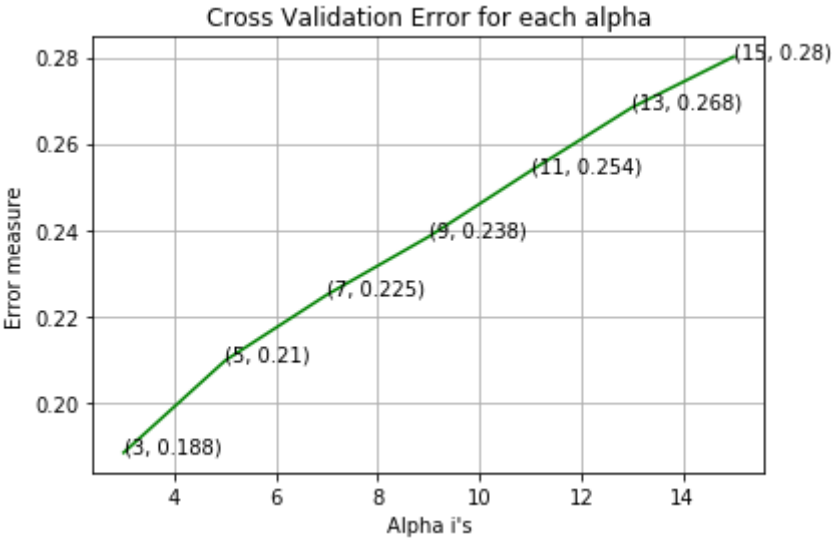
1  alpha = [x for x in [3,5,7,9,11,13,15]]
2  cv_log_error_array=[]
3  for i in alpha:
4      k_cfl=KNeighborsClassifier(n_neighbors=i)
5      k_cfl.fit(X_train,y_train)
6      sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
7      sig_clf.fit(X_train, y_train)
8      predict_y = sig_clf.predict_proba(X_cv)
9      cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cfl.classes_, eps=1e-15))
10
11  for i in zip(alpha,cv_log_error_array):
12      print ('log_loss for k = ',i[0],'is',i[1])
13
14  best_alpha = np.argmin(cv_log_error_array)
15
16  fig, ax = plt.subplots()
17  ax.plot(alpha, cv_log_error_array,c='g')
18  for i, txt in enumerate(np.round(cv_log_error_array,3)):
19      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
20  plt.grid()
21  plt.title("Cross Validation Error for each alpha")
22  plt.xlabel("Alpha i's")
23  plt.ylabel("Error measure")
24  plt.show()
25
26  k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
27  k_cfl.fit(X_train,y_train)
28  sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
29  sig_clf.fit(X_train, y_train)
30
31  predict_y = sig_clf.predict_proba(X_train)
32  print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,predict_y,labels=k_cfl.classes_,eps=1e-15))
33  predict_y = sig_clf.predict_proba(X_cv)
34  print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv,predict_y,labels=k_cfl.classes_,eps=1e-15))
35  predict_y = sig_clf.predict_proba(X_test)
36  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test,predict_y,labels=k_cfl.classes_,eps=1e-15))
37  plot_confusion_matrix(y_test, sig_clf.predict(X_test),'KNN(n = {})'.format(alpha[best_alpha]))

```

```

log_loss for k = 3 is 0.18837038397077538
log_loss for k = 5 is 0.20974175439876952
log_loss for k = 7 is 0.2249754515332201
log_loss for k = 9 is 0.23842533912526098
log_loss for k = 11 is 0.2536408767785393
log_loss for k = 13 is 0.2683846685296005
log_loss for k = 15 is 0.2801513482426867

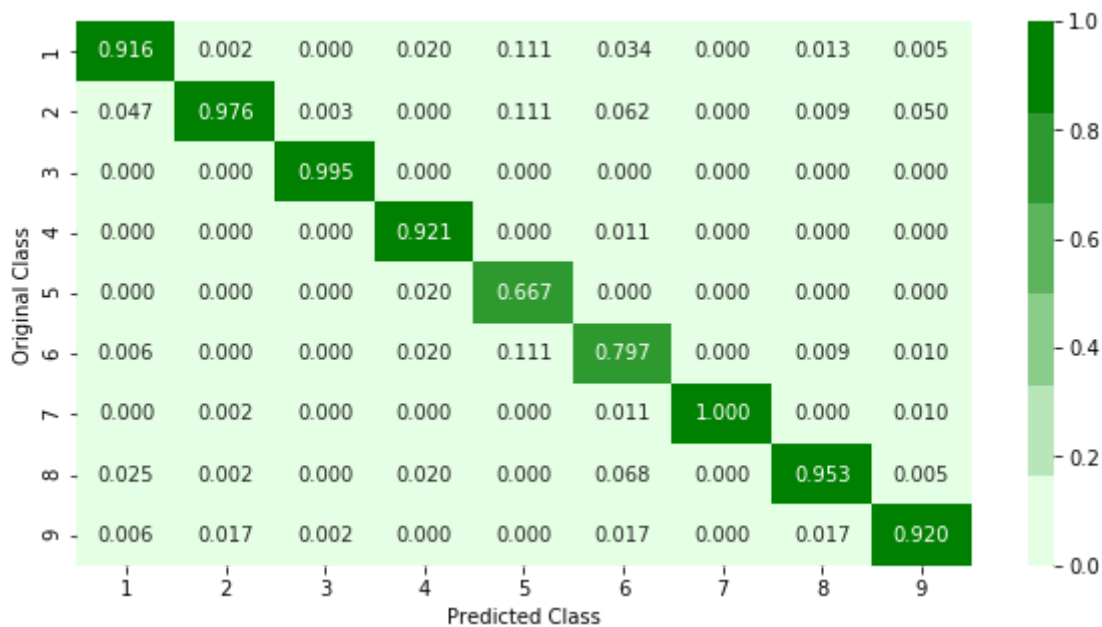
```



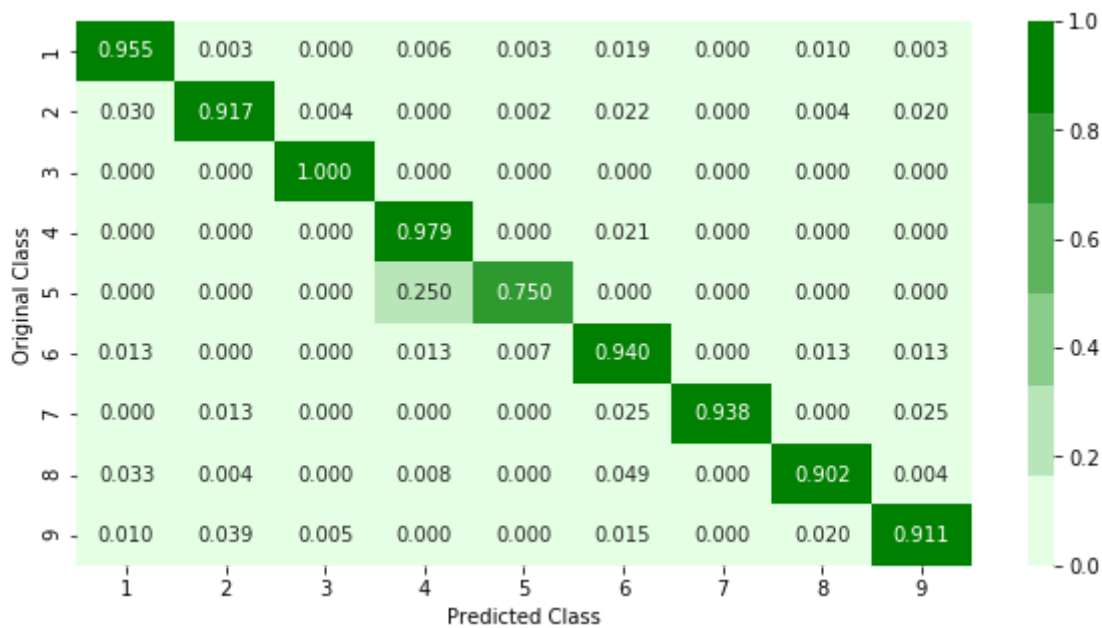
For values of best alpha = 3 The train log loss is: 0.11842650208618187  
For values of best alpha = 3 The cross validation log loss is: 0.18837038397077538  
For values of best alpha = 3 The test log loss is: 0.20877787985735408  
----- Confusion matrix -----  
-----



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



Recall matrix



Percentage of misclassified points using KNN(n = 3) Model: 5.29%

### 4.1.3. Logistic Regression

In [25]:

```

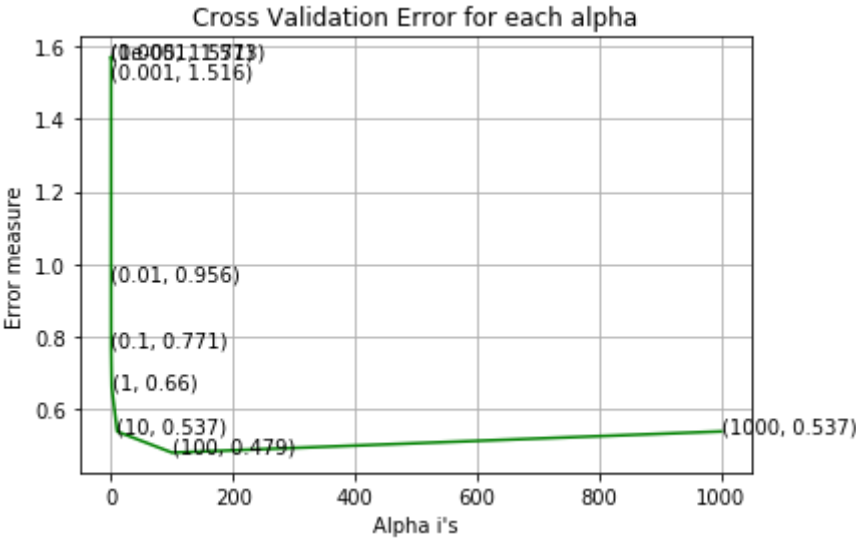
1  alpha = [10 ** x for x in range(-5, 4)]
2  cv_log_error_array=[]
3  for i in alpha:
4      logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
5      logisticR.fit(X_train,y_train)
6      sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
7      sig_clf.fit(X_train, y_train)
8      predict_y = sig_clf.predict_proba(X_cv)
9      cv_log_error_array.append(log_loss(y_cv, predict_y, labels=logisticR.classes_, eps:
10
11  for i in range(len(cv_log_error_array)):
12      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
13
14  best_alpha = np.argmin(cv_log_error_array)
15
16  fig, ax = plt.subplots()
17  ax.plot(alpha, cv_log_error_array,c='g')
18  for i, txt in enumerate(np.round(cv_log_error_array,3)):
19      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
20  plt.grid()
21  plt.title("Cross Validation Error for each alpha")
22  plt.xlabel("Alpha i's")
23  plt.ylabel("Error measure")
24  plt.show()
25
26  logisticR=LogisticRegression(penalty='l2',C=alpha[best_alpha],class_weight='balanced')
27  logisticR.fit(X_train,y_train)
28  sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
29  sig_clf.fit(X_train, y_train)
30
31
32  predict_y = sig_clf.predict_proba(X_train)
33  print ('log loss for train data',log_loss(y_train, predict_y, labels=logisticR.classes_
34  predict_y = sig_clf.predict_proba(X_cv)
35  print ('log loss for cv data',log_loss(y_cv, predict_y, labels=logisticR.classes_, eps:
36  predict_y = sig_clf.predict_proba(X_test)
37  print ('log loss for test data',log_loss(y_test, predict_y, labels=logisticR.classes_
38  plot_confusion_matrix(y_test, sig_clf.predict(X_test),'Logistic Regression (C = {})'.f

```

```

log_loss for c = 1e-05 is 1.5707695803563106
log_loss for c = 0.0001 is 1.5730050621309748
log_loss for c = 0.001 is 1.5162064471788645
log_loss for c = 0.01 is 0.9555554238450514
log_loss for c = 0.1 is 0.7711637984102886
log_loss for c = 1 is 0.6595735959566625
log_loss for c = 10 is 0.5368871361250238
log_loss for c = 100 is 0.47860793818742364
log_loss for c = 1000 is 0.5372097559827

```

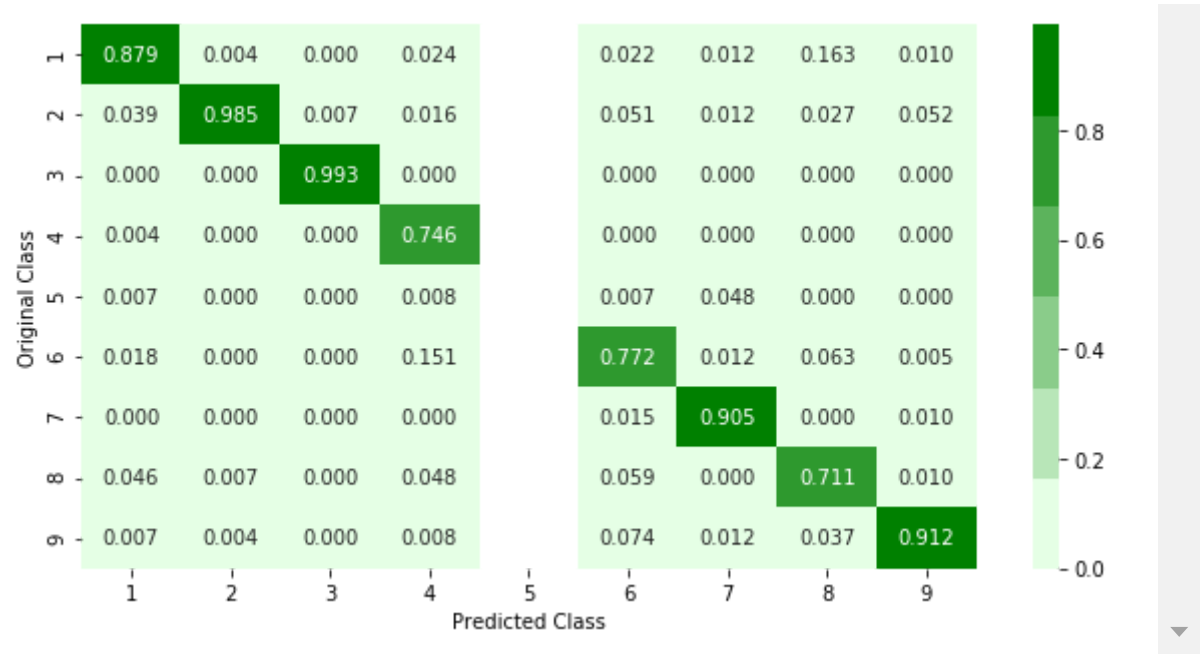


log loss for train data 0.4347852779128462  
log loss for cv data 0.47860793818742364  
log loss for test data 0.45486238918841615

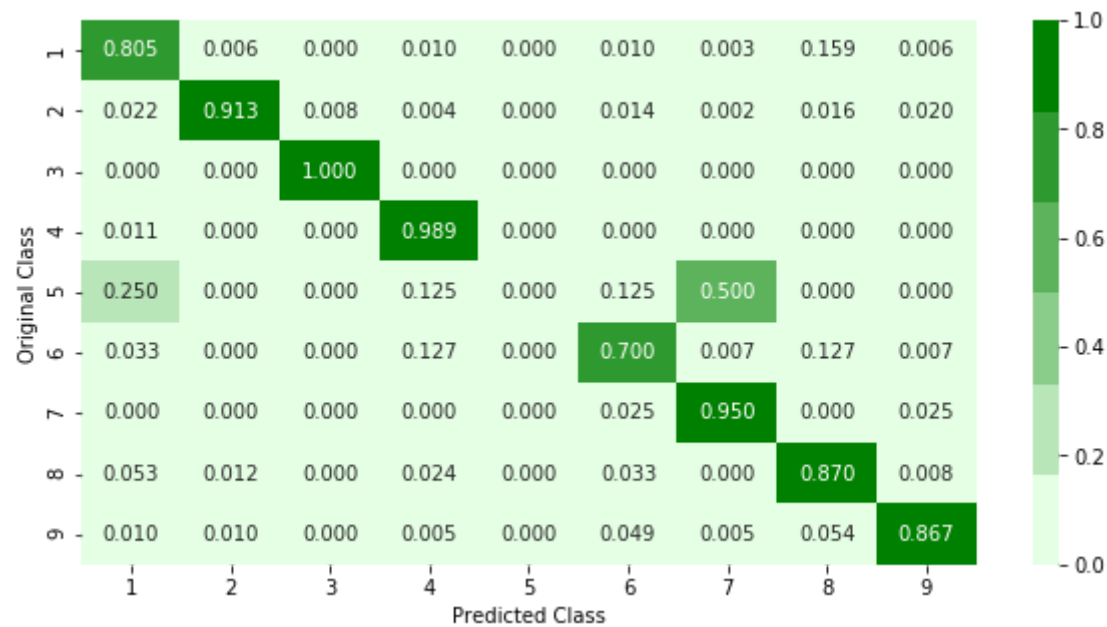
----- Confusion matrix -----  
-----



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



----- Recall matrix -----



Percentage of misclassified points using Logistic Regression (C = 100) Mode 1: 10.12%

4.1.4. Random Forest Classifier



In [26]:

```

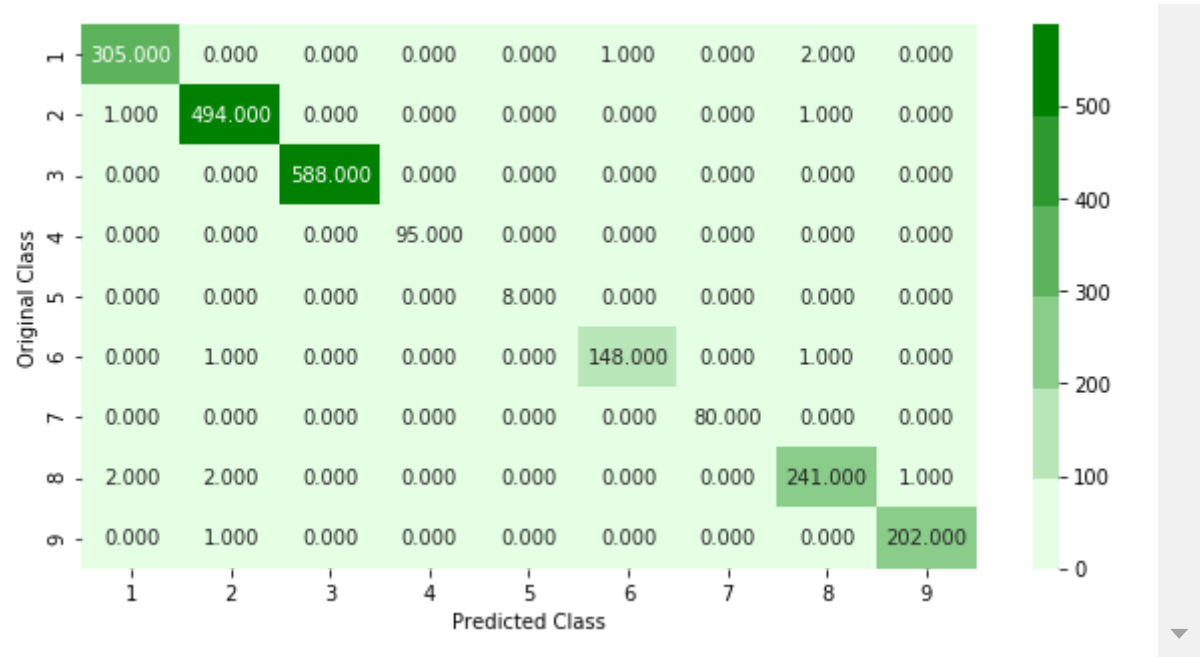
1  if not os.path.exists('RF_one_param_tuning.joblib'):
2      alpha=[10,50,100,500,1000,2000,3000]
3      cv_log_error_array=[]
4      train_log_error_array=[]
5      for i in alpha:
6          r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
7          r_cfl.fit(X_train,y_train)
8          sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
9          sig_clf.fit(X_train, y_train)
10         predict_y = sig_clf.predict_proba(X_cv)
11         cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps:
12
13     for i in range(len(cv_log_error_array)):
14         print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
15
16
17     best_alpha = np.argmin(cv_log_error_array)
18
19     fig, ax = plt.subplots()
20     ax.plot(alpha, cv_log_error_array,c='g')
21     for i, txt in enumerate(np.round(cv_log_error_array,3)):
22         ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
23     plt.grid()
24     plt.title("Cross Validation Error for each alpha")
25     plt.xlabel("Alpha i's")
26     plt.ylabel("Error measure")
27     plt.show()
28
29     r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs:
30     r_cfl.fit(X_train,y_train)
31     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
32     sig_clf.fit(X_train, y_train)
33
34     joblib.dump(sig_clf,'RF_one_param_tuning.joblib')
35 else:
36     sig_clf = joblib.load('RF_one_param_tuning.joblib')
37
38 predict_y = sig_clf.predict_proba(X_train)
39 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_l
40 predict_y = sig_clf.predict_proba(X_cv)
41 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
42 predict_y = sig_clf.predict_proba(X_test)
43 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_lo
44 plot_confusion_matrix(y_test, sig_clf.predict(X_test),'Random Forest (n_estimators = {

```

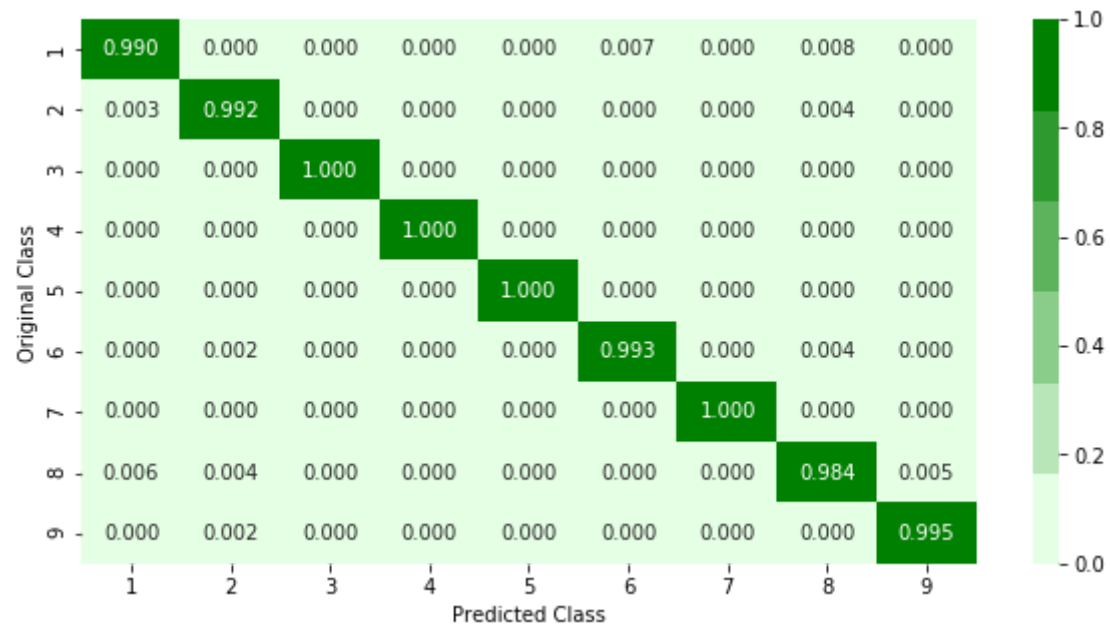
```

For values of best alpha = 100 The train log loss is: 0.044166783323594104
For values of best alpha = 100 The cross validation log loss is: 0.03739752
427035803
For values of best alpha = 100 The test log loss is: 0.04016687451345051
----- Confusion matrix -----
-----

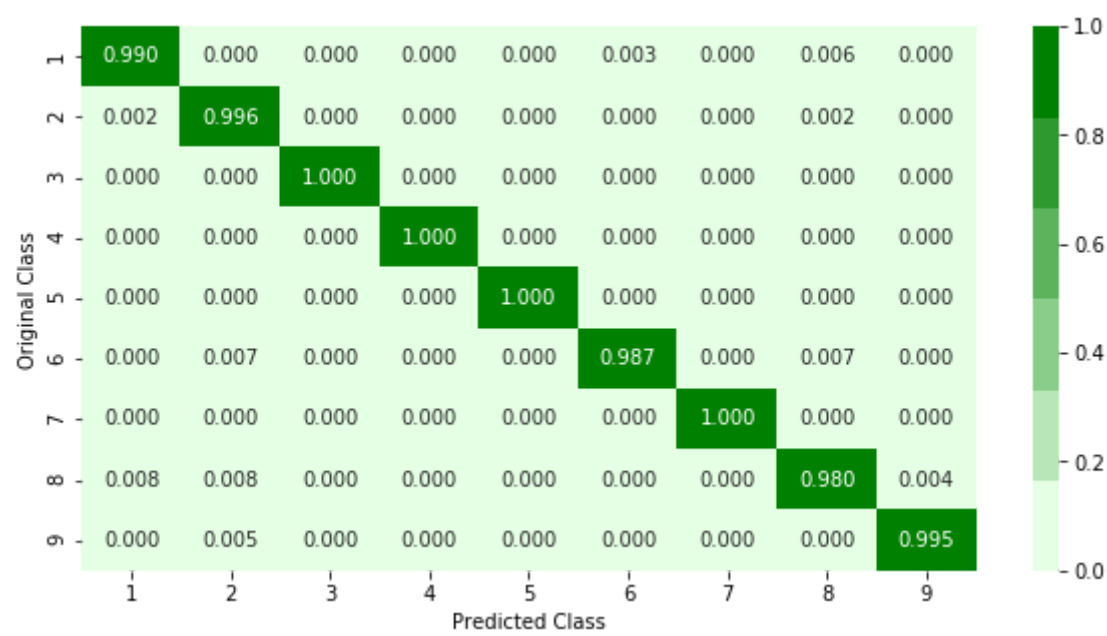
```



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



----- Recall matrix -----  
-----



Percentage of misclassified points using Random Forest (n\_estimators = 100)  
Model: 0.6%

4.1.5. XgBoost Classification

In [27]:

```

1  if not os.path.exists('XGB_one_param_tuning.joblib'):
2      alpha=[10,50,100,500,1000,2000]
3      cv_log_error_array=[]
4      for i in alpha:
5          x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
6          x_cfl.fit(X_train,y_train)
7          sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
8          sig_clf.fit(X_train, y_train)
9          predict_y = sig_clf.predict_proba(X_cv)
10         cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_, eps:
11
12     for i in range(len(cv_log_error_array)):
13         print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
14
15
16     best_alpha = np.argmin(cv_log_error_array)
17
18     fig, ax = plt.subplots()
19     ax.plot(alpha, cv_log_error_array,c='g')
20     for i, txt in enumerate(np.round(cv_log_error_array,3)):
21         ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
22     plt.grid()
23     plt.title("Cross Validation Error for each alpha")
24     plt.xlabel("Alpha i's")
25     plt.ylabel("Error measure")
26     plt.show()
27
28     x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
29     x_cfl.fit(X_train,y_train)
30     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
31     sig_clf.fit(X_train, y_train)
32
33     joblib.dump(sig_clf,'XGB_one_param_tuning.joblib')
34 else:
35     sig_clf = joblib.load('XGB_one_param_tuning.joblib')
36
37 predict_y = sig_clf.predict_proba(X_train)
38 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_
39 predict_y = sig_clf.predict_proba(X_cv)
40 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
41 predict_y = sig_clf.predict_proba(X_test)
42 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_lo
43 plot_confusion_matrix(y_test, sig_clf.predict(X_test),'XGB Model (n_estimators = {})'.

```

For values of best alpha = 100 The train log loss is: 0.03208534964008829

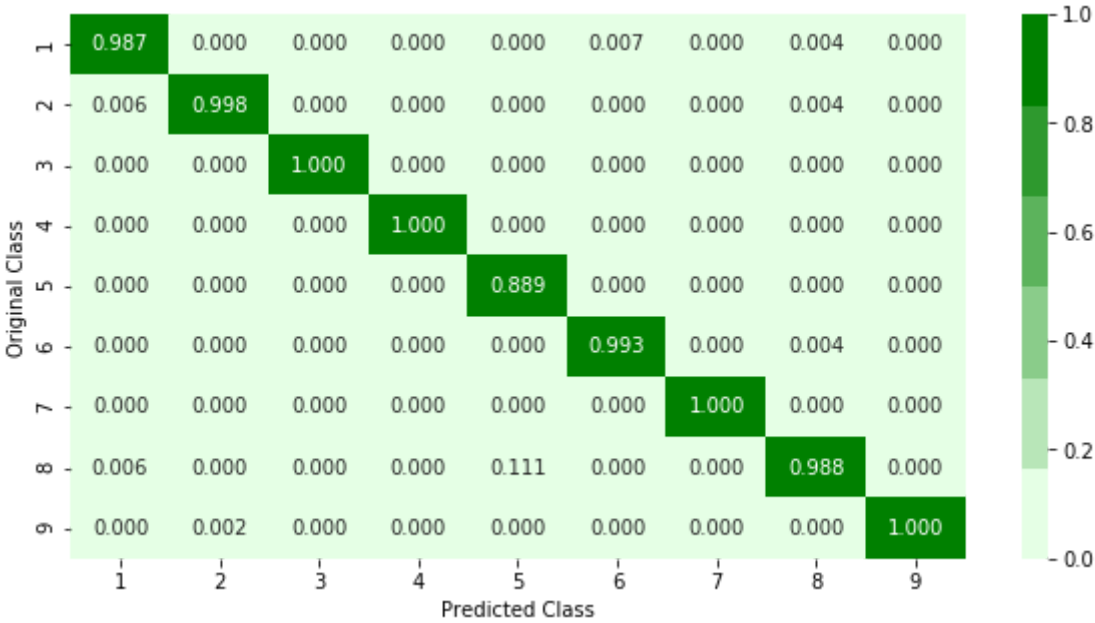
For values of best alpha = 100 The cross validation log loss is: 0.02210324  
001473583

For values of best alpha = 100 The test log loss is: 0.03295365564516063

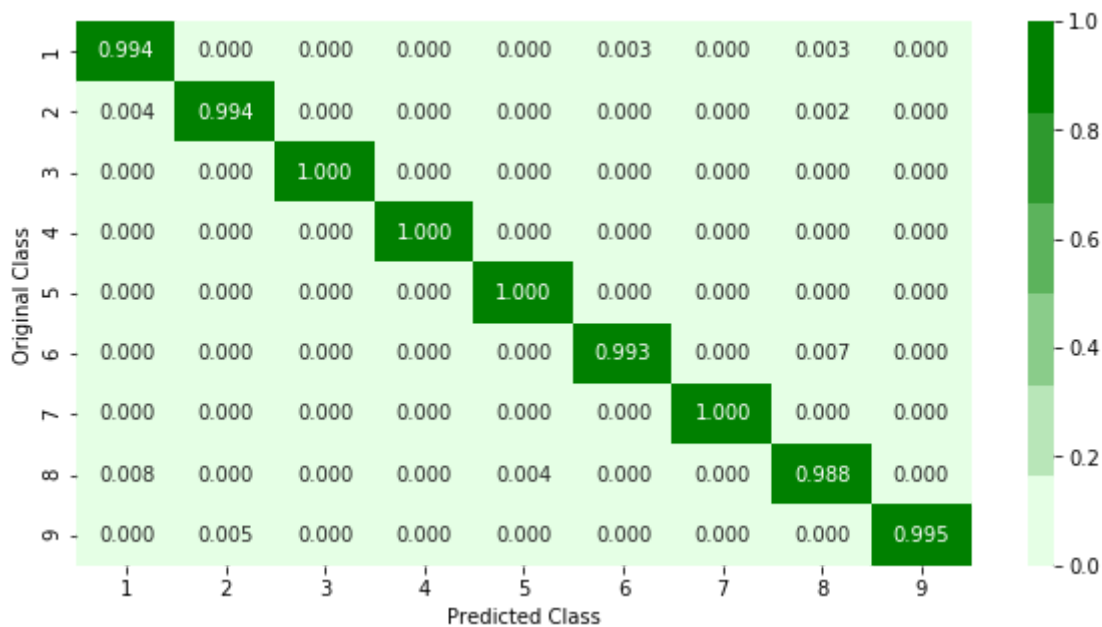
----- Confusion matrix -----  
-----



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



----- Recall matrix -----  
-----



Percentage of misclassified points using XGB Model ( $n_{\text{estimators}} = 100$ ) Mode  
1: 0.46%

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

In [28]:

```
1 # https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost
2 if not os.path.exists('XGB_multiple_param_tuning.joblib'):
3     x_cfl=XGBClassifier()
4
5     prams={
6         'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
7         'n_estimators':[100,200,500,1000,2000],
8         'max_depth':[3,5,10],
9         'colsample_bytree':[0.1,0.3,0.5,1],
10        'subsample':[0.1,0.3,0.5,1]
11    }
12    random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,n_iter=10,verbose=10)
13    random_cfl1.fit(X_train,y_train)
14    joblib.dump(random_cfl1, 'XGB_multiple_param_tuning.joblib')
15 else:
16     random_cfl1 = joblib.load('XGB_multiple_param_tuning.joblib')
```

In [29]:

```
1 random_cfl1.best_params_
```

Out[29]:

```
{'subsample': 1,  
 'n_estimators': 500,  
 'max_depth': 5,  
 'learning_rate': 0.2,  
 'colsample_bytree': 1}
```

In [30]:

```
1 random_cfl1.best_estimator_
```

Out[30]:

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,  
              colsample_bytree=1, gamma=0, learning_rate=0.2, max_delta_step=0,  
              max_depth=5, min_child_weight=1, missing=nan, n_estimators=500,  
              n_jobs=1, nthread=None, objective='multi:softprob', random_state=0,  
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,  
              silent=True, subsample=1)
```

In [31]:

```

1 if not os.path.exists('XGB_multiple_param_tuning_Calibrated.joblib'):
2     c_cfl = CalibratedClassifierCV(random_cfl1.best_estimator_ ,method='sigmoid')
3     c_cfl.fit(X_train,y_train)
4     joblib.dump(c_cfl, 'XGB_multiple_param_tuning_Calibrated.joblib')
5 else:
6     c_cfl = joblib.load('XGB_multiple_param_tuning_Calibrated.joblib')
7
8 predict_y = c_cfl.predict_proba(X_train)
9 print ("The train log loss is:",log_loss(y_train, predict_y))
10 predict_y = c_cfl.predict_proba(X_cv)
11 print("The cross validation log loss is:",log_loss(y_cv, predict_y))
12 predict_y = c_cfl.predict_proba(X_test)
13 print("The test log loss is:",log_loss(y_test, predict_y))
14 plot_confusion_matrix(y_test, c_cfl.predict(X_test),'XGB Model (with more parameters to

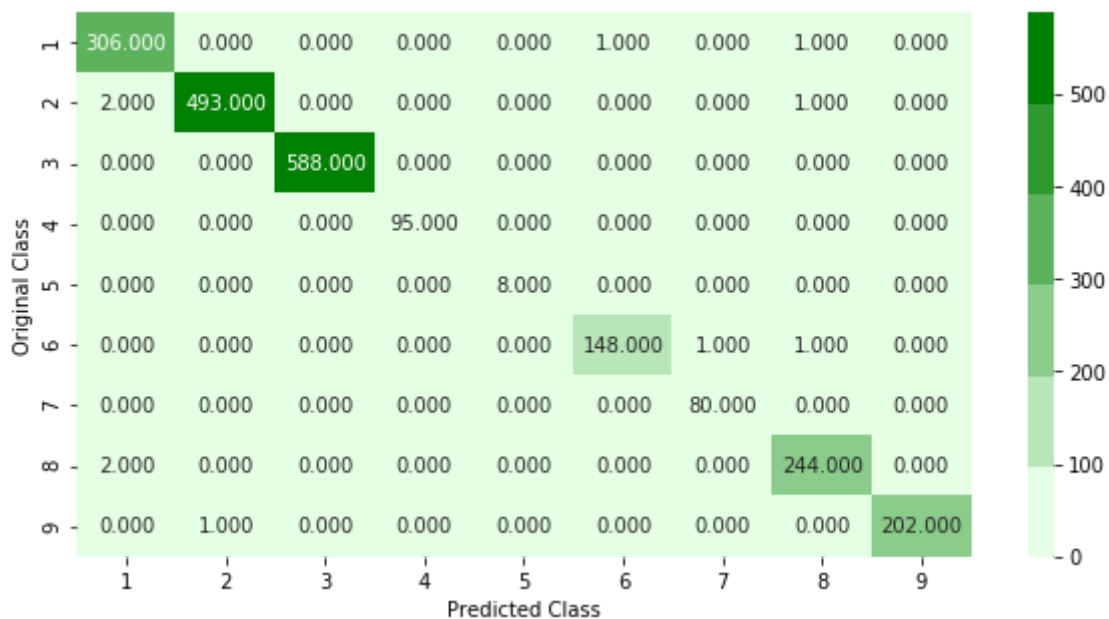
```

The train log loss is: 0.032408503225918525

The cross validation log loss is: 0.02121710172377088

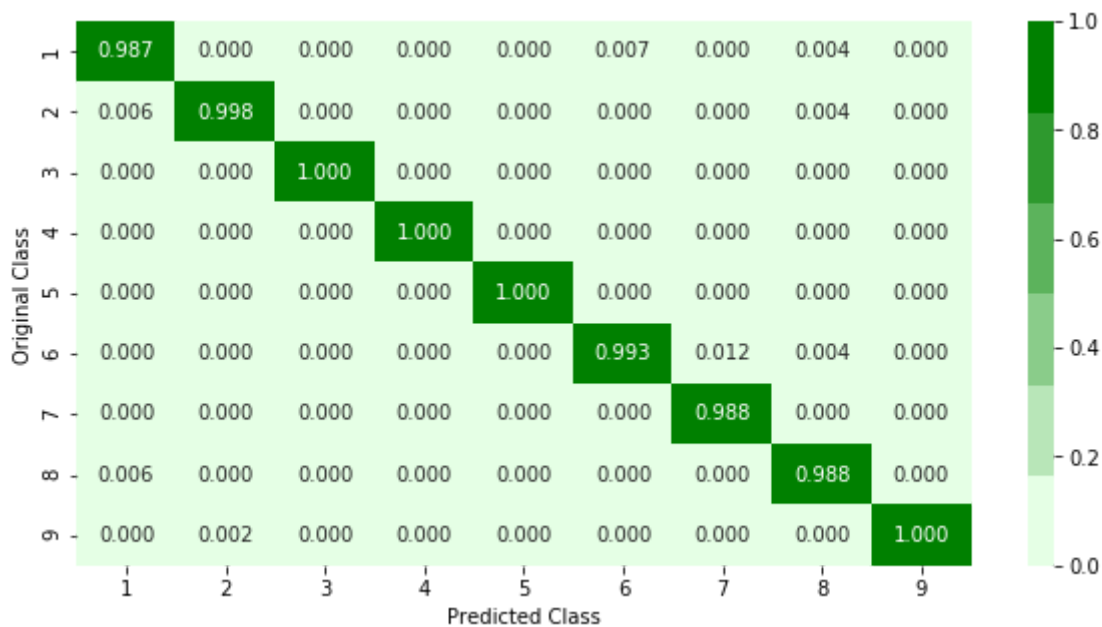
The test log loss is: 0.03332660642651036

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

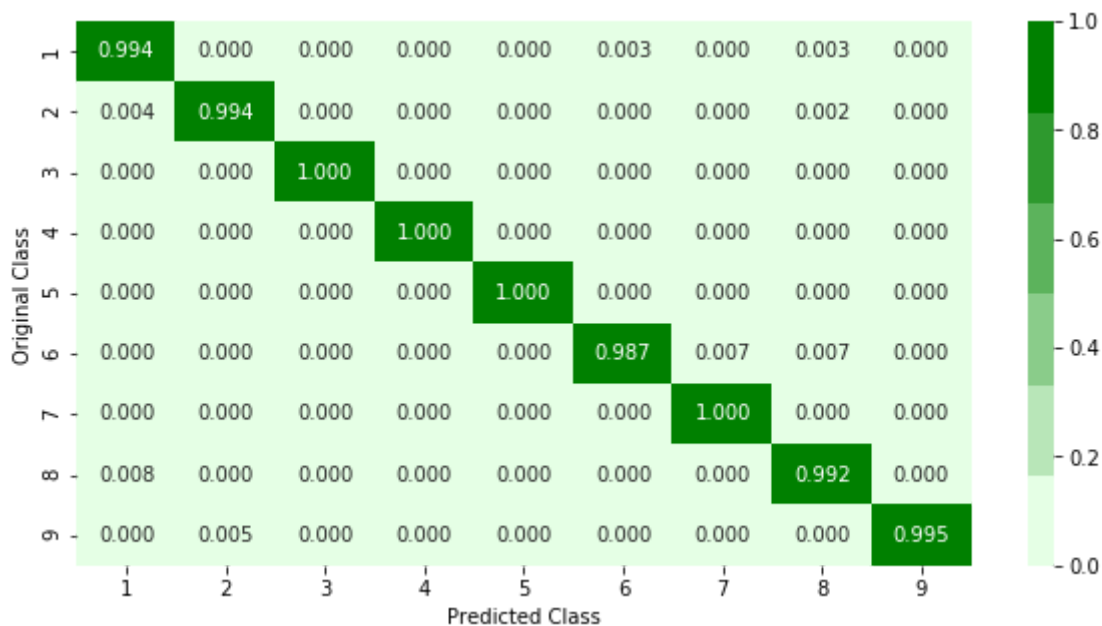


----- Precision matrix -----  
 -----





Recall matrix



Percentage of misclassified points using XGB Model (with more parameters tuned) Model: 0.46%

## 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/videos/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 [32]:

```
1 # Create the directories for parallel processing:
2 folder_1 = 'amsfiles_batch1'
3 folder_2 = 'amsfiles_batch2'
4 folder_3 = 'amsfiles_batch3'
5 folder_4 = 'amsfiles_batch4'
6 folder_5 = 'output'
7 for i in [folder_1, folder_2, folder_3, folder_4, folder_5]:
8     if not os.path.isdir(i):
9         os.makedirs(i)
10 # =====
11 source_dir = 'asmFiles' # change to the asmfiles directory (puth this back to testing i
12 files = os.listdir(source_dir)
13 total_files_per_batch = len(files)//4 # since my computer has 4 cores I create 4 equal
14 total_files_per_batch
15 for destination_dir in [folder_1, folder_2, folder_3, folder_4]:
16     if len(os.listdir(destination_dir)) == 0:
17         print('{} transfer started...'.format(destination_dir))
18         for i in range(total_files_per_batch):
19             file = files.pop(files.index(r.choice(files))) #pick a random file with equ
20             src = os.path.join(source_dir, file)
21             dst = os.path.join(destination_dir, file)
22             shutil.copy2(src=src, dst=dst)
23         print('{} transfer complete'.format(destination_dir))
24     else:
25         print('Files already moved to the destination directory')
```

Files already moved to the destination directory  
Files already moved to the destination directory  
Files already moved to the destination directory  
Files already moved to the destination directory

In [33]:

```

1  #http://flint.cs.yale.edu/cs421/papers/x86-asm/asm.html
2
3  def firstprocess():
4      #The prefixes tells about the segments that are present in the asm files
5      #There are 450 segments(approx) present in all asm files.
6      #this prefixes are best segments that gives us best values.
7      #https://en.wikipedia.org/wiki/Data_segment
8
9      prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:',
10                 '.edata:', '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE']
11      #this are opcodes that are used to get best results
12      #https://en.wikipedia.org/wiki/X86_instruction_listings
13
14      opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub',
15                 'inc', 'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call',
16                 'shl', 'ror', 'rol', 'jnb', 'jz', 'rtn', 'lea', 'movzx']
17      #best keywords that are taken from different blogs
18      keywords = ['.dll', 'std:', ':dword']
19      #Below taken registers are general purpose registers and special registers
20      #All the registers which are taken are best
21      registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
22
23      #symbols
24      symbols=['[', '-', '+', '*', '?', '@', ']']
25
26
27      asmfile1 = open(os.path.join(folder_5, 'asmfile1.txt'), "w+") # open the file stream
28
29      files = os.listdir(folder_1)
30
31      for file in files:
32
33          #filling the values with zeros into the arrays
34          prefixescount = np.zeros(len(prefixes), dtype=int)
35          opcodescount = np.zeros(len(opcodes), dtype=int)
36          keywordcount = np.zeros(len(keywords), dtype=int)
37          registerscount = np.zeros(len(registers), dtype=int)
38          symbolscount = np.zeros(len(symbols), dtype=int)
39
40          features = list()
41
42          filename = file.split('.')[0]
43
44          asmfile1.write(filename + ",")
45
46          # https://docs.python.org/3/Library/codecs.html#codecs.ignore_errors
47          # https://docs.python.org/3/Library/codecs.html#codecs.Codec.encode
48
49          with codecs.open(os.path.join(folder_1, file), encoding='cp1252', errors='replac
50                          for lines in fli:
51                              # https://www.tutorialspoint.com/python3/string_rstrip.htm
52
53                              line = lines.rstrip().split()
54                              l = line[0]
55                              #counting the prefixes in each and every line
56                              for i in range(len(prefixes)):
57                                  if prefixes[i] in line[0]:

```

```

58         prefixescount[i]+=1
59
60     line = line[1:]
61     #counting the opcodes in each and every line
62     for i in range(len(opcodes)):
63         if any(opcodes[i] == li for li in line):
64             features.append(opcodes[i])
65             opcodescount[i]+=1
66
67     #counting registers in the line
68     for i in range(len(registers)):
69         for li in line:
70             # we will use registers only in 'text' and 'CODE' segments
71             if registers[i] in li and ('text' in l or 'CODE' in l):
72                 registerscount[i]+=1
73
74     #counting keywords in the line
75     for i in range(len(keywords)):
76         for li in line:
77             if keywords[i] in li:
78                 keywordcount[i]+=1
79
80     for i in range(len(symbols)):
81         for item in line:
82             if symbols[i] in item:
83                 symbolscount[i]+=1
84
85
86     #pushing the values into the file after reading whole file
87     for prefix in prefixescount:
88         asmfile1.write(str(prefix)+",")
89     for opcode in opcodescount:
90         asmfile1.write(str(opcode)+",")
91     for register in registerscount:
92         asmfile1.write(str(register)+",")
93     for key in keywordcount:
94         asmfile1.write(str(key)+",")
95     for symbol in symbolscount:
96         asmfile1.write(str(symbol)+",")
97     asmfile1.write("\n")
98     asmfile1.close()
99
100 #same as above
101 def secondprocess():
102     #The prefixes tells about the segments that are present in the asm files
103     #There are 450 segments(approx) present in all asm files.
104     #this prefixes are best segments that gives us best values.
105     #https://en.wikipedia.org/wiki/Data_segment
106
107     prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:',
108                '.edata:', '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE']
109     #this are opcodes that are used to get best results
110     #https://en.wikipedia.org/wiki/X86_instruction_listings
111
112     opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub',
113                'inc', 'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call',
114                'shl', 'ror', 'rol', 'jnb', 'jz', 'rtn', 'lea', 'movzx']
115     #best keywords that are taken from different blogs
116     keywords = ['.dll', 'std:', ':dword']
117     #Below taken registers are general purpose registers and special registers
118     #All the registers which are taken are best

```

```

119 registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
120
121 #symbols
122 symbols=['[','-','+','*','?','@','']
123
124 asmfile2 = open(os.path.join(folder_5,'asmfile2.txt'),'w+') # open the file stream
125
126 files = os.listdir(folder_2)
127
128 for file in files:
129
130     #filling the values with zeros into the arrays
131     prefixescount = np.zeros(len(prefixes),dtype=int)
132     opcodescount = np.zeros(len(opcodes),dtype=int)
133     keywordcount = np.zeros(len(keywords),dtype=int)
134     registerscount = np.zeros(len(registers),dtype=int)
135     symbolscount = np.zeros(len(symbols),dtype=int)
136
137     features = list()
138
139     filename = file.split('.')[0]
140
141     asmfile2.write(filename + ",")
142
143     # https://docs.python.org/3/Library/codecs.html#codecs.ignore_errors
144     # https://docs.python.org/3/Library/codecs.html#codecs.Codec.encode
145
146     with codecs.open(os.path.join(folder_2,file),encoding='cp1252',errors='replac
147         for lines in fli:
148             # https://www.tutorialspoint.com/python3/string_rstrip.htm
149
150             line = lines.rstrip().split()
151             l = line[0]
152             #counting the prefixs in each and every line
153             for i in range(len(prefixes)):
154                 if prefixes[i] in line[0]:
155                     prefixescount[i]+=1
156
157             line = line[1:]
158             #counting the opcodes in each and every line
159             for i in range(len(opcodes)):
160                 if any(opcodes[i] == li for li in line):
161                     features.append(opcodes[i])
162                     opcodescount[i]+=1
163
164             #counting registers in the line
165             for i in range(len(registers)):
166                 for li in line:
167                     # we will use registers only in 'text' and 'CODE' segments
168                     if registers[i] in li and ('text' in l or 'CODE' in l):
169                         registerscount[i]+=1
170
171             #counting keywords in the line
172             for i in range(len(keywords)):
173                 for li in line:
174                     if keywords[i] in li:
175                         keywordcount[i]+=1
176
177             for i in range(len(symbols)):
178                 for item in line:
179                     if symbols[i] in item:

```

```

180         symbolscount[i]+=1
181
182     #pushing the values into the file after reading whole file
183     for prefix in prefixescount:
184         asmfile2.write(str(prefix)+",")
185     for opcode in opcodescount:
186         asmfile2.write(str(opcode)+",")
187     for register in registerscount:
188         asmfile2.write(str(register)+",")
189     for key in keywordcount:
190         asmfile2.write(str(key)+",")
191     for symbol in symbolscount:
192         asmfile2.write(str(symbol)+",")
193     asmfile2.write("\n")
194     asmfile2.close()
195
196 # same as above
197 def thirdprocess():
198     #The prefixes tells about the segments that are present in the asm files
199     #There are 450 segments(approx) present in all asm files.
200     #this prefixes are best segments that gives us best values.
201     #https://en.wikipedia.org/wiki/Data_segment
202
203     prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:',
204                '.edata:', '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE']
205     #this are opcodes that are used to get best results
206     #https://en.wikipedia.org/wiki/X86_instruction_Listings
207
208     opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub',
209               'inc', 'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call',
210               'shl', 'ror', 'rol', 'jnb', 'jz', 'rtn', 'lea', 'movzx']
211     #best keywords that are taken from different blogs
212     keywords = ['.dll', 'std:', ':dword']
213     #Below taken registers are general purpose registers and special registers
214     #All the registers which are taken are best
215     registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
216
217     #symbols
218     symbols=['[', '-', '+', '*', '?', '@', ']']
219
220     asmfile3 = open(os.path.join(folder_5, 'asmfile3.txt'), "w+") # open the file stream
221
222     files = os.listdir(folder_3)
223
224     for file in files:
225
226         #filling the values with zeros into the arrays
227         prefixescount = np.zeros(len(prefixes), dtype=int)
228         opcodescount = np.zeros(len(opcodes), dtype=int)
229         keywordcount = np.zeros(len(keywords), dtype=int)
230         registerscount = np.zeros(len(registers), dtype=int)
231         symbolscount = np.zeros(len(symbols), dtype=int)
232
233         features = list()
234
235         filename = file.split('.')[0]
236
237         asmfile3.write(filename + ",")
238
239         # https://docs.python.org/3/Library/codecs.html#codecs.ignore_errors
240         # https://docs.python.org/3/Library/codecs.html#codecs.Codec.encode

```

```

241
242 with codecs.open(os.path.join(folder_3,file),encoding='cp1252',errors ='replac
243     for lines in fli:
244         # https://www.tutorialspoint.com/python3/string\_rstrip.htm
245
246         line = lines.rstrip().split()
247         l = line[0]
248         #counting the prefixes in each and every line
249         for i in range(len(prefixes)):
250             if prefixes[i] in line[0]:
251                 prefixescount[i]+=1
252
253         line = line[1:]
254         #counting the opcodes in each and every line
255         for i in range(len(opcodes)):
256             if any(opcodes[i] == li for li in line):
257                 features.append(opcodes[i])
258                 opcodescount[i]+=1
259
260         #counting registers in the line
261         for i in range(len(registers)):
262             for li in line:
263                 # we will use registers only in 'text' and 'CODE' segments
264                 if registers[i] in li and ('text' in l or 'CODE' in l):
265                     registerscount[i]+=1
266
267         #counting keywords in the line
268         for i in range(len(keywords)):
269             for li in line:
270                 if keywords[i] in li:
271                     keywordcount[i]+=1
272
273         for i in range(len(symbols)):
274             for item in line:
275                 if symbols[i] in item:
276                     symbolscount[i]+=1
277
278         #pushing the values into the file after reading whole file
279         for prefix in prefixescount:
280             asmfile3.write(str(prefix)+",")
281         for opcode in opcodescount:
282             asmfile3.write(str(opcode)+",")
283         for register in registerscount:
284             asmfile3.write(str(register)+",")
285         for key in keywordcount:
286             asmfile3.write(str(key)+",")
287         for symbol in symbolscount:
288             asmfile3.write(str(symbol)+",")
289         asmfile3.write("\n")
290     asmfile3.close()
291
292 # same as above
293 def fourthprocess():
294     #The prefixes tells about the segments that are present in the asm files
295     #There are 450 segments(approx) present in all asm files.
296     #this prefixes are best segments that gives us best values.
297     #https://en.wikipedia.org/wiki/Data\_segment
298
299     prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:',
300                '.edata:', '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE']
301     #this are opcodes that are used to get best results

```



```

302 #https://en.wikipedia.org/wiki/X86_instruction_listings
303
304 opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'ret', 'nop', 'sub',
305            'inc', 'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call',
306            'shl', 'ror', 'rol', 'jnb', 'jz', 'rtn', 'lea', 'movzx']
307 #best keywords that are taken from different blogs
308 keywords = ['.dll', 'std:', ':dword']
309 #Below taken registers are general purpose registers and special registers
310 #All the registers which are taken are best
311 registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
312
313 #symbols
314 symbols=['[', '-', '+', '*', '?', '@', ',']
315
316 asmfile4 = open(os.path.join(folder_5, 'asmfile4.txt'), "w+") # open the file stream
317
318 files = os.listdir(folder_4)
319
320 for file in files:
321
322     #filling the values with zeros into the arrays
323     prefixescount = np.zeros(len(prefixes), dtype=int)
324     opcodescount = np.zeros(len(opcodes), dtype=int)
325     keywordcount = np.zeros(len(keywords), dtype=int)
326     registerscount = np.zeros(len(registers), dtype=int)
327     symbolscount = np.zeros(len(symbols), dtype=int)
328
329     features = list()
330
331     filename = file.split('.')[0]
332
333     asmfile4.write(filename + ",")
334
335     # https://docs.python.org/3/Library/codecs.html#codecs.ignore_errors
336     # https://docs.python.org/3/Library/codecs.html#codecs.Codec.encode
337
338     with codecs.open(os.path.join(folder_4, file), encoding='cp1252', errors='replac
339         for lines in fli:
340             # https://www.tutorialspoint.com/python3/string_rstrip.htm
341
342             line = lines.rstrip().split()
343             l = line[0]
344             #counting the prefixs in each and every line
345             for i in range(len(prefixes)):
346                 if prefixes[i] in line[0]:
347                     prefixescount[i] += 1
348
349             line = line[1:]
350             #counting the opcodes in each and every line
351             for i in range(len(opcodes)):
352                 if any(opcodes[i] == li for li in line):
353                     features.append(opcodes[i])
354                     opcodescount[i] += 1
355
356             #counting registers in the line
357             for i in range(len(registers)):
358                 for li in line:
359                     # we will use registers only in 'text' and 'CODE' segments
360                     if registers[i] in li and ('text' in l or 'CODE' in l):
361                         registerscount[i] += 1
362

```



```
363         #counting keywords in the line
364         for i in range(len(keywords)):
365             for li in line:
366                 if keywords[i] in li:
367                     keywordcount[i]+=1
368
369         for i in range(len(symbols)):
370             for item in line:
371                 if symbols[i] in item:
372                     symbolscount[i]+=1
373
374     #pushing the values into the file after reading whole file
375     for prefix in prefixescount:
376         asmfile4.write(str(prefix)+",")
377     for opcode in opcodescount:
378         asmfile4.write(str(opcode)+",")
379     for register in registerscount:
380         asmfile4.write(str(register)+",")
381     for key in keywordcount:
382         asmfile4.write(str(key)+",")
383     for symbol in symbolscount:
384         asmfile4.write(str(symbol)+",")
385     asmfile4.write("\n")
386     asmfile4.close()
387
388 # def main():
389 #     #the below code is used for multiprocessing
390 #     #the number of process depends upon the number of cores present System
391 #     #process is used to call multiprocessing
392 #     manager=multiprocessing.Manager()
393
394 #     p1=mp.Process(target=firstprocess)
395 #     p2=mp.Process(target=secondprocess)
396 #     p3=mp.Process(target=thirdprocess)
397 #     p4=mp.Process(target=fourthprocess)
398
399 #     #p1.start() is used to start the thread execution
400 #     p1.start()
401 #     p2.start()
402 #     p3.start()
403 #     p4.start()
404 #     #After completion all the threads are joined
405 #     p1.join()
406 #     p2.join()
407 #     p3.join()
408 #     p4.join()
409
410 def main():
411     try:
412         print('feature extraction starting...')
413         firstprocess()
414         print('batch 1 processing done')
415         secondprocess()
416         print('batch 2 processing done')
417         thirdprocess()
418         print('batch 3 processing done')
419         fourthprocess()
420         print('process complete')
421     except:
422         print('Error')
```

In [34]:

```

1 if len(os.listdir(folder_5)) == 0:
2     start = datetime.now()
3     main()
4     print(datetime.now() - start)
5 else:
6     print('Output files already created')

```

Output files already created

In [35]:

```

1 if not os.path.exists('asmoutputfile.csv'):
2     files = os.listdir('output')
3     asmoutputfile = pd.DataFrame()
4     for file in files:
5         asmoutputfile = pd.concat(objs=[asmoutputfile,pd.read_csv(os.path.join('output',
6         asmoutputfile.columns = ['ID', 'HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:',
7         '.CODE', 'jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'ret',
8         'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'j',
9         'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip', 'extra', '']
10    asmoutputfile.drop(columns='extra', inplace=True)
11    asmoutputfile.to_csv('asmoutputfile.csv', header=True, index=False)

```

In [36]:

```

1 # asmoutputfile.csv(output generated from the above two cells) will contain all the ex
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.drop(columns=']', inplace=True)
7 result_asm.head()

```

Out[36]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc
0	01kcPWA9K2BOxQeS5Rju	19	744	0	127	57	0	323	0	:
1	02IOCvYEy8mjiuAQHax3	17	838	0	41	77840	0	0	0	:
2	04EjldbPV5e1XroFOpiN	19	60476	0	349	3760	0	0	0	:
3	05EeG39MTRrl6VY21DPd	17	11119	0	323	1047	0	3385	0	:
4	05LHG8fR3iPn6aglo9z7	17	0	0	178	0	0	511	0	:

5 rows × 59 columns

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

```

1 #file sizes of byte files
2 files=os.listdir('asmFiles')
3 filenames=Y['ID'].tolist()
4 class_y=Y['Class'].tolist()
5 class_bytes=[]
6 sizebytes=[]
7 num_lines = []
8 fnames=[]
9 for file in tqdm(files):
10     statinfo=os.stat('asmFiles\\'+file)
11     with codecs.open('asmFiles\\'+file,'r',encoding='cp1252',errors='replace') as con
12         num_of_lines = len(content.readlines())
13         # split the file name at '.' and take the first part of it i.e the file name
14         file=file.split('.')[0]
15         if any(file == filename for filename in filenames):
16             i=filenames.index(file)
17             class_bytes.append(class_y[i])
18             # converting into Mb's
19             sizebytes.append(statinfo.st_size/(1024.0*1024.0))
20             fnames.append(file)
21             num_lines.append(num_of_lines)
22 asm_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes, 'filelines':num_lines, 'Clas
23 print (asm_size_byte.head())

```

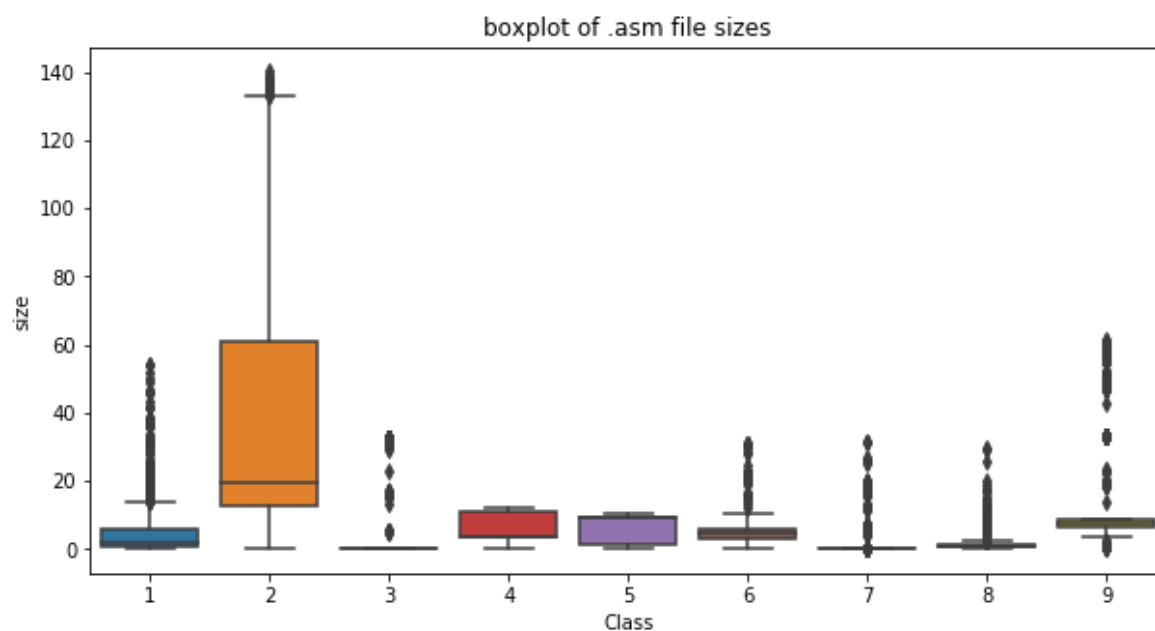
```
100%|███████████| 
10868/10868 [34:15<00:00, 14.27it/s]
```

	ID	size	filelines	Class
0	01azqd4InC7m9JpocGv5	56.229886	1392154	9
1	01IsoiSMh5gxyDYTL4CB	13.999378	161528	2
2	01jsnpXSAlgW6aPeDxrU	8.507785	70960	9
3	01kcPWA9K2B0xQeS5Rju	0.078190	1276	1
4	01SuzwMJEIXsK7A8dQb1	0.996723	15282	8

#### 4.2.1.2.1 Distribution of .asm file sizes

In [38]:

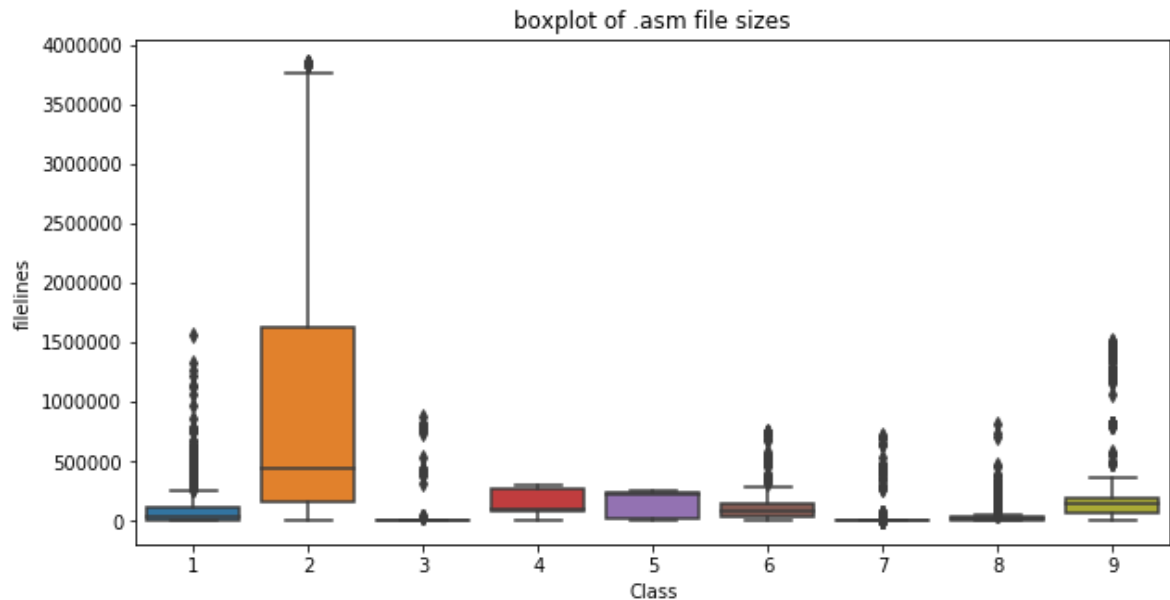
```
1 #boxplot of asm files
2 plt.figure(figsize = (10,5))
3 ax = sns.boxplot(x="Class", y="size", data=asm_size_byte)
4 plt.title("boxplot of .asm file sizes")
5 plt.show()
```



#### 4.2.1.2.2 Distribution of file row numbers of .asm files

In [39]:

```
1 #boxplot of asm files
2 plt.figure(figsize = (10,5))
3 ax = sns.boxplot(x="Class", y="filelines", data=asm_size_byte)
4 plt.title("boxplot of .asm file sizes")
5 plt.show()
```



In [40]:

```
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='ID', how='left')
5 result_asm['.dataweighted'] = result_asm['.data:']*4
6 result_asm.head()
```

(10868, 59)  
(10868, 4)

Out[40]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc
0	01kcPWA9K2BOxQeS5Rju	19	744	0	127	57	0	323	0	:
1	02lOCvYEy8mjiuAQHax3	17	838	0	41	77840	0	0	0	:
2	04EjldbPV5e1XroFOpiN	19	60476	0	349	3760	0	0	0	:
3	05EeG39MTRrl6VY21DPd	17	11119	0	323	1047	0	3385	0	:
4	05LHG8fR3iPn6aglo9z7	17	0	0	178	0	0	511	0	:

5 rows × 62 columns



In [41]:

```
1 # we normalize the data each column
2 result_asm = normalize(result_asm)
3 result_asm.head()
```

Out[41]:

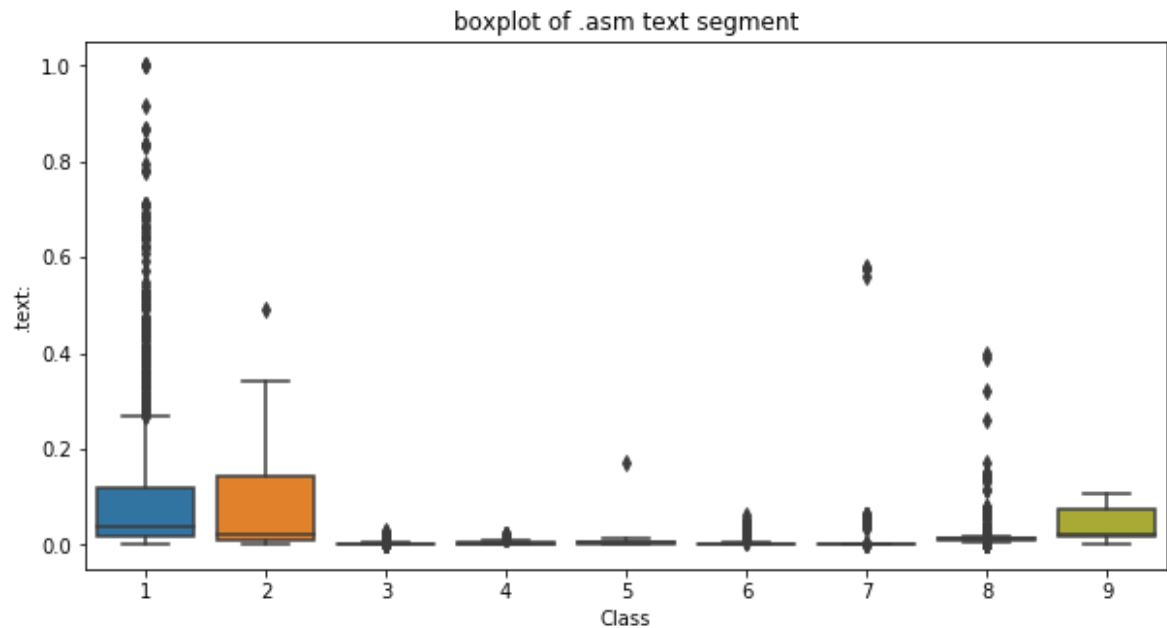
	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.ec
0	01kcPWA9K2BOxQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084	
1	02IOCvYEy8mjiuAQHax3	0.096045	0.001230	0.0	0.000246	0.030916	0.0	0.000000	
2	04EjldbPV5e1XroFOpiN	0.107345	0.088779	0.0	0.002091	0.001493	0.0	0.000000	
3	05EeG39MTRrI6VY21DPd	0.096045	0.016323	0.0	0.001935	0.000416	0.0	0.000882	
4	05LHG8fR3iPn6aglo9z7	0.096045	0.000000	0.0	0.001066	0.000000	0.0	0.000133	

5 rows × 62 columns

4.2.2 Univariate analysis on asm file features

In [42]:

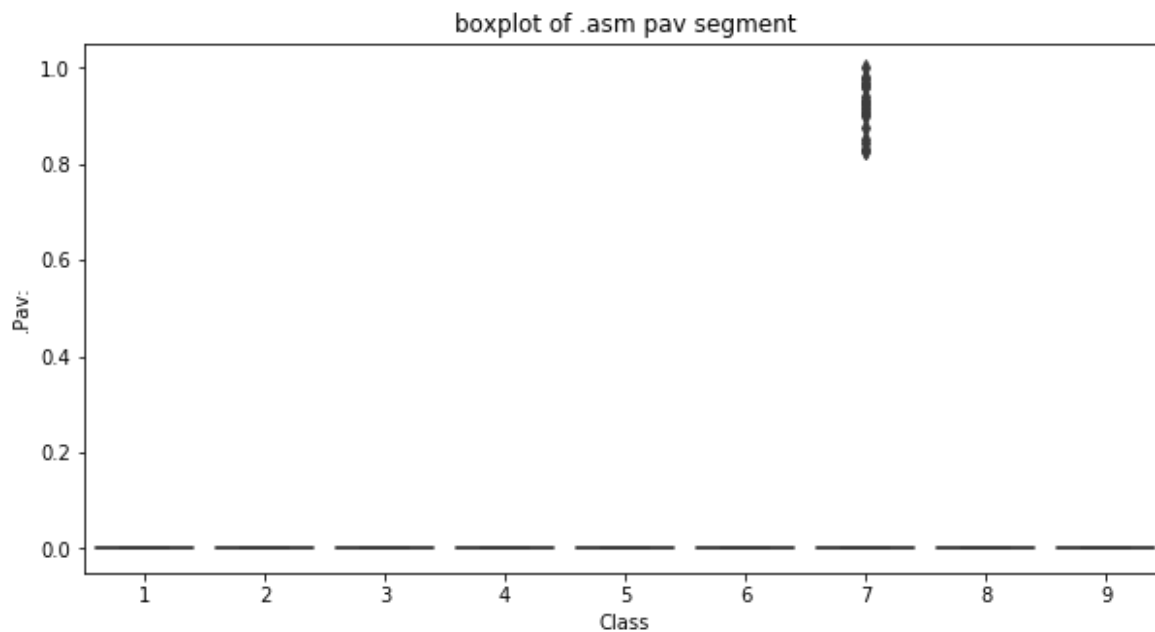
```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y=".text:", data=result_asm)
3 plt.title("boxplot of .asm text segment")
4 plt.show()
```



The plot is between Text and class  
Class 1,2 and 9 can be easily separated

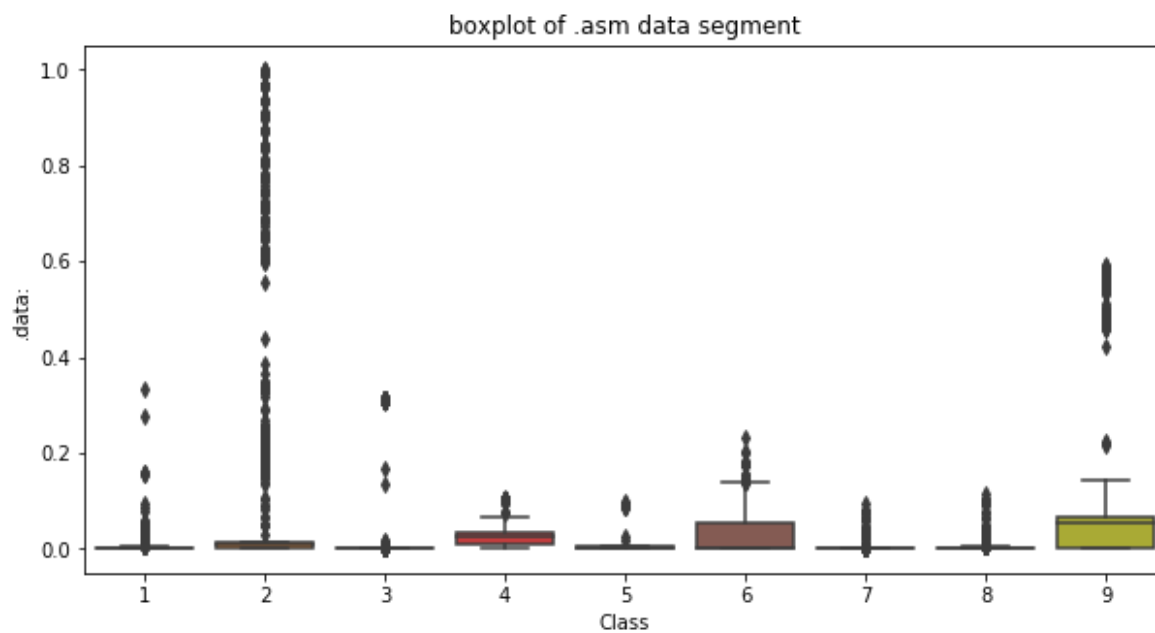
In [43]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y=".Pav:", data=result_asm)
3 plt.title("boxplot of .asm pav segment")
4 plt.show()
```



In [44]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y=".data:", data=result_asm)
3 plt.title("boxplot of .asm data segment")
4 plt.show()
```

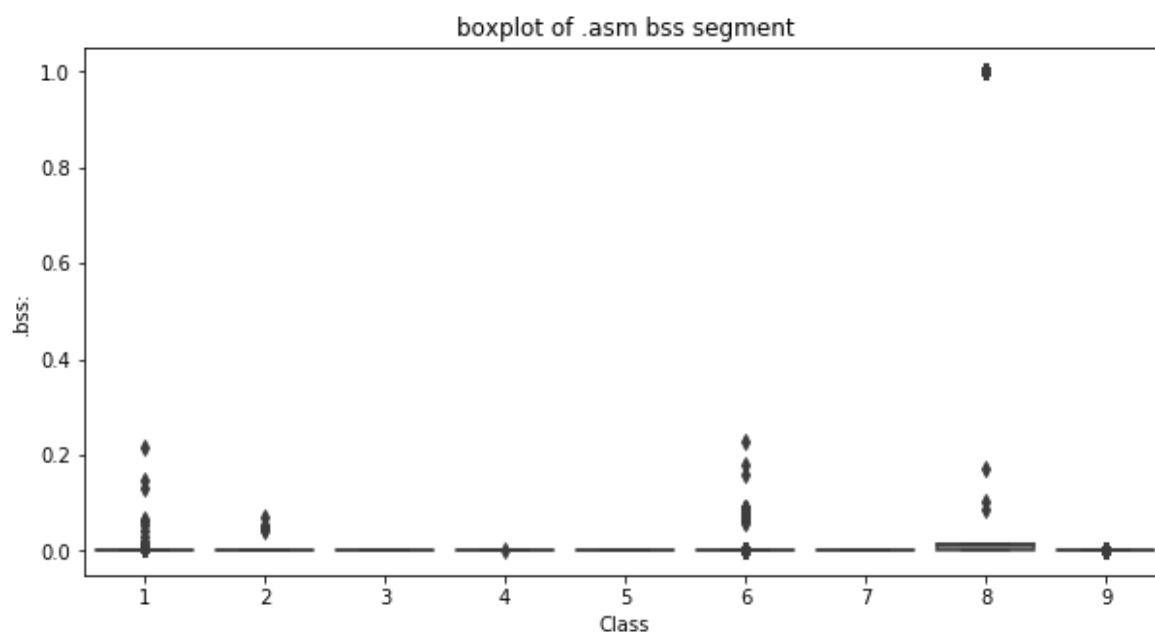


The plot is between data segment and class label

class 6 and class 9 can be easily separated from given points

In [45]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y=".bss:", data=result_asm)
3 plt.title("boxplot of .asm bss segment")
4 plt.show()
```



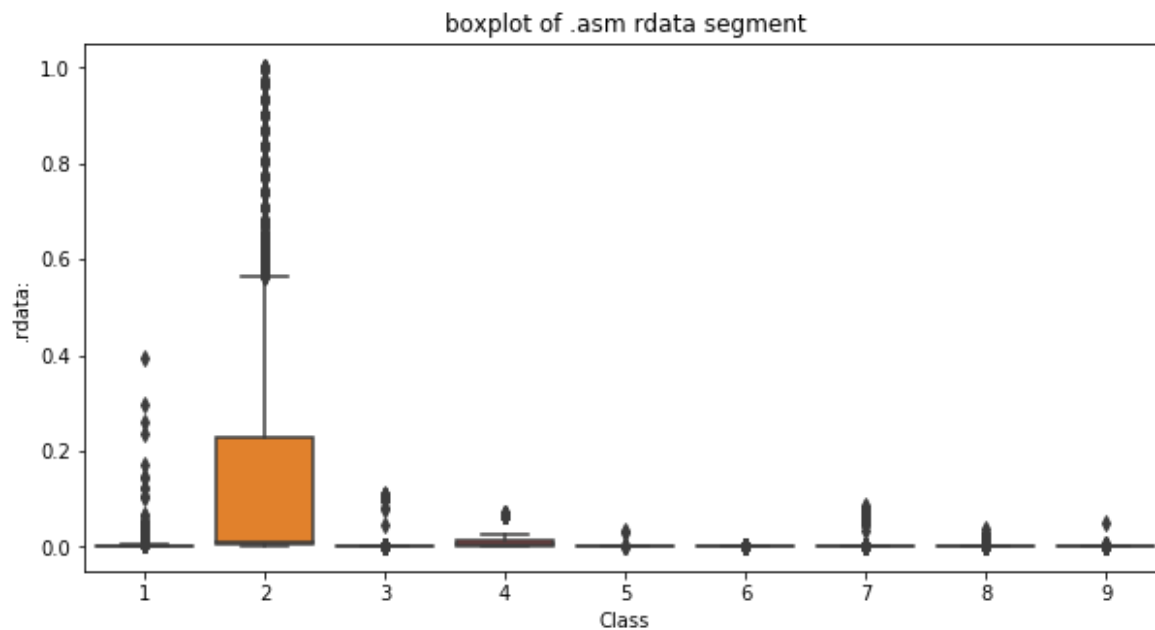
plot between bss segment and class label

very less number of files are having bss segment



In [46]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y=".rdata:", data=result_asm)
3 plt.title("boxplot of .asm rdata segment")
4 plt.show()
```

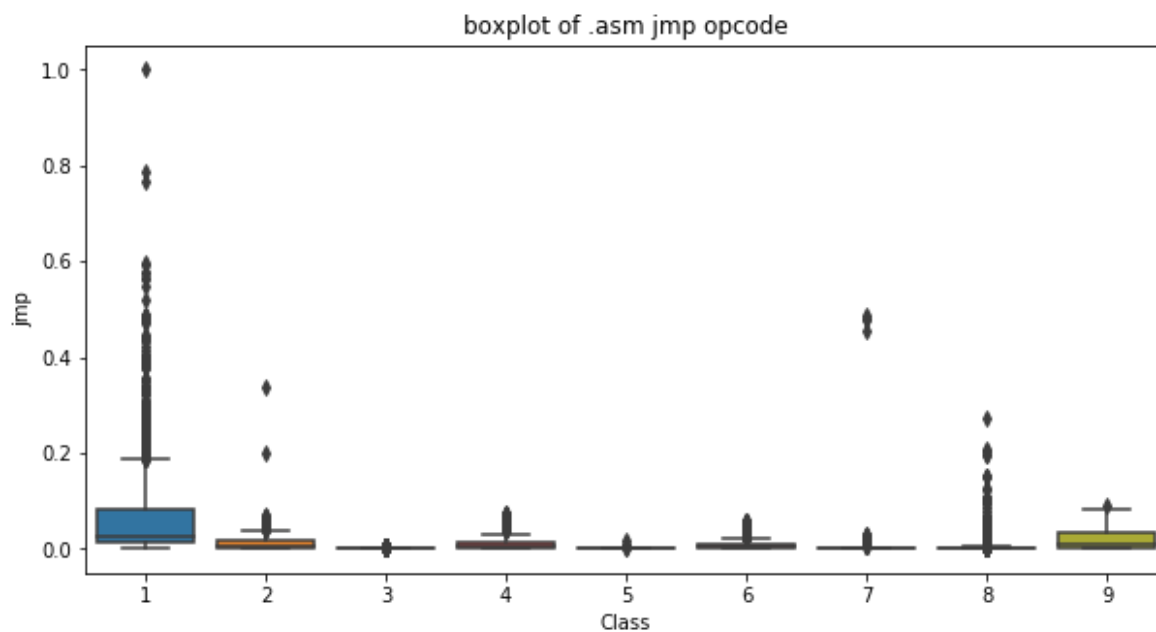


Plot between rdata segment and Class segment

Class 2 can be easily separated 75 percentile files are having 1M rdata lines

In [47]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y="jmp", data=result_asm)
3 plt.title("boxplot of .asm jmp opcode")
4 plt.show()
```

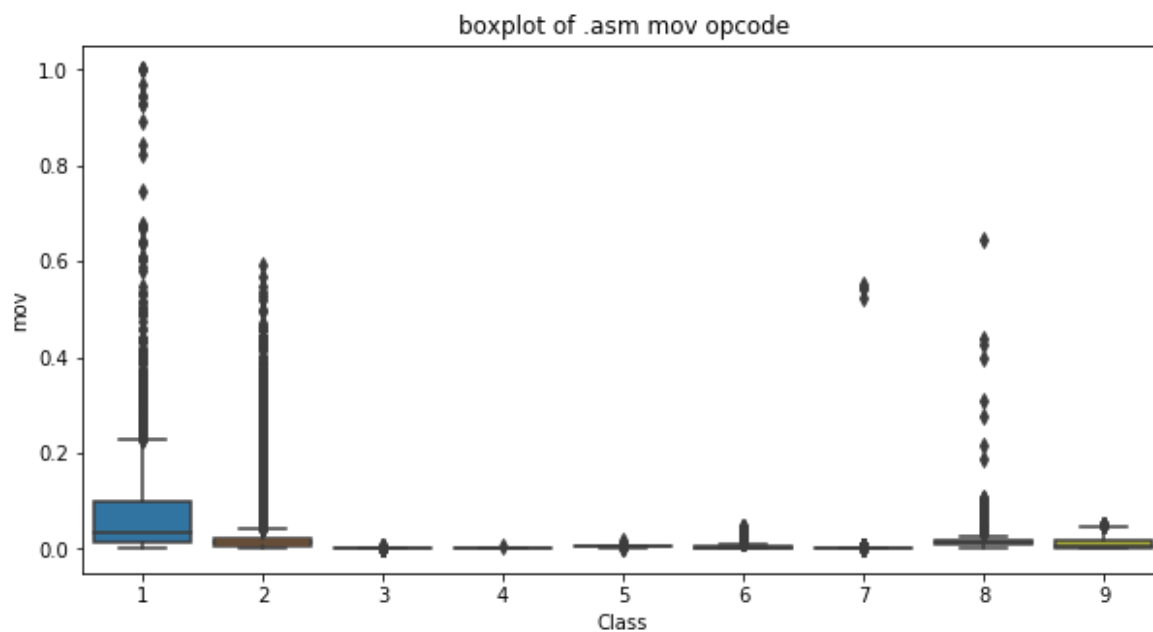


plot between jmp and Class label

Class 1 is having frequency of 2000 approx in 75 percentile of files

In [48]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y="mov", data=result_asm)
3 plt.title("boxplot of .asm mov opcode")
4 plt.show()
```

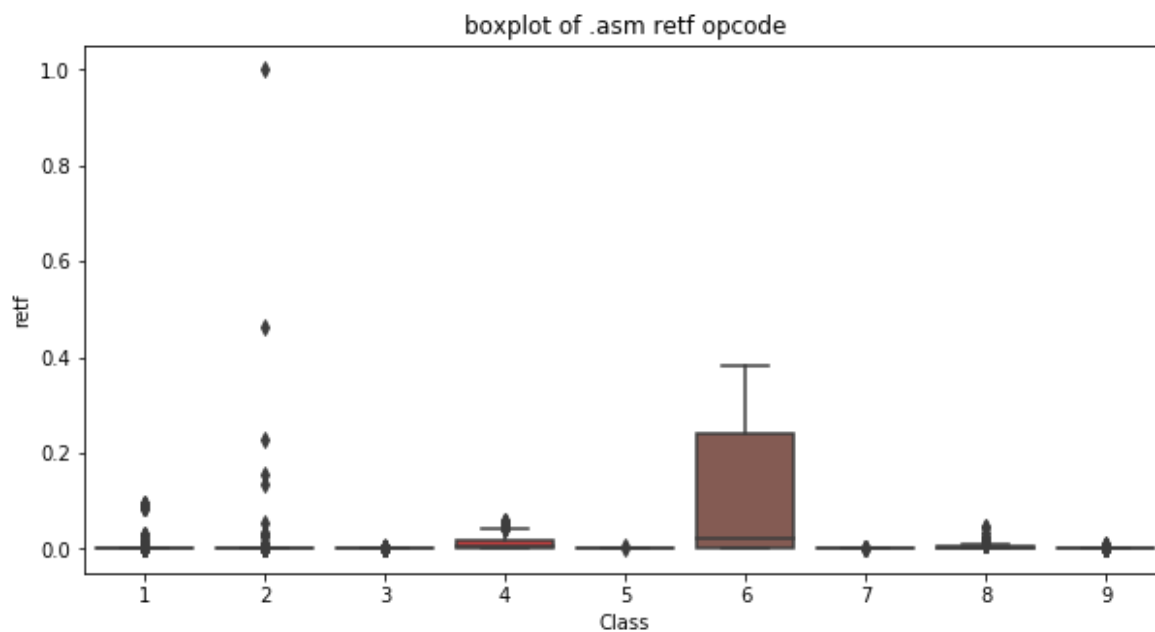


plot between Class label and mov opcode

Class 1 is having frequency of 2000 approx in 75 perentile of files

In [49]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y="retf", data=result_asm)
3 plt.title("boxplot of .asm retf opcode")
4 plt.show()
```



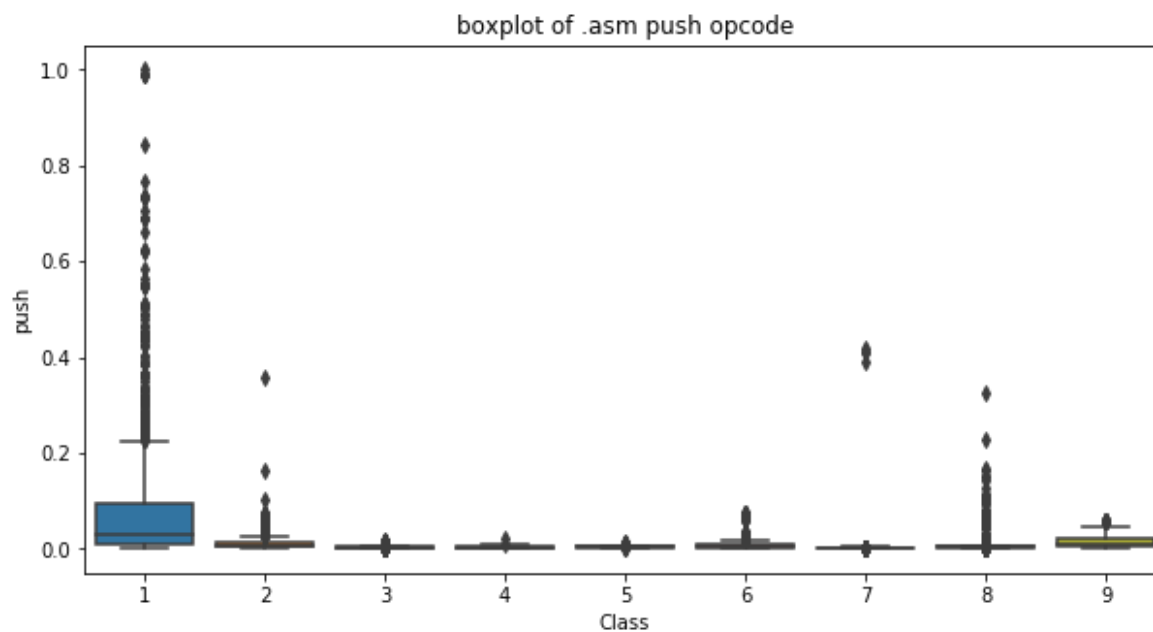
plot between Class label and retf

Class 6 can be easily separated with opcode retf

The frequency of retf is approx of 250.

In [50]:

```
1 plt.figure(figsize = (10,5))
2 ax = sns.boxplot(x="Class", y="push", data=result_asm)
3 plt.title("boxplot of .asm push opcode")
4 plt.show()
```



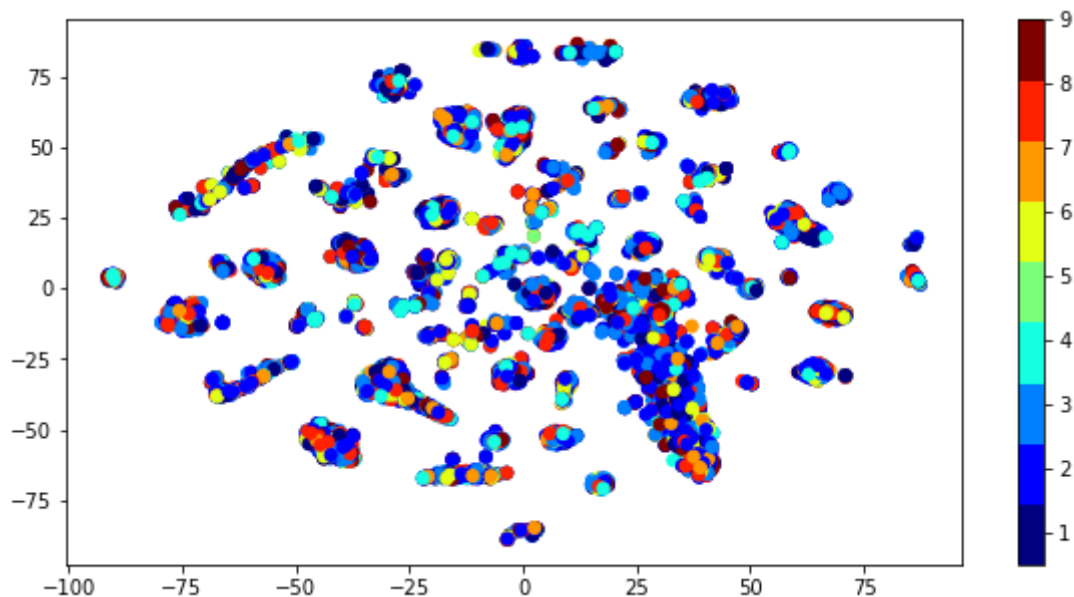
plot between push opcode and Class label

Class 1 is having 75 percentile files with push opcodes of frequency 1000

## 4.2.2 Multivariate Analysis on .asm file features

In [51]:

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

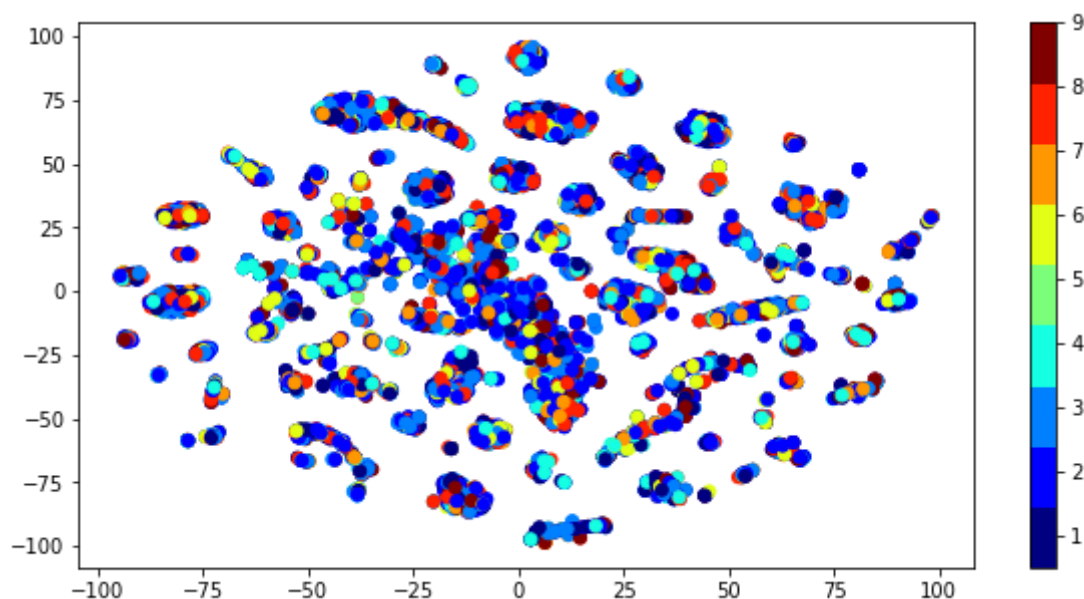


In [52]:

```

1 # by univariate analysis on the .asm file features we are getting very negligible info
2 # 'rtn', '.BSS:', '.CODE' features, so here we are trying multivariate analysis after re
3 # the plot looks very messy
4 plt.figure(figsize = (10,5))
5 xtsne=TSNE(perplexity=30)
6 results=xtsne.fit_transform(result_asm.drop(['ID','Class', 'rtn', '.BSS:', '.CODE','si
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()

```



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 [53]:

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

In [54]:

```
1 asm_x.rename(columns = {'[': 'bracket', '-': 'minussign', '+': 'plussign', '*': 'mulsign', '?'
```



In [55]:

```
1 X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x, asm_y, strat:  
2 X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_a:
```



In [56]:

```
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
bracket	False
minussign	False
plussign	False
mulsign	False
questionmark	False
atsign	False

```
size          False
filelines     False
.dataweighted False
dtype: bool
```

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

### 4.4.1 K-Nearest Neighbors

In [57]:

```

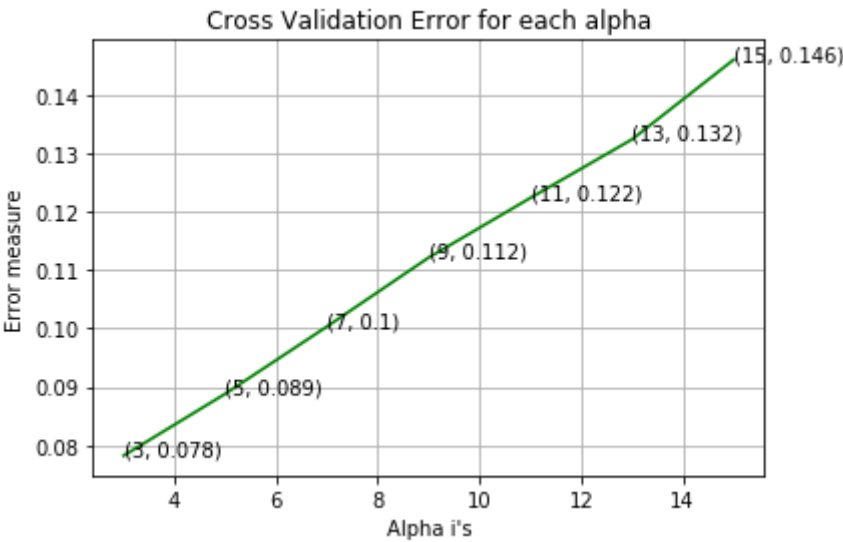
1  alpha = [x for x in [3,5,7,9,11,13,15]]
2  cv_log_error_array=[]
3  for i in alpha:
4      k_cfl=KNeighborsClassifier(n_neighbors=i)
5      k_cfl.fit(X_train_asm,y_train_asm)
6      sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
7      sig_clf.fit(X_train_asm,y_train_asm)
8      predict_y = sig_clf.predict_proba(X_cv_asm)
9      cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=k_cfl.classes_, eps=
10
11  for i in zip(alpha,cv_log_error_array):
12      print ('log_loss for k = ',i[0],'is',i[1])
13
14  best_alpha = np.argmin(cv_log_error_array)
15
16  fig, ax = plt.subplots()
17  ax.plot(alpha, cv_log_error_array,c='g')
18  for i, txt in enumerate(np.round(cv_log_error_array,3)):
19      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
20  plt.grid()
21  plt.title("Cross Validation Error for each alpha")
22  plt.xlabel("Alpha i's")
23  plt.ylabel("Error measure")
24  plt.show()
25
26  k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
27  k_cfl.fit(X_train_asm,y_train_asm)
28  sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
29  sig_clf.fit(X_train_asm, y_train_asm)
30
31  predict_y = sig_clf.predict_proba(X_train_asm)
32  print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_
33  predict_y = sig_clf.predict_proba(X_cv_asm)
34  print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
35  predict_y = sig_clf.predict_proba(X_test_asm)
36  print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_lo
37  plot_confusion_matrix(y_test_asm, sig_clf.predict(X_test_asm),'KNN(n = {})'.format(alp

```

```

log_loss for k = 3 is 0.07822634324887906
log_loss for k = 5 is 0.08882553980508383
log_loss for k = 7 is 0.10034291193361768
log_loss for k = 9 is 0.1121531486151857
log_loss for k = 11 is 0.1223683462308549
log_loss for k = 13 is 0.1324030327413686
log_loss for k = 15 is 0.14605447214817008

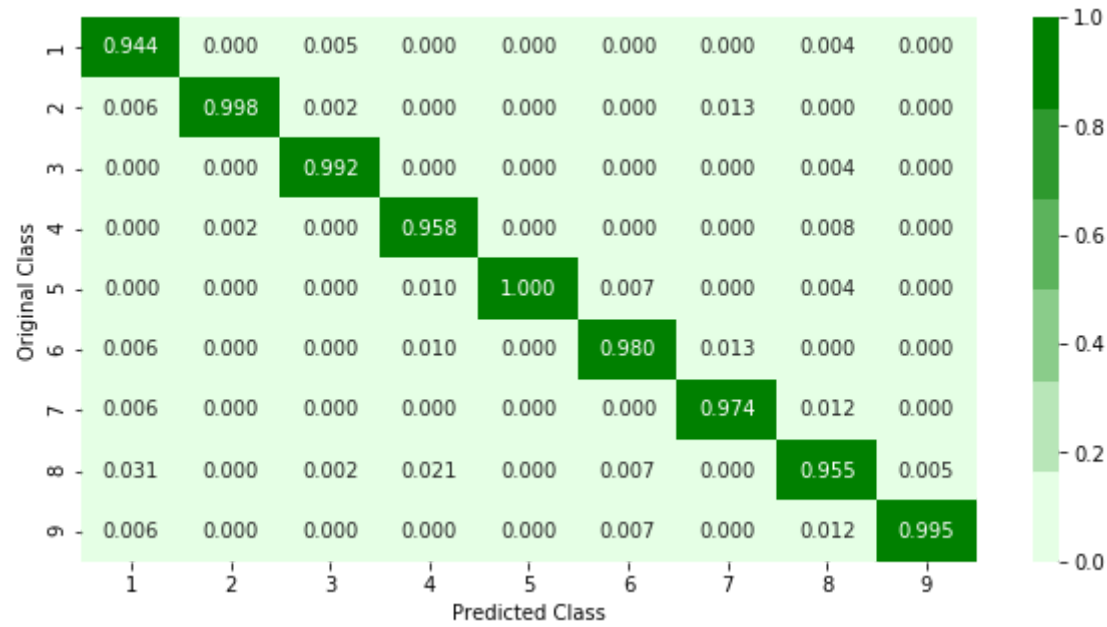
```



For values of best alpha = 3 The train log loss is: 0.051252139577566594  
For values of best alpha = 3 The cross validation log loss is: 0.07822634324887906  
For values of best alpha = 3 The test log loss is: 0.08829219350799569  
----- Confusion matrix -----  
-----



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



----- Recall matrix -----



Percentage of misclassified points using KNN(n = 3) Model: 2.07%

4.4.2 Logistic Regression

In [58]:

```

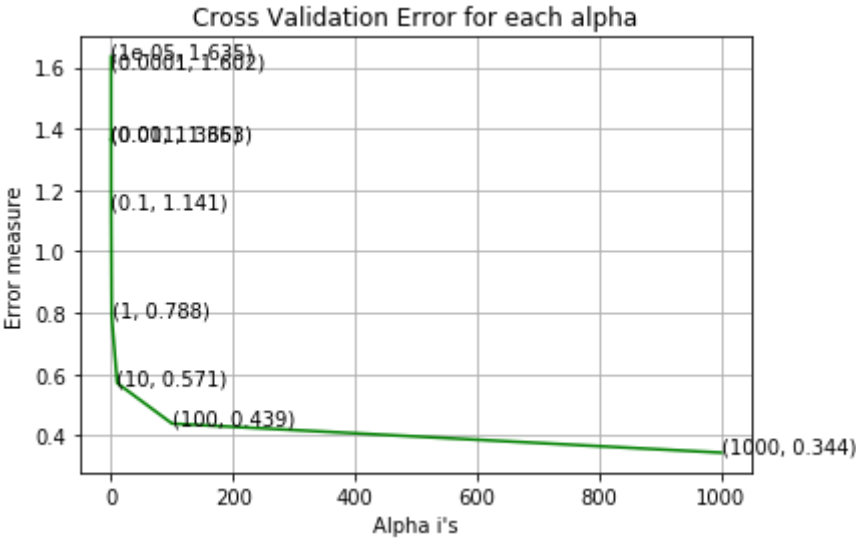
1  alpha = [10 ** x for x in range(-5, 4)]
2  cv_log_error_array=[]
3  for i in alpha:
4      logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
5      logisticR.fit(X_train_asm,y_train_asm)
6      sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
7      sig_clf.fit(X_train_asm,y_train_asm)
8      predict_y = sig_clf.predict_proba(X_cv_asm)
9      cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.classes_,
10
11  for i in range(len(cv_log_error_array)):
12      print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
13
14  best_alpha = np.argmin(cv_log_error_array)
15
16  fig, ax = plt.subplots()
17  ax.plot(alpha, cv_log_error_array,c='g')
18  for i, txt in enumerate(np.round(cv_log_error_array,3)):
19      ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
20  plt.grid()
21  plt.title("Cross Validation Error for each alpha")
22  plt.xlabel("Alpha i's")
23  plt.ylabel("Error measure")
24  plt.show()
25
26  logisticR=LogisticRegression(penalty='l2',C=alpha[best_alpha],class_weight='balanced')
27  logisticR.fit(X_train_asm,y_train_asm)
28  sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
29  sig_clf.fit(X_train_asm,y_train_asm)
30
31
32  predict_y = sig_clf.predict_proba(X_train_asm)
33  print ('log loss for train data',log_loss(y_train_asm, predict_y, labels=logisticR.classes_))
34  predict_y = sig_clf.predict_proba(X_cv_asm)
35  print ('log loss for cv data',log_loss(y_cv_asm, predict_y, labels=logisticR.classes_))
36  predict_y = sig_clf.predict_proba(X_test_asm)
37  print ('log loss for test data',log_loss(y_test_asm, predict_y, labels=logisticR.classes_))
38  plot_confusion_matrix(y_test_asm, sig_clf.predict(X_test_asm),'Logistic Regression (C = '

```

```

log_loss for c = 1e-05 is 1.6353018273097308
log_loss for c = 0.0001 is 1.601693313639587
log_loss for c = 0.001 is 1.3634943675721063
log_loss for c = 0.01 is 1.3649332524062547
log_loss for c = 0.1 is 1.1408752185228466
log_loss for c = 1 is 0.7875378732127429
log_loss for c = 10 is 0.5707092238819212
log_loss for c = 100 is 0.43859802179195917
log_loss for c = 1000 is 0.3438629524532219

```

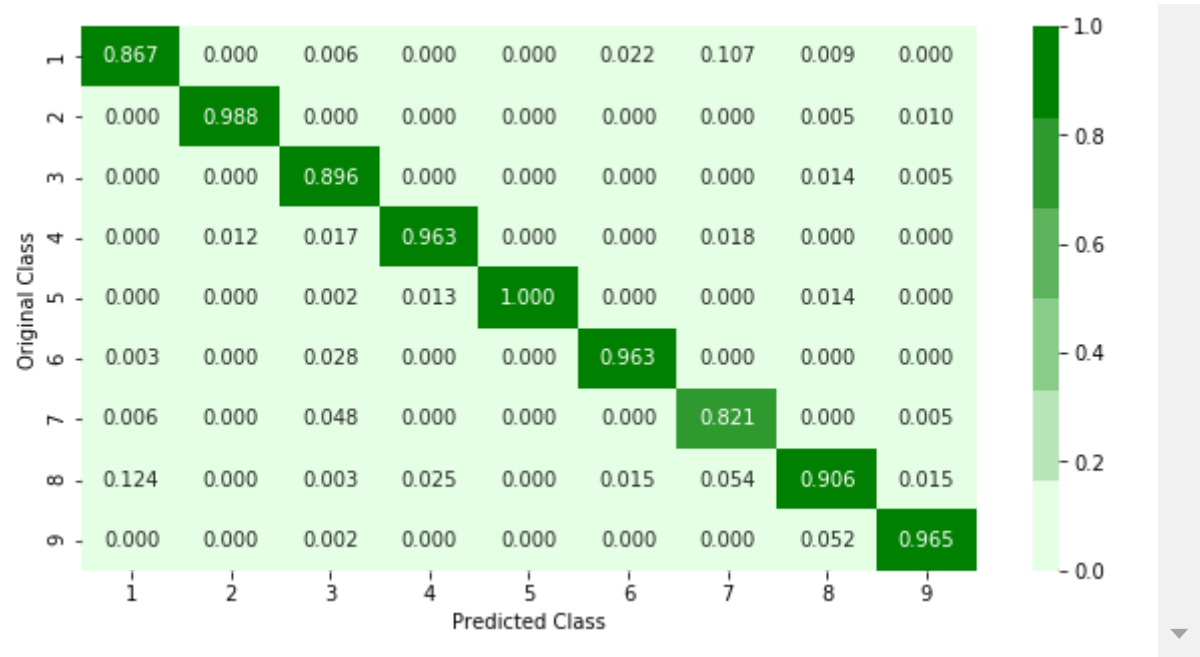


log loss for train data 0.3281465886689375  
log loss for cv data 0.3438629524532219  
log loss for test data 0.3478471628021207

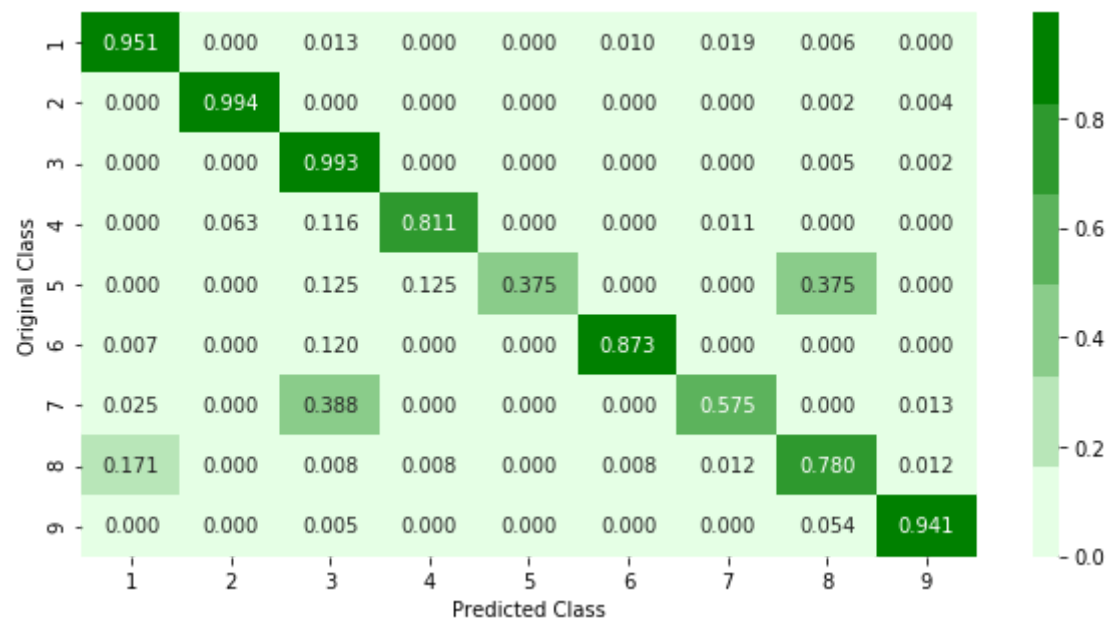
----- Confusion matrix -----  
-----



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



----- Recall matrix -----  
-----



Percentage of misclassified points using Logistic Regression (C = 1000) Mode 1: 7.54%

4.4.3 Random Forest Classifier



In [59]:

```

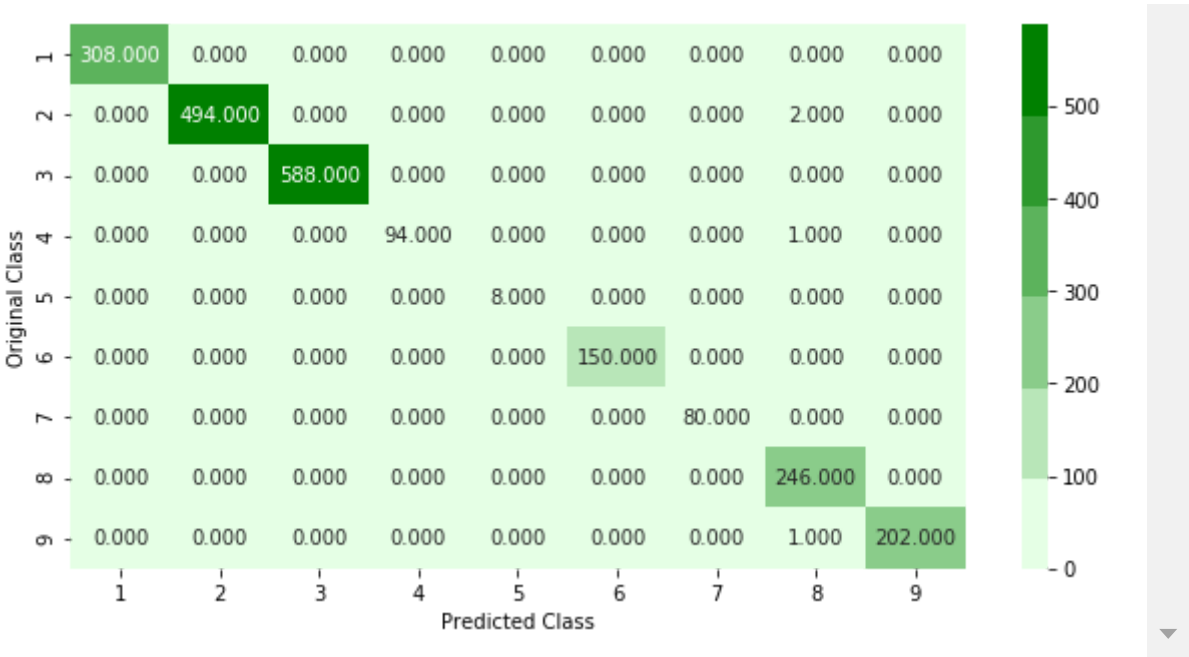
1  if not os.path.exists('RF_one_param_tuning_asm.joblib'):
2      alpha=[10,50,100,500,1000,2000,3000]
3      cv_log_error_array=[]
4      train_log_error_array=[]
5      for i in alpha:
6          r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
7          r_cfl.fit(X_train_asm,y_train_asm)
8          sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
9          sig_clf.fit(X_train_asm,y_train_asm)
10         predict_y = sig_clf.predict_proba(X_cv_asm)
11         cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_,
12
13     for i in range(len(cv_log_error_array)):
14         print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
15
16
17     best_alpha = np.argmin(cv_log_error_array)
18
19     fig, ax = plt.subplots()
20     ax.plot(alpha, cv_log_error_array,c='g')
21     for i, txt in enumerate(np.round(cv_log_error_array,3)):
22         ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
23     plt.grid()
24     plt.title("Cross Validation Error for each alpha")
25     plt.xlabel("Alpha i's")
26     plt.ylabel("Error measure")
27     plt.show()
28
29     r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
30     r_cfl.fit(X_train_asm,y_train_asm)
31     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
32     sig_clf.fit(X_train_asm,y_train_asm)
33
34     joblib.dump(sig_clf,'RF_one_param_tuning_asm.joblib')
35 else:
36     sig_clf = joblib.load('RF_one_param_tuning_asm.joblib')
37
38 predict_y = sig_clf.predict_proba(X_train_asm)
39 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train_asm, predict_y, labels=r_cfl.classes_))
40 predict_y = sig_clf.predict_proba(X_cv_asm)
41 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_))
42 predict_y = sig_clf.predict_proba(X_test_asm)
43 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test_asm, predict_y, labels=r_cfl.classes_))
44 plot_confusion_matrix(y_test_asm, sig_clf.predict(X_test_asm),'Random Forest (n_estimators=1000)')

```

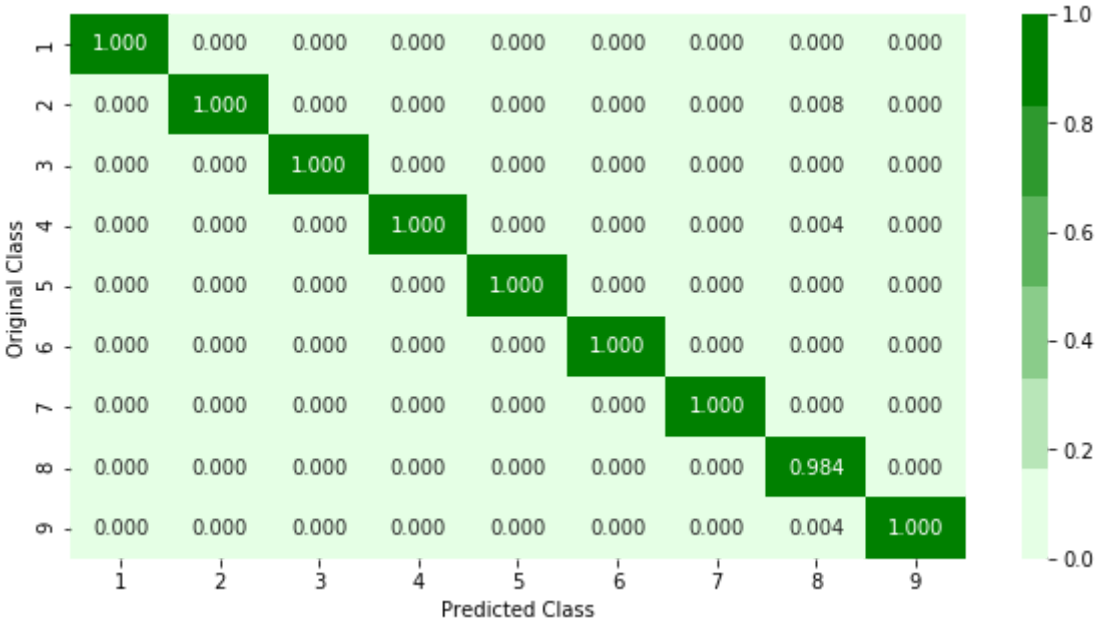
```

For values of best alpha = 1000 The train log loss is: 0.02130763882281661
For values of best alpha = 1000 The cross validation log loss is: 0.0217788
30232004408
For values of best alpha = 1000 The test log loss is: 0.023422488727851302
----- Confusion matrix -----
-----

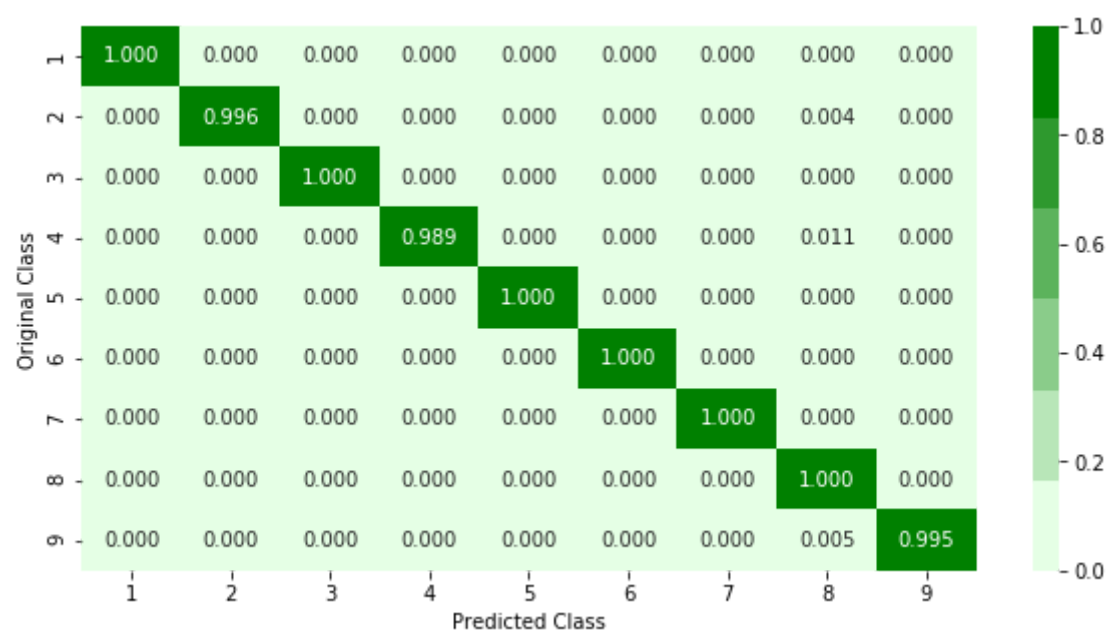
```



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



----- Recall matrix -----  
-----



Percentage of misclassified points using Random Forest (n\_estimators = 1000)  
Model: 0.18%

4.4.4 XgBoost Classifier

In [60]:

```

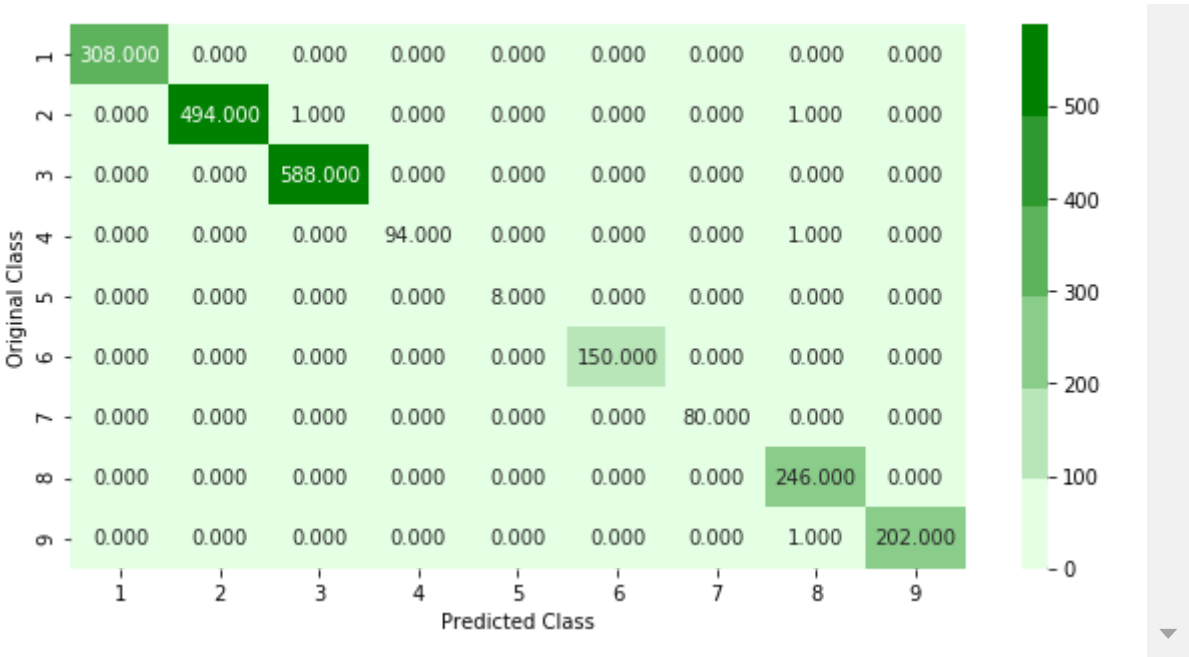
1  if not os.path.exists('XGB_one_param_tuning_asm.joblib'):
2      alpha=[10,50,100,500,1000,2000]
3      cv_log_error_array=[]
4      for i in alpha:
5          x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
6          x_cfl.fit(X_train_asm,y_train_asm)
7          sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
8          sig_clf.fit(X_train_asm,y_train_asm)
9          predict_y = sig_clf.predict_proba(X_cv_asm)
10         cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=x_cfl.classes_,
11
12     for i in range(len(cv_log_error_array)):
13         print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
14
15
16     best_alpha = np.argmin(cv_log_error_array)
17
18     fig, ax = plt.subplots()
19     ax.plot(alpha, cv_log_error_array,c='g')
20     for i, txt in enumerate(np.round(cv_log_error_array,3)):
21         ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
22     plt.grid()
23     plt.title("Cross Validation Error for each alpha")
24     plt.xlabel("Alpha i's")
25     plt.ylabel("Error measure")
26     plt.show()
27
28     x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
29     x_cfl.fit(X_train_asm,y_train_asm)
30     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
31     sig_clf.fit(X_train_asm,y_train_asm)
32
33     joblib.dump(sig_clf,'XGB_one_param_tuning_asm.joblib')
34 else:
35     sig_clf = joblib.load('XGB_one_param_tuning_asm.joblib')
36
37 predict_y = sig_clf.predict_proba(X_train_asm)
38 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_
39 predict_y = sig_clf.predict_proba(X_cv_asm)
40 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
41 predict_y = sig_clf.predict_proba(X_test_asm)
42 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_lo
43 plot_confusion_matrix(y_test_asm, sig_clf.predict(X_test_asm),'XGB Model (n_estimators

```

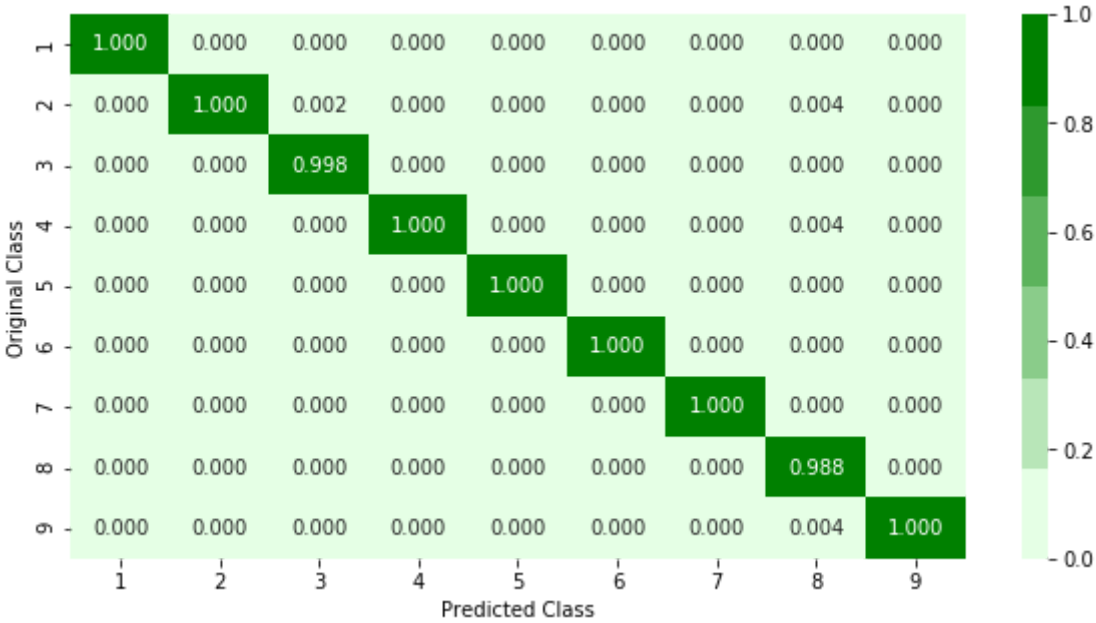
```

For values of best alpha = 1000 The train log loss is: 0.016417298234131933
For values of best alpha = 1000 The cross validation log loss is: 0.0165180
29039938188
For values of best alpha = 1000 The test log loss is: 0.01880318519615657
----- Confusion matrix -----
-----

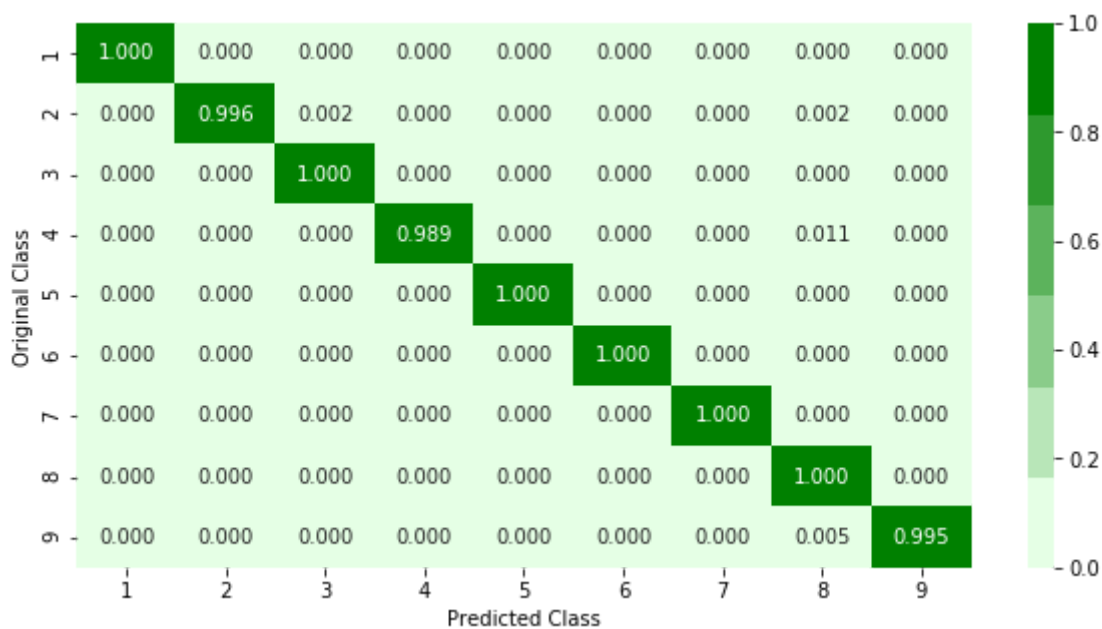
```



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



----- Recall matrix -----  
-----



Percentage of misclassified points using XGB Model (n\_estimators = 1000) Model: 0.18%

#### 4.4.5 Xgboost Classifier with best hyperparameters

In [61]:

```
1 if not os.path.exists('XGB_multiple_param_tuning_asm.joblib'):
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,n_iter=10,verbose=10)
12    random_cfl1.fit(X_train_asm,y_train_asm)
13    joblib.dump(random_cfl1, 'XGB_multiple_param_tuning_asm.joblib')
14 else:
15     random_cfl1 = joblib.load('XGB_multiple_param_tuning_asm.joblib')
```

In [62]:

```

1 if not os.path.exists('XGB_multiple_param_tuning_Calibrated_asm.joblib'):
2     c_cfl = CalibratedClassifierCV(random_cfl1.best_estimator_ ,method='sigmoid')
3     c_cfl.fit(X_train_asm,y_train_asm)
4     joblib.dump(c_cfl, 'XGB_multiple_param_tuning_Calibrated_asm.joblib')
5 else:
6     c_cfl = joblib.load('XGB_multiple_param_tuning_Calibrated_asm.joblib')
7
8 predict_y = c_cfl.predict_proba(X_train_asm)
9 print ("The train log loss is:",log_loss(y_train_asm, predict_y))
10 predict_y = c_cfl.predict_proba(X_cv_asm)
11 print("The cross validation log loss is:",log_loss(y_cv_asm, predict_y))
12 predict_y = c_cfl.predict_proba(X_test_asm)
13 print("The test log loss is:",log_loss(y_test_asm, predict_y))
14 plot_confusion_matrix(y_test_asm, c_cfl.predict(X_test_asm),'XGB Model (with more para

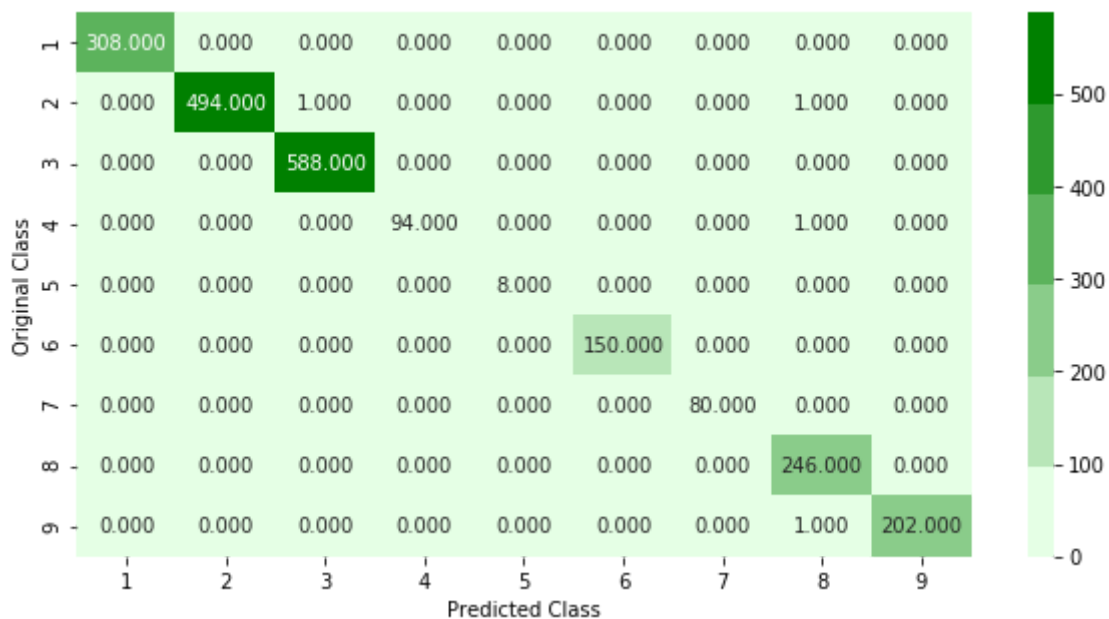
```

The train log loss is: 0.01717220458840472

The cross validation log loss is: 0.01829959740822653

The test log loss is: 0.019137559973987937

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



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



----- Recall matrix -----



Percentage of misclassified points using XGB Model (with more parameters tuned) Model: 0.18%

## 4.5. Machine Learning models on features of both .asm and .bytes files

### 4.5.1. Merging both asm and byte file features



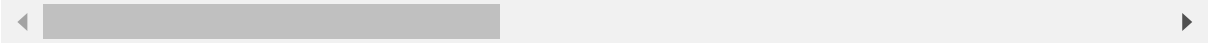
In [63]:

```
1 result.head()
```

Out[63]:

	ID	0	1	2	3	4	5	
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.0020
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.0047
2	01jsnpXSAlgW6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.0050
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.0003
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.0001

5 rows × 266 columns



In [64]:

```
1 result.rename(columns={'filelines':'bytefilelines'},inplace=True)
```

In [65]:

```
1 result_asm.head()
```

Out[65]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.ec
0	01kcPWA9K2BOxQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084	
1	02IOCvYEy8mjiuAQHax3	0.096045	0.001230	0.0	0.000246	0.030916	0.0	0.000000	
2	04EjldbPV5e1XroFOpiN	0.107345	0.088779	0.0	0.002091	0.001493	0.0	0.000000	
3	05EeG39MTRrl6VY21DPd	0.096045	0.016323	0.0	0.001935	0.000416	0.0	0.000882	
4	05LHG8fR3iPn6aglo9z7	0.096045	0.000000	0.0	0.001066	0.000000	0.0	0.000133	

5 rows × 62 columns



In [66]:

```
1 result_asm.rename(columns = {'filelines':'asmfilelines','[':'bracket','-':'minussign',
```

In [67]:

```
1 print(result.shape)
2 print(result_asm.shape)
```

```
(10868, 266)
(10868, 62)
```

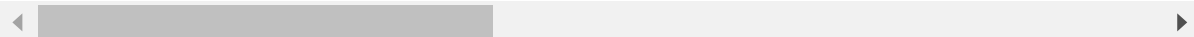
In [68]:

```
1 result_x = pd.merge(result,result_asm.drop(['Class'], axis=1),on='ID', how='left')
2 result_y = result_x['Class']
3 result_x = result_x.drop(['Class'],axis=1)
4 result_x.head()
```

Out[68]:

	ID	0	1	2	3	4	5	
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.0020
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.0047
2	01jsnpXSAIgw6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.0050
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.0003
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.0001

5 rows × 325 columns



Dchad more features

In [69]:

```

1 dchad2 = pd.DataFrame()
2 for file in os.listdir('dchad/dchad2'):
3     if file[-3:] == 'csv':
4         filelocation = 'dchad/dchad2/'+file
5         dchad2 = pd.concat(objs=[dchad2,pd.read_csv(filelocation)],axis=1)
6 dchad2.rename(columns={'filename':'ID'},inplace=True)
7 dchad2.head()

```

Out[69]:

	ID	edx	esi	es	ds	ss	cs	ah	al	ax	...	ASM_964	ASM_972
0	01IsoiSMh5gxyDYTI4CB	750	496	3	0	0	0	8	224	49	...	32	49
1	01SuzwMJEIXsK7A8dQbl	1121	24	3	1	4	2	6	22	7	...	48	9
2	01azqd4lnC7m9JpocGv5	1493	1900	0	0	0	0	1	398	0	...	48	9
3	01jsnpXSAlgw6aPeDxrU	525	4	0	0	0	0	0	0	0	...	48	9
4	01kcPWA9K2BOxQeS5Rju	23	35	0	0	0	0	0	3	0	...	48	89

5 rows × 1276 columns



In [70]:

```

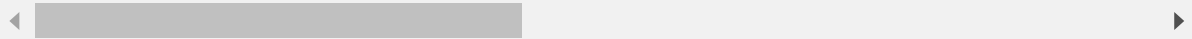
1 dchad2_fix = pd.DataFrame()
2 for column in dchad2.iteritems():
3     if column[0] in dchad2_fix.columns.values:
4         continue
5     else:
6         dchad2_fix[column[0]] = column[1]
7 dchad2_fix.head()

```

Out[70]:

	ID	edx	esi	es	ds	ss	cs	ah	al	ax	...	train_byte_p1	train_
0	01IsoiSMh5gxyDYTI4CB	750	496	3	0	0	0	8	224	49	...	1.0	(
1	01SuzwMJEIXsK7A8dQbl	1121	24	3	1	4	2	6	22	7	...	1.0	(
2	01azqd4lnC7m9JpocGv5	1493	1900	0	0	0	0	1	398	0	...	1.0	(
3	01jsnpXSAlgw6aPeDxrU	525	4	0	0	0	0	0	0	0	...	1.0	(
4	01kcPWA9K2BOxQeS5Rju	23	35	0	0	0	0	0	3	0	...	1.0	(

5 rows × 653 columns



In [71]:

```
1 dchad2_fix = normalize(dchad2_fix)
```

In [72]:

```

1 result_x = pd.merge(result_x,dchad2_fix,on='ID', how='left')
2 result_x = result_x.drop(labels=['ID'],axis=1)
3 result_x = result_x.loc[:,[
4     '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', '0a', '0b', '0c',
5     '0d', '0e', '0f', '10', '11', '12', '13', '14', '15', '16', '17',
6     '18', '19', '1a', '1b', '1c', '1d', '1e', '1f', '20', '21', '22',
7     '23', '24', '25', '26', '27', '28', '29', '2a', '2b', '2c', '2d',
8     '2e', '2f', '30', '31', '32', '33', '34', '35', '36', '37', '38',
9     '39', '3a', '3b', '3c', '3d', '3e', '3f', '40', '41', '42', '43',
10    '44', '45', '46', '47', '48', '49', '4a', '4b', '4c', '4d', '4e',
11    '4f', '50', '51', '52', '53', '54', '55', '56', '57', '58', '59',
12    '5a', '5b', '5c', '5d', '5e', '5f', '60', '61', '62', '63', '64',
13    '65', '66', '67', '68', '69', '6a', '6b', '6c', '6d', '6e', '6f',
14    '70', '71', '72', '73', '74', '75', '76', '77', '78', '79', '7a',
15    '7b', '7c', '7d', '7e', '7f', '80', '81', '82', '83', '84', '85',
16    '86', '87', '88', '89', '8a', '8b', '8c', '8d', '8e', '8f', '90',
17    '91', '92', '93', '94', '95', '96', '97', '98', '99', '9a', '9b',
18    '9c', '9d', '9e', '9f', 'a0', 'a1', 'a2', 'a3', 'a4', 'a5', 'a6',
19    'a7', 'a8', 'a9', 'aa', 'ab', 'ac', 'ad', 'ae', 'af', 'b0', 'b1',
20    'b2', 'b3', 'b4', 'b5', 'b6', 'b7', 'b8', 'b9', 'ba', 'bb', 'bc',
21    'bd', 'be', 'bf', 'c0', 'c1', 'c2', 'c3', 'c4', 'c5', 'c6', 'c7',
22    'c8', 'c9', 'ca', 'cb', 'cc', 'cd', 'ce', 'cf', 'd0', 'd1', 'd2',
23    'd3', 'd4', 'd5', 'd6', 'd7', 'd8', 'd9', 'da', 'dc',
24    'de', 'df', 'e0', 'e1', 'e2', 'e3', 'e4', 'e5', 'e6', 'e7',
25    'e8', 'e9', 'ea', 'eb', 'ec', 'ed', 'ee', 'ef', 'f0', 'f1', 'f2',
26    'f3', 'f4', 'f5', 'f6', 'f7', 'f8', 'f9', 'fa', 'fb', 'fc', 'fd',
27    'fe', 'ff', '??',
28    'bytefilelines', 'vertex_count',
29    'delta_max', 'density', 'HEADER:', '.idata:', '.rsrc:', '.reloc:',
30    'retf', '.dll', 'eip', 'asmfilelines', '.dataweighted',
31    'esi', 'es', 'ds', 'ss', 'cs', 'ah', 'al', 'ax', 'bh', 'bl', 'bx',
32    'ch', 'cl', 'cx', 'dh', 'dx', 'cdq', 'cld', 'cli',
33    'cmc', 'const', 'daa', 'faddp', 'fchs', 'fdiv', 'fdivp', 'fild', 'fistp', 'fword',
34    'jge', 'jl', 'jno', 'jo',
35    'mul', 'neg', 'not', 'out', 'popf', 'proc', 'rcl', 'rep', 'retn',
36    'sal', 'sar', 'sbb', 'setnle', 'setnz', 'setz',
37    'shld', 'stc', 'sti',
38    'wait', 'wcslen', '__vbaUII2', 'send',
39    'IsDBCSLeadByte', 'SetWindowLongW', 'time', 'GetWindowLongW',
40    'ValidateRect', 'socket', 'memmove', 'connect', 'lstrcpyW',
41    'wcscmp', 'lstrcpw', 'wcscpy', 'GetUserDefaultLangID', 'Escape', 'LoadLibrary',
42    'WriteFile', 'GetModuleFileNameA', 'CloseHandle', 'RegCloseKey',
43    'VirtualFree', 'GetLastError', 'FreeLibrary',
44    'MultiByteToWideChar', 'RegQueryValueExA', 'GetTickCount',
45    'GetStartupInfoA', 'RegOpenKeyExA', 'GetCurrentProcess',
46    'EnterCriticalSection', 'LeaveCriticalSection', 'DeleteCriticalSection',
47    'InitializeCriticalSection', 'LocalAlloc', 'RtlUnwind', 'lstrlenA',
48    'GetDC', 'UnhandledExceptionFilter', 'GetLocaleInfoA',
49    'TlsSetValue', 'TlsGetValue', 'RaiseException', 'TerminateProcess',
50    'CharNextA', 'HeapAlloc', 'LocalFree', 'GetCPInfo', 'GetFileType',
51    'HeapFree', 'DispatchMessageA', 'ShowWindow', 'DestroyWindow',
52    'GetSystemMetrics', 'InterlockedIncrement', 'InterlockedDecrement',
53    'LoadStringA', 'GetThreadLocale', 'SendMessageA', 'DeleteObject',
54    'GetClientRect', 'GetWindowRect', 'wsprintfA',
55    'SetLastError', 'QueryPerformanceCounter', 'GlobalUnlock',
56    'TranslateMessage', 'HeapReAlloc', 'GetDeviceCaps', 'SetFocus',
57    'CoTaskMemFree', '_initterm', 'CreatePopupMenu',

```

```

58 'GetForegroundWindow', 'HeapSize', 'exit', 'free',
59 'GetModuleHandleW', '_except_handler3', 'malloc', 'GetMenu',
60 'VariantClear', 'GetSystemTime', 'GetMessagePos',
61 'SetRect', 'memset', 'lstrlenW', 'memcpy', 'IsWindow',
62 'EnableWindow', 'SetWindowPos', 'lstrcpynA', 'GetOEMCP',
63 'SysFreeString', 'GetStringTypeW', 'lstrcpyA', 'GetWindowLongA',
64 'LCMapStringW', 'LCMapStringA', 'SetHandleCount', 'GetStringTypeA',
65 'GetEnvironmentStringsW', 'GetSysColor', 'RegisterClassA',
66 'InvalidateRect', 'InterlockedExchange',
67 'TlsAlloc', 'lstrcatA', 'CompareStringA',
68 'ReleaseDC', 'CoCreateInstance', 'GetWindowTextA', 'lstrcmpiA',
69 'GetDlgItem', 'TlsFree', 'GetParent', 'GetWindow', 'SetMenu',
70 '__p_fmode', 'CompareStringW', '__getmainargs', 'sprintf',
71 '_XcptFilter', 'DrawEdge',
72 'GetIconInfo', '_controlfp', 'strstr', 'IsValidCodePage',
73 'VariantInit', 'strncpy', '__CxxFrameHandler', 'GetCursor',
74 'ReleaseMutex', 'strlen', '__vbaGenerateBoundsError', 'strchr',
75 'strcpy', 'Rectangle', 'fclose', 'rand', 'strcat',
76 'GetMessageTime', 'SendMessageW', 'SendDlgItemMessageA',
77 'LoadStringW', 'Virtual', 'Offset', 'Import', 'var', 'UINT',
78 'LONG', 'BOOL', 'WORD', 'BYTES', 'large', 'short', 'dd.1',
79 'dw.1', 'ptr', 'DATA', 'FUNCTION', 'byte', 'word', 'char',
80 'stdcall', 'arg', 'asc', 'align', 'WinMain', 'unk', 'cookie',
81 'off', 'nullsub', 'DllEntryPoint', 'dll', 'HMENU', 'void',
82 'HDC', 'HWND', 'entry', 'Software', '__imp_', 'case', 'assume', 'entropy'
83 ]]

```

In [73]:

```

1 for column in result_x.iteritems():
2     if len(set(column[1].isnull())) == 1 and list(set(column[1].isnull()))[0] == True:
3         result_x.drop(labels=[column[0]],axis=1,inplace=True)

```

In [74]:

```
1 result_x.shape
```

Out[74]:

(10868, 510)

In [75]:

```

1 # from sklearn.feature_selection import SelectKBest, chi2
2 # bestfeatures = SelectKBest(score_func=chi2,k=100)
3 # result_x = bestfeatures.fit_transform(result_x,result_y)
4 # result_x = pd.DataFrame(result_x)

```

In [76]:

```
1 pp.ProfileReport(result_x)
```

Out[76]:

# Overview

## Dataset info

Number of variables	510
Number of observations	10868
Total Missing (%)	0.0%
Total size in memory	42.4 MiB
Average record size in memory	4.0 KiB

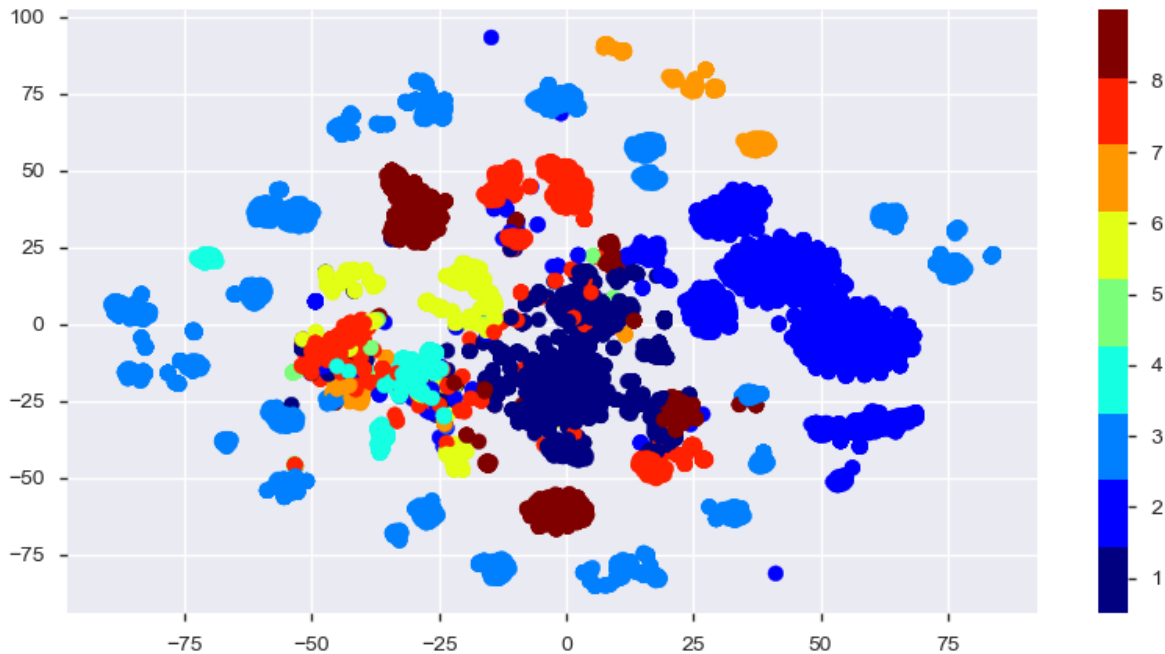
## Variables types

Numeric	509
Categorical	0

4.5.2. Multivariate Analysis on final fearures

In [77]:

```
1 plt.figure(figsize = (10,5))
2 xtsne=TSNE(perplexity=50)
3 results=xtsne.fit_transform(result_x)
4 vis_x = results[:, 0]
5 vis_y = results[:, 1]
6 plt.scatter(vis_x, vis_y, c=result_y, cmap=plt.cm.get_cmap("jet", 9))
7 plt.colorbar(ticks=range(9))
8 plt.clim(0.5, 9)
9 plt.show()
```



### 4.5.3. Train and Test split

In [78]:

```
1 X_train, X_test_merge, y_train, y_test_merge = train_test_split(result_x, result_y, str
2 X_train_merge, X_cv_merge, y_train_merge, y_cv_merge = train_test_split(X_train, y_tra
```

### 4.5.4. Random Forest Classifier on final features

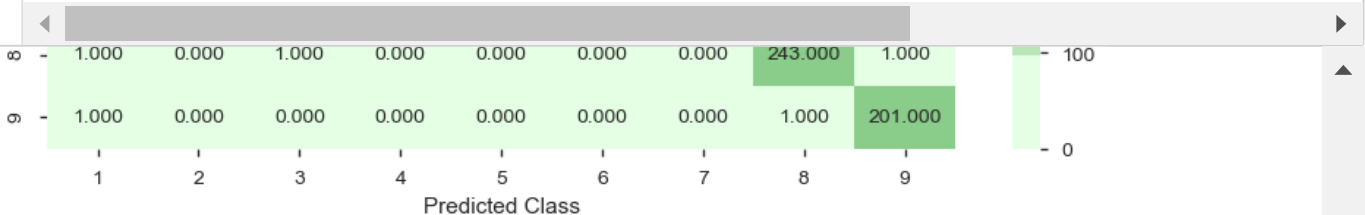


In [79]:

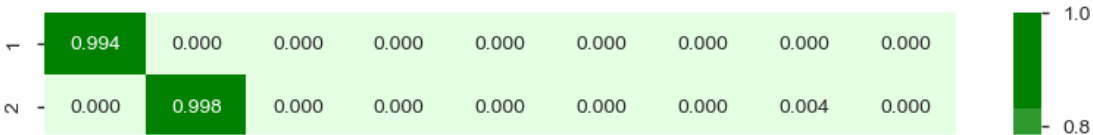
```

1  if not os.path.exists('RF_one_param_tuning_merge.joblib'):
2      alpha=[500,1000,5000,7000]
3      cv_log_error_array=[]
4      train_log_error_array=[]
5      for i in alpha:
6          r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
7          r_cfl.fit(X_train_merge,y_train_merge)
8          sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
9          sig_clf.fit(X_train_merge,y_train_merge)
10         predict_y = sig_clf.predict_proba(X_cv_merge)
11         cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=r_cfl.classes_))
12
13     for i in range(len(cv_log_error_array)):
14         print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
15
16
17     best_alpha = np.argmin(cv_log_error_array)
18
19     fig, ax = plt.subplots()
20     ax.plot(alpha, cv_log_error_array,c='g')
21     for i, txt in enumerate(np.round(cv_log_error_array,3)):
22         ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
23     plt.grid()
24     plt.title("Cross Validation Error for each alpha")
25     plt.xlabel("Alpha i's")
26     plt.ylabel("Error measure")
27     plt.show()
28
29     r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
30     r_cfl.fit(X_train_merge,y_train_merge)
31     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
32     sig_clf.fit(X_train_merge,y_train_merge)
33
34     joblib.dump(sig_clf,'RF_one_param_tuning_merge.joblib')
35 else:
36     sig_clf = joblib.load('RF_one_param_tuning_merge.joblib')
37
38 predict_y = sig_clf.predict_proba(X_train_merge)
39 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train_merge, predict_y))
40 predict_y = sig_clf.predict_proba(X_cv_merge)
41 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv_merge, predict_y))
42 predict_y = sig_clf.predict_proba(X_test_merge)
43 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test_merge, predict_y))
44 plot_confusion_matrix(y_test_merge, sig_clf.predict(X_test_merge),'Random Forest (n_estimators=7000)')

```



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



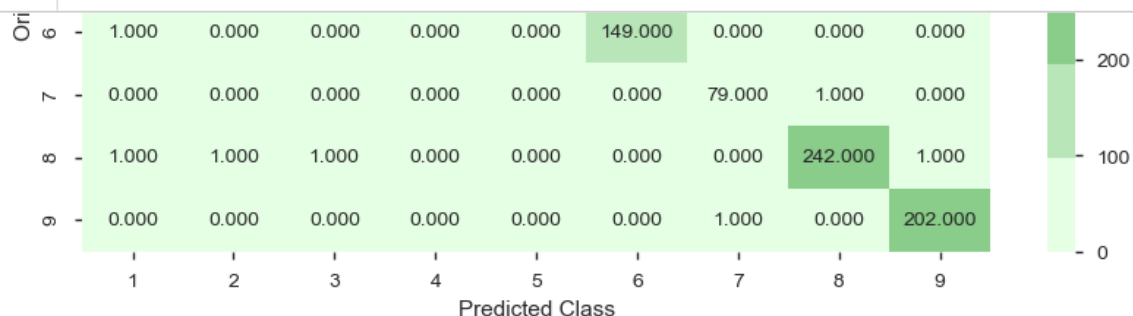
4.5.5. XgBoost Classifier on final features

In [80]:

```

1  if not os.path.exists('XGB_one_param_tuning_merge.joblib'):
2      alpha=[500,1000,5000]
3      cv_log_error_array=[]
4      for i in alpha:
5          x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
6          x_cfl.fit(X_train_merge,y_train_merge)
7          sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
8          sig_clf.fit(X_train_merge,y_train_merge)
9          predict_y = sig_clf.predict_proba(X_cv_merge)
10         cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=x_cfl.classes_)
11
12     for i in range(len(cv_log_error_array)):
13         print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
14
15
16     best_alpha = np.argmin(cv_log_error_array)
17
18     fig, ax = plt.subplots()
19     ax.plot(alpha, cv_log_error_array,c='g')
20     for i, txt in enumerate(np.round(cv_log_error_array,3)):
21         ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
22     plt.grid()
23     plt.title("Cross Validation Error for each alpha")
24     plt.xlabel("Alpha i's")
25     plt.ylabel("Error measure")
26     plt.show()
27
28     x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
29     x_cfl.fit(X_train_merge,y_train_merge)
30     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
31     sig_clf.fit(X_train_merge,y_train_merge)
32
33     joblib.dump(sig_clf,'XGB_one_param_tuning_merge.joblib')
34 else:
35     sig_clf = joblib.load('XGB_one_param_tuning_merge.joblib')
36
37 predict_y = sig_clf.predict_proba(X_train_merge)
38 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_
39 predict_y = sig_clf.predict_proba(X_cv_merge)
40 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
41 predict_y = sig_clf.predict_proba(X_test_merge)
42 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_lo
43 plot_confusion_matrix(y_test_merge, sig_clf.predict(X_test_merge),'XGB Model (n_estima

```



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



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

In [81]:

```
1 if not os.path.exists('XGB_multiple_param_tuning_merge.joblib'):
2     x_cfl=XGBClassifier()
3
4     prams={
5         'learning_rate':[0.001,0.01,0.1],
6         'n_estimators':[500,1000,2000,5000,7000],
7         'max_depth':[3,5,7,9]
8     }
9     random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,n_iter=10,verbose=10)
10    random_cfl1.fit(X_train_merge,y_train_merge)
11    joblib.dump(random_cfl1, 'XGB_multiple_param_tuning_merge.joblib')
12 else:
13     random_cfl1 = joblib.load('XGB_multiple_param_tuning_merge.joblib')
```

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

```
[Parallel(n_jobs=-1)]: Done    2 tasks      | elapsed: 31.7min
[Parallel(n_jobs=-1)]: Done    9 tasks      | elapsed: 106.1min
[Parallel(n_jobs=-1)]: Done   19 out of   30 | elapsed: 176.5min remaining: 10
2.2min
[Parallel(n_jobs=-1)]: Done   23 out of   30 | elapsed: 204.7min remaining: 6
2.3min
[Parallel(n_jobs=-1)]: Done   27 out of   30 | elapsed: 209.1min remaining: 2
3.2min
[Parallel(n_jobs=-1)]: Done   30 out of   30 | elapsed: 286.7min finished
```

In [82]:

```
1 random_cfl1.best_params_
```

Out[82]:

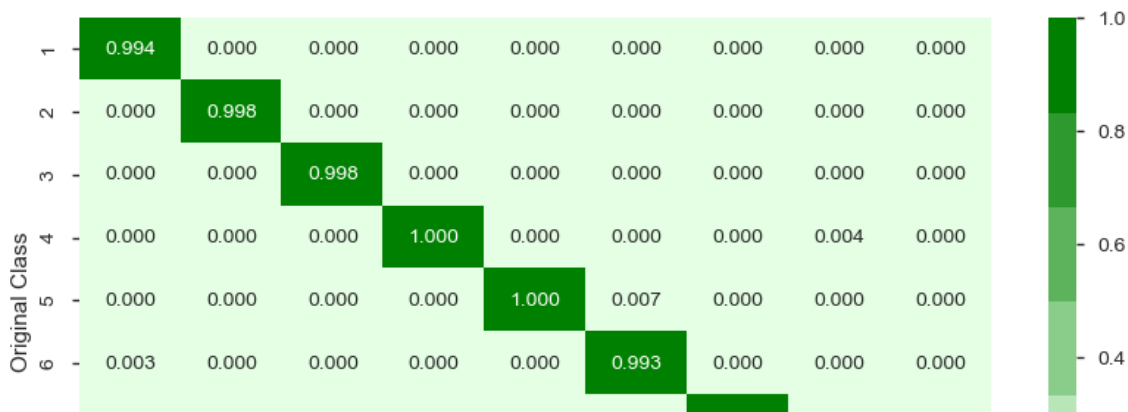
```
{'n_estimators': 7000, 'max_depth': 3, 'learning_rate': 0.1}
```

In [83]:

```

1 if not os.path.exists('XGB_multiple_param_tuning_Calibrated_merge.joblib'):
2     c_cfl = CalibratedClassifierCV(random_cfl1.best_estimator_ ,method='sigmoid')
3     c_cfl.fit(X_train_merge,y_train_merge)
4     joblib.dump(c_cfl, 'XGB_multiple_param_tuning_Calibrated_merge.joblib')
5 else:
6     c_cfl = joblib.load('XGB_multiple_param_tuning_Calibrated_merge.joblib')
7
8 predict_y = c_cfl.predict_proba(X_train_merge)
9 print ("The train log loss is:",log_loss(y_train_merge, predict_y))
10 predict_y = c_cfl.predict_proba(X_cv_merge)
11 print("The cross validation log loss is:",log_loss(y_cv_merge, predict_y))
12 predict_y = c_cfl.predict_proba(X_test_merge)
13 print("The test log loss is:",log_loss(y_test_merge, predict_y))
14 plot_confusion_matrix(y_test_merge, c_cfl.predict(X_test_merge),'XGB Model (with more

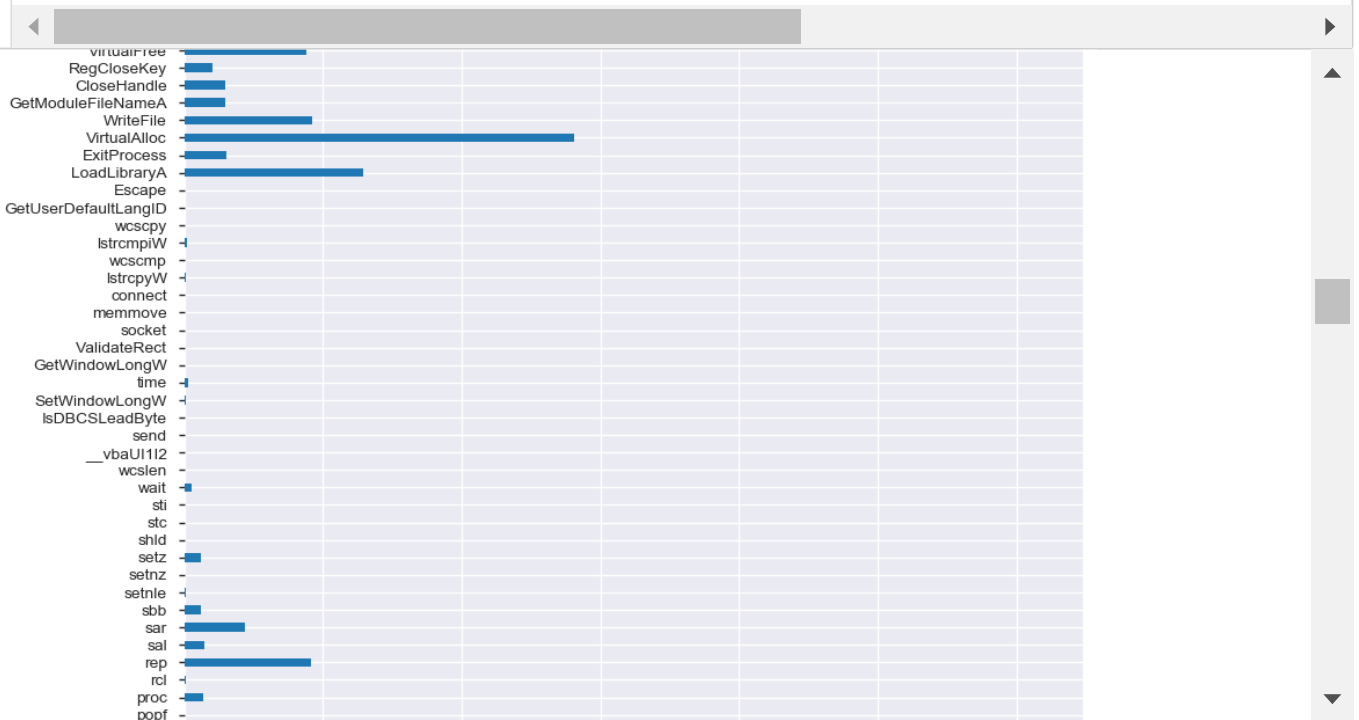
```



The feautre importances achieved by the XGBoost model plotted

In [84]:

```
1 feature_importance = pd.DataFrame(dict(zip(list(X_train_merge.columns.values), list(ran
2 feature_importance = feature_importance.transpose()
3 feature_importance.rename(columns={1: 'score'}, inplace=True)
4 feature_importance.plot.barh(figsize=(10,100));
```



## 5. MLP Softmax

1. Add bi-grams and n-gram features on byte files and improve the log-loss
2. Using the 'dchad' github account (<https://github.com/dchad/malware-detection>), decrease the logloss to  $\leq 0.01$
3. Watch the video ( <https://www.youtube.com/watch?v=VLQTRILGz5Y> ) that was in reference section and implement the image features to improve the logloss

For this I read through the documentation provided and tried out a few different things.

- (1) I fixed up all the code and reran everything from file processing to models (run time took a couple of days)
- (2) I tried the n-gram features, but struggle with memory (so I omitted this approach)
- (3) I used additional features from the document here: <https://arxiv.org/pdf/1511.04317.pdf> (<https://arxiv.org/pdf/1511.04317.pdf>) and added them in (['-', '+', '\*', '?', '@', ',']) during the processing of the asm files
- (4) I used the feature importances from the XGBoost model and weighted the two most important features more (please see bar chart above)
- (5) I performed hyperparameter tuning
- (6) I built a MLP softmax running GPU's with adam optimizer to see if it could possibly reach a more global optimum
- (7) I summarized the three top models in a pretty table

Conclusion:

I could not get the log loss below 0.01.

Since the heuristic is random I could only get the log loss between 0.01 - 0.04

In [85]:

```
1 # just something I tried for ngrams but I had memory issues:
2 if not os.path.exists('ngrambytesfile.csv'):
3     files = os.listdir('byteFiles')
4     bytefiledict = dict()
5     counter = 0
6     for file in files:
7         filename = file.split('.')[0]
8         filetext = ''
9         with open(os.path.join('byteFiles',file),'r') as content:
10             content = content.readlines()
11             for line in content:
12                 filetext = filetext + ' '.join(line)
13             if bytefiledict.get(filename) == None:
14                 bytefiledict[filename] = filetext.replace('\n','').strip()
15                 f = pd.DataFrame(bytefiledict,index=bytefiledict.keys()).reset_index()
16                 f.columns = ['ID','text']
17                 f.to_csv('ngrambytesfile.csv',index=False,header=False,mode='a')
18                 bytefiledict = dict()
19                 counter += 1
20                 if counter%2000 == 0:
21                     print(str(counter), ' documents processed')
22 print(str(counter), ' documents processed')
23 print('DONE!')
```

One hot encoding on the response variable

In [86]:

```
1 from sklearn.preprocessing import LabelBinarizer #MultilabelBinarizer
2 onehotencoding = LabelBinarizer()
3 y_train_onehot = onehotencoding.fit_transform(y_train_merge)
4 y_cv_onehot = onehotencoding.fit_transform(y_cv_merge)
5 y_test_onehot = onehotencoding.fit_transform(y_test_merge)
```

In [87]:

```
1 X_train_merge.head()
```

Out[87]:

	0	1	2	3	4	5	6	7	
8505	0.006351	0.012047	0.003451	0.003336	0.003887	0.003366	0.003716	0.005765	0.00531
4179	0.038119	0.001487	0.000328	0.000391	0.000386	0.000481	0.000273	0.000403	0.00048
10631	0.030318	0.026681	0.009145	0.008028	0.009156	0.008561	0.008355	0.013154	0.01458
5787	0.270462	0.005530	0.001695	0.001919	0.001777	0.001734	0.001890	0.002615	0.00308
9627	0.270160	0.005959	0.001927	0.001617	0.002290	0.002226	0.001608	0.002538	0.00315

5 rows × 510 columns



Multilayer Perceptron with softmax layer  
L1 - 317 relu activation units  
L2 - 300 relu activation units  
L3 - 9 softmax activation units

Performed batch normalization and dropout in between the layers  
batch size = 2000 and 150 epochs



In [89]:

```

1 import keras
2
3 # Create model:
4 model = keras.models.Sequential()
5
6 # Add Layers:
7 model.add(keras.layers.Dense(510,activation='relu',kernel_initializer=keras.initializer
8 model.add(keras.layers.BatchNormalization()) #L2
9 model.add(keras.layers.Dropout(0.5))
10 model.add(keras.layers.Dense(300,activation='relu',kernel_initializer=keras.initializer
11 model.add(keras.layers.BatchNormalization()) #L4
12 model.add(keras.layers.Dropout(0.5))
13 model.add(keras.layers.Dense(9,activation='softmax')) #L5
14
15 model.summary()
16
17 model.compile(optimizer='adam',loss='categorical_crossentropy',metrics=['accuracy'])
18 history = model.fit(X_train_merge,y_train_onehot,batch_size=1000,epochs=150,verbose=1,
19 model_output = model.evaluate(X_test_merge,y_test_onehot)

```

```

6955/6955 [=====] - ETA: 0s - loss: 0.0016 - acc:
1.000 - ETA: 0s - loss: 0.0026 - acc: 0.999 - ETA: 0s - loss: 0.0034 - ac
c: 0.998 - ETA: 0s - loss: 0.0060 - acc: 0.998 - ETA: 0s - loss: 0.0050 -
acc: 0.998 - ETA: 0s - loss: 0.0046 - acc: 0.998 - 1s 162us/step - loss:
0.0048 - acc: 0.9986 - val_loss: 0.0585 - val_acc: 0.9908
Epoch 149/150
6955/6955 [=====] - ETA: 0s - loss: 0.0050 - acc:
0.996 - ETA: 0s - loss: 0.0057 - acc: 0.996 - ETA: 0s - loss: 0.0046 - ac
c: 0.997 - ETA: 0s - loss: 0.0048 - acc: 0.997 - ETA: 0s - loss: 0.0042 -
acc: 0.997 - ETA: 0s - loss: 0.0041 - acc: 0.997 - 1s 146us/step - loss:
0.0037 - acc: 0.9980 - val_loss: 0.0584 - val_acc: 0.9902
Epoch 150/150
6955/6955 [=====] - ETA: 0s - loss: 0.0028 - acc:
0.999 - ETA: 0s - loss: 0.0032 - acc: 0.999 - ETA: 0s - loss: 0.0038 - ac
c: 0.998 - ETA: 0s - loss: 0.0036 - acc: 0.998 - ETA: 0s - loss: 0.0036 -
acc: 0.998 - ETA: 0s - loss: 0.0032 - acc: 0.999 - 1s 151us/step - loss:
0.0029 - acc: 0.9991 - val_loss: 0.0588 - val_acc: 0.9908
2174/2174 [=====] - ETA:  - ETA:  - ETA:  - ETA:
- ETA:  - ETA:  - ETA:  - ETA:  - ETA:  - 1s 239us/step

```

Model training results

In [90]:

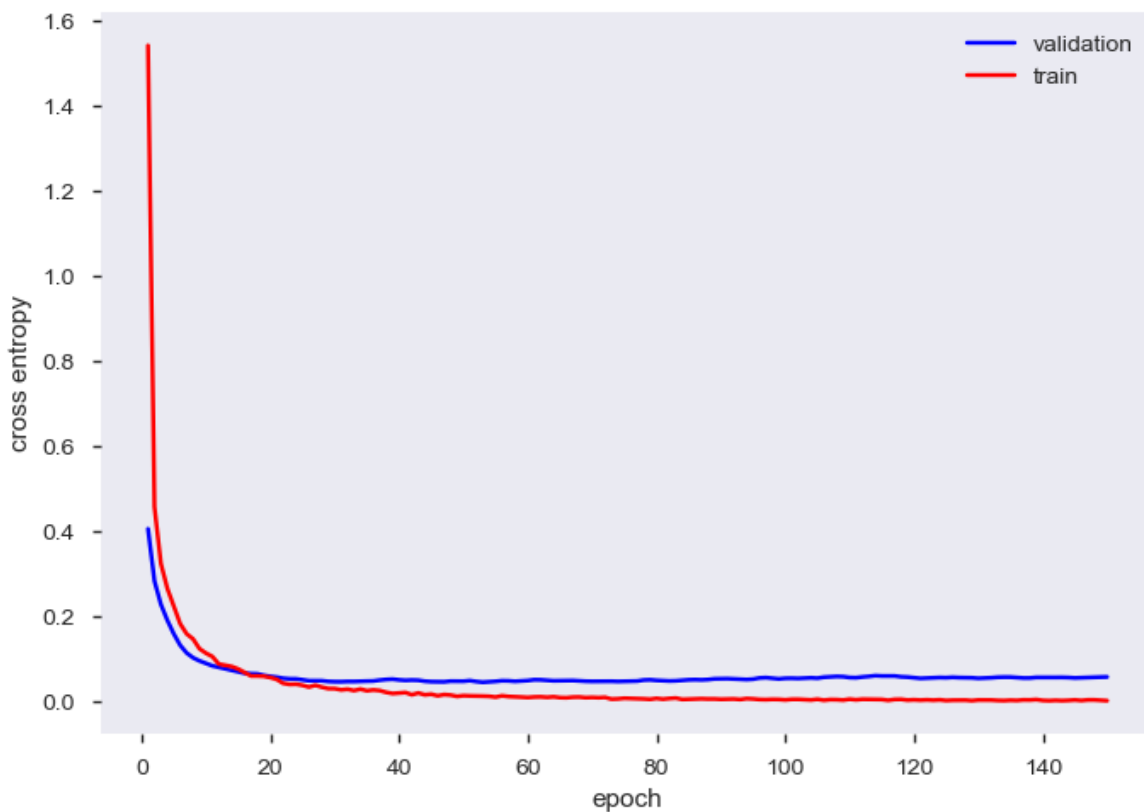
```

1 print("Cross entropy: ",model_output[0])
2 print("Accuracy: ",model_output[1])
3 def plt_dynamic(x,vy,ty,ax,colors=['b']):
4     ax.plot(x,vy,'b',label='validation')
5     ax.plot(x,ty,'r',label='train')
6     plt.legend()
7     plt.grid()
8     fig.canvas.draw()
9
10 fig,ax = plt.subplots(1,1)
11 ax.set_xlabel('epoch')
12 ax.set_ylabel('cross entropy')
13 x=list(range(1,150+1))
14 vy=history.history['val_loss']
15 ty=history.history['loss']
16 plt_dynamic(x,vy,ty,ax)

```

Cross entropy: 0.04784118044415784

Accuracy: 0.9917203311867525



In [91]:

```

1 table = PrettyTable(['Model','Test Log Loss','Validation Log Loss','Train Log Loss', ''])

```

In [92]:

```
1 table.add_row(['Random Forest', '0.02', '0.03', '0.01', '99.1%'])
2 table.add_row(['XGBoost', '0.01', '0.01', '0.01', '99.5%'])
3 table.add_row(['MLP (Softmax)', '0.04', '0.01', '0.01', '98.8%'])
```

In [93]:

```
1 print(table)
```

Model	Test Log Loss	Validation Log Loss	Train Log Loss	Accuracy
Random Forest	0.02	0.03	0.01	99.1%
XGBoost	0.01	0.01	0.01	99.5%
MLP (Softmax)	0.04	0.01	0.01	98.8%