

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 its 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 consists 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 given data
- Types of Malware:
 1. Ramnit

1. [Xanadu](#)
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:nothing, ds:_data,
s:nothing, gs:nothing
.text:00401000 56          push    esi
                            lea     eax, [esp+8]
.text:00401001 8D 44 24 08  push    eax
                            mov     esi, ecx
.text:00401005 50          call    ???
                            mov     dword ptr [esi], offset c
f_42BB08
.text:00401006 8B F1          mov     eax, esi
.text:00401008 E8 1C 1B 00 00  pop    esi
                            retn   4
.text:0040100D C7 06 08 BB 42 00
                            mov     dword ptr [esi], offset c
f_42BB08
.text:00401013 8B C6          mov     eax, esi
.text:00401015 5E          pop    esi
                            retn   4
.text:00401016 C2 04 00
.text:00401016
; -----
----- align 10h
.text:00401019 CC CC CC CC CC CC CC CC
.text:00401020 C7 01 08 BB 42 00          mov     dword ptr [ecx], offset c
f_42BB08
.text:00401026 E9 26 1C 00 00          jmp    sub_402C51
.text:00401026
; -----
----- align 10h
.text:0040102B CC CC CC CC CC CC
.text:00401030 56          push    esi
                            mov     esi, ecx
.text:00401031 8B F1          mov     dword ptr [esi], offset c
f_42BB08
.text:00401033 C7 06 08 BB 42 00
                            call    sub_402C51
                            test   byte ptr [esp+8], 1
                            jz    short loc_40104E
.text:00401043 74 09          push    esi
                            call    ??3@YAXPAX@Z      ; operator delete(void *)
.text:00401045 56
.text:00401046 E8 6C 1E 00 00          add    esp, 4
.delete(void *)
.text:0040104B 83 C4 04
.text:0040104E
.loc_40104E:           ; CODE XREF:
.text:0040104E j
.text:0040104E 8B C6          mov     eax, esi
.text:00401050 5E          pop    esi
                            retn   4
.text:00401051 C2 04 00
.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 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 11 20 80 04 80 10 00 20 00 00 25 00
00401080 00 00 01 00 00 04 00 10 02 C1 80 80 00 20 20 00
00401090 08 A0 01 01 44 28 00 00 08 10 20 00 02 08 00 00
004010A0 00 40 00 00 00 34 40 40 00 04 00 08 80 08 00 08
004010B0 10 00 40 00 68 02 40 04 E1 00 28 14 00 08 20 0A
004010C0 06 01 02 00 40 00 00 00 00 00 00 20 00 02 00 04
004010D0 80 18 90 00 00 10 A0 00 45 09 00 10 04 40 44 82
004010E0 90 00 26 10 00 00 04 00 82 00 00 00 20 40 00 00
004010F0 B4 00 00 40 00 02 20 25 08 00 00 00 00 00 00 00 00
00401100 08 00 00 50 00 08 40 50 00 02 06 22 08 85 30 00
00401110 00 80 00 80 60 00 09 00 04 20 00 00 00 00 00 00 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>

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learing Objectives and Constraints

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

Constraints:

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

2.3. Train and Test Dataset

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

2.4. Useful blogs, videos and reference papers

<http://blog.kaggle.com/2015/05/26/microsoft-malware-winners-interview-1st-place-no-to-overfitting/>
<https://arxiv.org/pdf/1511.04317.pdf>

First place solution in Kaggle competition: <https://www.youtube.com/watch?v=VLQTRILGz5Y>

<https://github.com/dchad/malware-detection>
<http://vizsec.org/files/2011/Nataraj.pdf>
https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu_pIB6ua?dl=0
" Cross validation is more trustworthy than domain knowledge."

3. Exploratory Data Analysis

In [2]:

```
import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use('nbAgg')
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process# this is used for multithreading
import multiprocessing
import codecs# this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from tqdm import tqdm
import scipy
```

In [2]:

```
#lists all the files in train folder
#folder contains .asm and .bytes files
os.listdir('train')
```

Out[2]:

```
['01azqd4InC7m9JpocGv5.asm',
 '01azqd4InC7m9JpocGv5.bytes',
 '01IsoiSMh5gxyDYTl4CB.asm',
 '01IsoiSMh5gxyDYTl4CB.bytes',
 '01jsnpXSAlg6aPeDxrU.asm',
 '01jsnpXSAlg6aPeDxrU.bytes',
 '01kcPWA9K2B0xQeS5Rju.asm',
 '01kcPWA9K2B0xQeS5Rju.bytes',
 '01SuzwMJEIXsK7A8dQbl.asm',
 '01SuzwMJEIXsK7A8dQbl.bytes',
 '02IOcvYEy8mjuiAQHax3.asm',
 '02IOcvYEy8mjuiAQHax3.bytes',
 '02JqQ7H3yEoD8viYWlmS.asm',
 '02JqQ7H3yEoD8viYWlmS.bytes',
 '02K5GMYITj7bBoAisEmD.asm',
 '02K5GMYITj7bBoAisEmD.bytes',
 '02mlBLHZTDFXGa7Nt6cr.asm',
 '02mlBLHZTDFXGa7Nt6cr.bytes',
 '02MRILoE6rNhmt7FUi45.asm',
 '02MRILoE6rNhmt7FUi45.bytes',
 '02zcUmKV16Lya5xqnPGB.asm',
 '02zcUmKV16Lya5xqnPGB.bytes',
 '03nJaQV6K2ObICUmyWoR.asm',
 '03nJaQV6K2ObICUmyWoR.bytes',
 '04BfoQRA6XEshiNuI7pF.asm',
 '04BfoQRA6XEshiNuI7pF.bytes']
```


'0B2RwKm6dqyfjUWDN10a.bytes',
'0b5LqcWix3J4fGIEhXQu.asm',
'0b5LqcWix3J4fGIEhXQu.bytes',
'0BEsCP7NAUy8XmkenHWG.asm',
'0BEsCP7NAUy8XmkenHWG.bytes',
'0BFIPv1rO83whtpMYyAs.asm',
'0BFIPv1rO83whtpMYyAs.bytes',
'0BIdbVDEgmPwjYF4xzir.asm',
'0BIdbVDEgmPwjYF4xzir.bytes',
'0bjN3Kgw5OATSreRmEdi.asm',
'0bjN3Kgw5OATSreRmEdi.bytes',
'0BKcmNv4iGY2hsVSaXJ6.asm',
'0BKcmNv4iGY2hsVSaXJ6.bytes',
'04EjIdbPV5e1XroFOpiN.asm',
'06osXqPUVM1HbvBGNncT.asm',
'0AnoOZDNbPXIr2MRBSCJ.asm',
'0BLbmzJRkjNynCgQIdtV.asm',
'0CzL6rfwaTqGOu9eghBt.asm',
'0DTs2PhZfCwEv7q8349K.asm',
'0fvnGU7dkbr8iEhZuMcP.asm',
'0gSm7QZu5x6MBvVzUncH.asm',
'0HICT7RtjaVQzcNOeMgS.asm',
'0iBaz3krsQ8HuA7cGDSt.asm',
'0JAx9gzbC54Q61XBrqc7.asm',
'0KyDiQb1whgaSrmlx58J.asm',
'0mlhuKGpCc6OB4zwrLy.asm',
'0BLbmzJRkjNynCgQIdtV.bytes',
'0bN6ODYWw2xeCQBn3tEg.asm',
'0bN6ODYWw2xeCQBn3tEg.bytes',
'0BY2iPs03bEmudlUzpfq.asm',
'0BY2iPs03bEmudlUzpfq.bytes',
'0BZQIJak6Pu2tyAXfrzR.asm',
'0BZQIJak6Pu2tyAXfrzR.bytes',
'0C4aVbN58O1nAigFJt9z.asm',
'0C4aVbN58O1nAigFJt9z.bytes',
'0cdnSIvN489sFUwY1rMQ.asm',
'0cdnSIvN489sFUwY1rMQ.bytes',
'0cfGJLYgE6ROaZH7KT1h.asm',
'0cfGJLYgE6ROaZH7KT1h.bytes',
'0cfIE39ihRN02rkZOw5H.asm',
'0cfIE39ihRN02rkZOw5H.bytes',
'0cGWK6VvCkm7O2AxDjtW.asm',
'0cGWK6VvCkm7O2AxDjtW.bytes',
'0cH8YeO15ZywEhPrJvmj.asm',
'0cH8YeO15ZywEhPrJvmj.bytes',
'0co46B8IkPt2UN3HSaw7.asm',
'0co46B8IkPt2UN3HSaw7.bytes',
'0CPaAXtyswrBq83D6VEg.asm',
'0CPaAXtyswrBq83D6VEg.bytes',
'0Cq4wfhlrKBjiut11YAZ.asm',
'0Cq4wfhlrKBjiut11YAZ.bytes',
'0csgzpwdL3FbZEJu6DjO.asm',
'0csgzpwdL3FbZEJu6DjO.bytes',
'0cTu2bkfAOAJqIhYUWFK.asm',
'0cTu2bkfAOAJqIhYUWFK.bytes',
'0CzL6rfwaTqGOu9eghBt.bytes',
'0czUXKSCiGY2j5mxLdWa.asm',
'0czUXKSCiGY2j5mxLdWa.bytes',
'0D9IedmC1viTPugLRWX6.asm',
'0D9IedmC1viTPugLRWX6.bytes',
'0daTri9PSkeEsVHu5Dhw.asm',
'0daTri9PSkeEsVHu5Dhw.bytes',
'0dauMIK4ATfybzqUgNLc.asm',
'0dauMIK4ATfybzqUgNLc.bytes',
'0DbLeKSoxu47wjqVHsi9.asm',
'0DbLeKSoxu47wjqVHsi9.bytes',
'0df4cbsTBCn1VGW81QRv.asm',
'0df4cbsTBCn1VGW81QRv.bytes',
'0dhL8Jvcswa7U1qHiDS5.asm',
'0dhL8Jvcswa7U1qHiDS5.bytes',
'0Dk7Wd8MERu3b5rmQzCK.asm',
'0Dk7Wd8MERu3b5rmQzCK.bytes',
'0dkuzUXLTEFwW71vP5bS.asm',
'0dkuzUXLTEFwW71vP5bS.bytes',
'0DM3hS6Gg2QVkb1fZydv.asm',
'0DM3hS6Gg2QVkb1fZydv.bytes',
'0dnTix1lMYzDUp svEVrGc.asm',
.....

'0dnTixlMYzDUp svEVrGc.bytes',
'0DNVFKwYlcjO7bTfJ5p1.asm',
'0DNVFKwYlcjO7bTfJ5p1.bytes',
'0DqUX5rkg3IbMY6BLGCE.asm',
'0DqUX5rkg3IbMY6BLGCE.bytes',
'0Dt p59Av1RLifoK1Udm7.asm',
'0Dt p59Av1RLifoK1Udm7.bytes',
'0Dt s2PhZfCwEv7q8349K.bytes',
'0EAdHtLDypMcwjTFJziC.asm',
'0EAdHtLDypMcwjTFJziC.bytes',
'0eaNKwl uUm kYdIvZ923c.asm',
'0eaNKwl uUm kYdIvZ923c.bytes',
'0EL7OGZKozbiNCVP61gk.asm',
'0EL7OGZKozbiNCVP61gk.bytes',
'0eN9lyQfwmTVk7C2ZoYp.asm',
'0eN9lyQfwmTVk7C2ZoYp.bytes',
'0Eo9qT6idXHDMebwmvPA.asm',
'0Eo9qT6idXHDMebwmvPA.bytes',
'0evDQX7AVfC1ZTJEK1tg.asm',
'0evDQX7AVfC1ZTJEK1tg.bytes',
'0F4qIH aR7xOrm19Set3o.asm',
'0F4qIH aR7xOrm19Set3o.bytes',
'0FdOadWrfBU6TqwCRYxA.asm',
'0FdOadWrfBU6TqwCRYxA.bytes',
'0fGuCWgTraQ6nEmLPN8q.asm',
'0fGuCWgTraQ6nEmLPN8q.bytes',
'0fh nXI9ESr4jgWmkiaTe.asm',
'0fh nXI9ESr4jgWmkiaTe.bytes',
'0fHVZKeTE6iRb1PIQ4au.asm',
'0fHVZKeTE6iRb1PIQ4au.bytes',
'0FKerJ118x0c3jdoyg4A.asm',
'0FKerJ118x0c3jdoyg4A.bytes',
'0FOXjzmnD9CUMVcS1Eqh.asm',
'0FOXjzmnD9CUMVcS1Eqh.bytes',
'0Fu9oETtMW4zlg1ZrUy6.asm',
'0Fu9oETtMW4zlg1ZrUy6.bytes',
'0fvnGU7dkbr8iEhZuMcP.bytes',
'0fxgjYEC1PL1BDbcshzJ.asm',
'0fxgjYEC1PL1BDbcshzJ.bytes',
'0G2RV1chB1Ibk t6JqA5Q.asm',
'0G2RV1chB1Ibk t6JqA5Q.bytes',
'0G4hwobLuAzv11PWYfmd.asm',
'0G4hwobLuAzv11PWYfmd.bytes',
'0GbMkY1Nyt720zBjIcVh.asm',
'0GbMkY1Nyt720zBjIcVh.bytes',
'0gCmlyxw2UJvX7SNOGqu.asm',
'0gCmlyxw2UJvX7SNOGqu.bytes',
'0gcZkSFr7VnEmLPbTxUe.asm',
'0gcZkSFr7VnEmLPbTxUe.bytes',
'0gDsIvry1X5fPbG7cSBn.asm',
'0gDsIvry1X5fPbG7cSBn.bytes',
'0gHs6DEouiCPAc mWF rTX.asm',
'0gHs6DEouiCPAc mWF rTX.bytes',
'0giIqhw6e4mrHYzKF18T.asm',
'0giIqhw6e4mrHYzKF18T.bytes',
'0gkj92oIl eU4SYiCWpaM.asm',
'0gkj92oIl eU4SYiCWpaM.bytes',
'0GKp9ZJclxTABMuniOD2.asm',
'0GKp9ZJclxTABMuniOD2.bytes',
'0GKzFQ81IYXqUWkmfv26.asm',
'0GKzFQ81IYXqUWkmfv26.bytes',
'0gL3h5G6CszBV7RSinjJ.asm',
'0gL3h5G6CszBV7RSinjJ.bytes',
'0glscKoNakWL84EpunPe.asm',
'0glscKoNakWL84EpunPe.bytes',
'0gSm7QZu5x6MBvVzUncH.bytes',
'0Gu4misTcKynQD2O11Jx.asm',
'0Gu4misTcKynQD2O11Jx.bytes',
'0GUILi7xAlODwZ4YBenNM.asm',
'0GUILi7xAlODwZ4YBenNM.bytes',
'0gUpzkLVT73PCXx5WFRI.asm',
'0gUpzkLVT73PCXx5WFRI.bytes',
'0GuYe4J7oLwQ82xr3pWS.asm',
'0GuYe4J7oLwQ82xr3pWS.bytes',
'0GvcTdBQXWUJ2t7vjphN.asm',
'0GvcTdBQXWUJ2t7vjphN.bytes',
'0GvtWEPU BfDAcMbiYVSR.asm',

'0GvtWEPUBfDAcMbiYVSR.bytes',
'0gWUIudhwovMYb3NSnZA.asm',
'0gWUIudhwovMYb3NSnZA.bytes',
'0gxJ1YmwFUvnOzoM8N53.asm',
'0gxJ1YmwFUvnOzoM8N53.bytes',
'0H5OFk1Qm7wAPfE1qaJY.asm',
'0H5OFk1Qm7wAPfE1qaJY.bytes',
'0H63jydvIahOVqgx5Kfo.asm',
'0H63jydvIahOVqgx5Kfo.bytes',
'0hAlkjTR1Q6PewMczavb.asm',
'0hAlkjTR1Q6PewMczavb.bytes',
'0hBIiRpkMZtoYj31cDLA.asm',
'0hBIiRpkMZtoYj31cDLA.bytes',
'0HcZRmLi9VTpuQJCoXny.asm',
'0HcZRmLi9VTpuQJCoXny.bytes',
'0hH3JB2wM791YdsyuK5N.asm',
'0hH3JB2wM791YdsyuK5N.bytes',
'0HICT7RtjaVQzcNOeMgS.bytes',
'0HKFs3AXTt1IrO152eVu.asm',
'0HKFs3AXTt1IrO152eVu.bytes',
'0HKM38fmCR5DrxoIkBnQ.asm',
'0HKM38fmCR5DrxoIkBnQ.bytes',
'0Hlm4XgE1cQhC6BkMays.asm',
'0Hlm4XgE1cQhC6BkMays.bytes',
'0Hn2ojct97Wp1TbdNM4Q.asm',
'0Hn2ojct97Wp1TbdNM4Q.bytes',
'0Hrfce4X5YGESJPj19uL.asm',
'0Hrfce4X5YGESJPj19uL.bytes',
'0HVAAnMrp1LjKDmuoOJFY.asm',
'0HVAAnMrp1LjKDmuoOJFY.bytes',
'0hWRb28Umdgj7xcXOwtC.asm',
'0hWRb28Umdgj7xcXOwtC.bytes',
'0hZEqJ5eMVjU21HAG7Ii.asm',
'0hZEqJ5eMVjU21HAG7Ii.bytes',
'0hZqVRKkw7GfMdpaLLiN.asm',
'0hZqVRKkw7GfMdpaLLiN.bytes',
'0i4ENysvVrgFnbaHUuJK.asm',
'0i4ENysvVrgFnbaHUuJK.bytes',
'0i4FNJPQ8GuB3WU56LTS.asm',
'0i4FNJPQ8GuB3WU56LTS.bytes',
'0I4ZVvngsAatm8fzD3pk.asm',
'0I4ZVvngsAatm8fzD3pk.bytes',
'0iABvIkp3WHfgrJ79ymq.asm',
'0iABvIkp3WHfgrJ79ymq.bytes',
'0IALcuEiP9G6epb71Oom.asm',
'0IALcuEiP9G6epb71Oom.bytes',
'0iBaz3krsQ8HuA7cGDSt.bytes',
'0icJrNnmPvDqVQkC3we1.asm',
'0icJrNnmPvDqVQkC3we1.bytes',
'0IelgX5H2s14KutkEyNU.asm',
'0IelgX5H2s14KutkEyNU.bytes',
'0IMUK1Zs1Sm8LpGRkWhT.asm',
'0IMUK1Zs1Sm8LpGRkWhT.bytes',
'0iS3pwlgJco8XORD4TLq.asm',
'0iS3pwlgJco8XORD4TLq.bytes',
'0isdESDMzq2K8T6FLPcC.asm',
'0isdESDMzq2K8T6FLPcC.bytes',
'0itbI5mjJF28ocTkrUf9.asm',
'0itbI5mjJF28ocTkrUf9.bytes',
'0ItXLAUOhK8ZYdDf7HW4.asm',
'0ItXLAUOhK8ZYdDf7HW4.bytes',
'0Iv6U2hbcP1xeBitW5Oo.asm',
'0Iv6U2hbcP1xeBitW5Oo.bytes',
'0IyaidUKRqnt2PDFOzHT.asm',
'0IyaidUKRqnt2PDFOzHT.bytes',
'0IYZltU7uMpaco85PfKr.asm',
'0IYZltU7uMpaco85PfKr.bytes',
'0iTThuQ5KMb4RtAlrz6D.asm',
'0iTThuQ5KMb4RtAlrz6D.bytes',
'0J2pOclDKjadkL57eroz.asm',
'0J2pOclDKjadkL57eroz.bytes',
'0J61YGoWjV25TzxzeSluf.asm',
'0J61YGoWjV25TzxzeSluf.bytes',
'0jAopX6290wEH8WPkzVU.asm',
'0jAopX6290wEH8WPkzVU.bytes',
'0JAx9gzbC54Q61XBrqc7.bytes',
'0JAzwGUKORhFQWr3oldN.asm',

'0JAzwGUKORhFQWr3o1dN.bytes',
'0JBNEWmdi7GptrK5qYD9.asm',
'0JBNEWmdi7GptrK5qYD9.bytes',
'0JECiqrVNR1dgj67pZue.asm',
'0JECiqrVNR1dgj67pZue.bytes',
'0JfwyrEcBqaRzN9TgFMh.asm',
'0JfwyrEcBqaRzN9TgFMh.bytes',
'0jkmvR43UQ9yKxqXei61.asm',
'0jkmvR43UQ9yKxqXei61.bytes',
'0jKSsqXVHNucByZ916Ao.asm',
'0jKSsqXVHNucByZ916Ao.bytes',
'0JnvoeflBWWIcQa5GEPK.asm',
'0JnvoeflBWWIcQa5GEPK.bytes',
'0Job8TyN6VBGCrjAkzfP.asm',
'0Job8TyN6VBGCrjAkzfP.bytes',
'0JPAX13cjxewaTh6tRCi.asm',
'0JPAX13cjxewaTh6tRCi.bytes',
'0K4sTCLtrIJ5SinQbe7u.asm',
'0K4sTCLtrIJ5SinQbe7u.bytes',
'0K6yBUcTw3qjtNo4ZQpY.asm',
'0K6yBUcTw3qjtNo4ZQpY.bytes',
'0KgE6ksUeytoHf12cT4r.asm',
'0KgE6ksUeytoHf12cT4r.bytes',
'0KigmP9TLwJXNGz26tfO.asm',
'0KigmP9TLwJXNGz26tfO.bytes',
'0KLUAMqmPJhOwaYrbSCE.asm',
'0KLUAMqmPJhOwaYrbSCE.bytes',
'0KyDiQb1whgaSrmlx58J.bytes',
'0KZFcsOYR4MdPJf6VvGS.asm',
'0KZFcsOYR4MdPJf6VvGS.bytes',
'0kzRDUmBLHGd4YPj7hO6.asm',
'0kzRDUmBLHGd4YPj7hO6.bytes',
'015IobyKpuqcwo4NxfgD.asm',
'015IobyKpuqcwo4NxfgD.bytes',
'016fhCty3aSLDOgAjYQi.asm',
'016fhCty3aSLDOgAjYQi.bytes',
'0LAXajqhQy7po16dw8Tx.asm',
'0LAXajqhQy7po16dw8Tx.bytes',
'0LQSi5wnRZ3muIs6Mx9E.asm',
'0LQSi5wnRZ3muIs6Mx9E.bytes',
'0LVqv1HF8PuepodIiBUb.asm',
'0LVqv1HF8PuepodIiBUb.bytes',
'0LZkc7qeS39TUtVHuJB1.asm',
'0LZkc7qeS39TUtVHuJB1.bytes',
'0M7aSiE9csDzkmfKheVt.asm',
'0M7aSiE9csDzkmfKheVt.bytes',
'0m94tRnhgpsAUuY1L8KC.asm',
'0m94tRnhgpsAUuY1L8KC.bytes',
'0mcWyK6unLRGV8Hfr97Y.asm',
'0mcWyK6unLRGV8Hfr97Y.bytes',
'0meUjiuJvODcf3k9z4Iy.asm',
'0meUjiuJvODcf3k9z4Iy.bytes',
'0mfwtlekXE1poYAnqMRO.asm',
'0mfwtlekXE1poYAnqMRO.bytes',
'0mgFnqeLAMr7jthUYRTv.asm',
'0mgFnqeLAMr7jthUYRTv.bytes',
'0mlhuKGpCc6OB4zwrblY.bytes',
'0MmZ8j5pn2R3VG9wlxYi.asm',
'0MmZ8j5pn2R3VG9wlxYi.bytes',
'0MOorvEIRmZGhqQdc3TA.asm',
'0MOorvEIRmZGhqQdc3TA.bytes',
'0MpJYhdbf8T7InoqcXr1.asm',
'0MpJYhdbf8T7InoqcXr1.bytes',
'0MPV9Y8WNcFoyRZqQ76G.asm',
'0MPV9Y8WNcFoyRZqQ76G.bytes',
'0MQD6mnoy413zV8WPRYe.asm',
'0MQD6mnoy413zV8WPRYe.bytes',
'0MsQ16zDg9XVerbmfcDdU.asm',
'0MsQ16zDg9XVerbmfcDdU.bytes',
'0Mx7E5XgoRcJavdPAit8.asm',
'0Mx7E5XgoRcJavdPAit8.bytes',
'0NEsQ1DGnUMg3Bew7R1A.asm',
'0NEsQ1DGnUMg3Bew7R1A.bytes',
'0NiOTDde1ktgx954SJFE.asm',
'0NiOTDde1ktgx954SJFE.bytes',
'0njs7MJQObCY8ABgykiP.asm',
'0njs7MJQObCY8ABgykiP.bytes',

'0NmoAEOtIDdwiVr9PCBf.asm',
'0NmoAEOtIDdwiVr9PCBf.bytes',
'ONXFnJyOEhBAISKfiU67.asm',
'ONXFnJyOEhBAISKfiU67.bytes',
'0nxrvcZJUBNGM8Vg4SRf.asm',
'0nxrvcZJUBNGM8Vg4SRf.bytes',
'0NyfG Xt8nmlK72Q9Irhs.asm',
'0NyfG Xt8nmlK72Q9Irhs.bytes',
'0NZT2sFp4JKXqYGC8Rna.bytes',
'0odUVkrjp2B1n8NDS6bR.asm',
'0odUVkrjp2B1n8NDS6bR.bytes',
'0OfazrXvcTqVVi3dEGU9.asm',
'0OfazrXvcTqVVi3dEGU9.bytes',
'0OQ9MEykugoYZJrnj64d.asm',
'0OQ9MEykugoYZJrnj64d.bytes',
'0tC2dSxuGDzY6XHLkq3.asm',
'0tC2dSxuGDzY6XHLkq3.bytes',
'0OXtHMVdxJuvezcLRGZQ.asm',
'0OXtHMVdxJuvezcLRGZQ.bytes',
'0P6tEopzr2mIhkuOKiCD.asm',
'0P6tEopzr2mIhkuOKiCD.bytes',
'0p8RCBnvzTYM1Id9SoZL.asm',
'0p8RCBnvzTYM1Id9SoZL.bytes',
'0pER9ak46CKVjPoqGLcD.asm',
'0pER9ak46CKVjPoqGLcD.bytes',
'0PGoOfXAjVDBBytkRH4Um.asm',
'0PGoOfXAjVDBBytkRH4Um.bytes',
'0Phttp2LVcsCFMKkGgmRH.asm',
'0Phttp2LVcsCFMKkGgmRH.bytes',
'0PlfqyKM1JtYzx2me5FU.asm',
'0PlfqyKM1JtYzx2me5FU.bytes',
'0pLrSeQMbqfCc2IZHUAo.asm',
'0pLrSeQMbqfCc2IZHUAo.bytes',
'0pnRqtHF3dKI5MQmwT29.asm',
'0pnRqtHF3dKI5MQmwT29.bytes',
'0pQ1231KqiLJtFTTwjBfG.asm',
'0pQ1231KqiLJtFTTwjBfG.bytes',
'0pqtVjvLk89PdYr5MwZO.bytes',
'0PREAmkFZjrftD3Li wOV.asm',
'0PREAmkFZjrftD3Li wOV.bytes',
'0pTO4SVnWDehgU1YBAvq.asm',
'0pTO4SVnWDehgU1YBAvq.bytes',
'0PVAd2O41MGUbawzfXBT.asm',
'0PVAd2O41MGUbawzfXBT.bytes',
'0pXDqTLSebKWgaUOdnVE.asm',
'0pXDqTLSebKWgaUOdnVE.bytes',
'0Q4ALVSRnlHUBjyOb1sw.asm',
'0Q4ALVSRnlHUBjyOb1sw.bytes',
'0q6eABpEH SbX1IzGh7jc.asm',
'0q6eABpEH SbX1IzGh7jc.bytes',
'0qDikloYBvpQTcCOMyS4.asm',
'0qDikloYBvpQTcCOMyS4.bytes',
'0qeE63wvCGzSapRuW4TI.asm',
'0qeE63wvCGzSapRuW4TI.bytes',
'0Qfd8jy2JYX3U7gFBG5D.asm',
'0Qfd8jy2JYX3U7gFBG5D.bytes',
'0QIZ6hemJLErWdFlGnCc.asm',
'0QIZ6hemJLErWdFlGnCc.bytes',
'0qjuDC7Rh x9rHkLlItAp.asm',
'0qjuDC7Rh x9rHkLlItAp.bytes',
'0QK142N5HoBntC7RydAh.asm',
'0QK142N5HoBntC7RydAh.bytes',
'0qN2U7Ayorv9IbiTW4Gd.asm',
'0qN2U7Ayorv9IbiTW4Gd.bytes',
'0qNtRXvP8F2431bMmpse.asm',
'0qNtRXvP8F2431bMmpse.bytes',
'0qPGt4cRVk9NoiJgubf2.bytes',
'0qpMX5CZU8mgD2kaRNv v.asm',
'0qpMX5CZU8mgD2kaRNv v.bytes',
'0qPYfxdvNrBX6ASbZltF.asm',
'0qPYfxdvNrBX6ASbZltF.bytes',
'0QqWIUsiyj3AaPwJD81S.asm',
'0QqWIUsiyj3AaPwJD81S.bytes',
'0qV43bnIcBDPF1WlG6tJ.asm',
'0qV43bnIcBDPF1WlG6tJ.bytes',
'0QWUCHwNSFt1Ojr45JpT.asm',
'0QWUCHwNSFt1Ojr45JpT.bytes',

'0R3rxkTidWOuPI5MUbL6.asm',
'0R3rxkTidWOuPI5MUbL6.bytes',
'0raU6iGvjZREQJYB1HX5.asm',
'0raU6iGvjZREQJYB1HX5.bytes',
'0RCcpJ1DuKNMBSaQUgmA.asm',
'0RCcpJ1DuKNMBSaQUgmA.bytes',
'0Rgetc1wAfxlzHTGBoa7.asm',
'0Rgetc1wAfxlzHTGBoa7.bytes',
'0rgi8OfjvwKbA57IXVcE.asm',
'0rgi8OfjvwKbA57IXVcE.bytes',
'0rG1Kyvn5Vo67RaDhPti.asm',
'0rG1Kyvn5Vo67RaDhPti.bytes',
'0rgudc7PpbexCtBjNqWF.asm',
'0rgudc7PpbexCtBjNqWF.bytes',
'0rlkeuBLpWm15TyFRHUJ.asm',
'0rlkeuBLpWm15TyFRHUJ.bytes',
'0RLWbvBk9rEtVZ74KAen.asm',
'0RLWbvBk9rEtVZ74KAen.bytes',
'0RSW8EmilbznyYPIrgvD.bytes',
'0RvUg8rntNxyDHPJeEYG.asm',
'0RvUg8rntNxyDHPJeEYG.bytes',
'0S7z4qxYTHPUDO8fgIyA.asm',
'0S7z4qxYTHPUDO8fgIyA.bytes',
'0SaYqdFRA5G7rfktjT63.asm',
'0SaYqdFRA5G7rfktjT63.bytes',
'0ScbCK2aI4xrwoALDmNd.asm',
'0ScbCK2aI4xrwoALDmNd.bytes',
'0sdr2cfHLzNO9Ep4RSjt.asm',
'0sdr2cfHLzNO9Ep4RSjt.bytes',
'0sH5GdO3LpZwMbPaXjfN.asm',
'0sH5GdO3LpZwMbPaXjfN.bytes',
'0sh9tOC5GprAyciTmaXW.asm',
'0sh9tOC5GprAyciTmaXW.bytes',
'0SHJ2aOnMCUGdlZrtQDV.asm',
'0SHJ2aOnMCUGdlZrtQDV.bytes',
'0SHsZKUkI2iN1Rw4be6E.asm',
'0SHsZKUkI2iN1Rw4be6E.bytes',
'0SIqXMkDBWKE8s9pQaOv.asm',
'0SIqXMkDBWKE8s9pQaOv.bytes',
'0SkbpFZGvraL6fQXTzol.asm',
'0SkbpFZGvraL6fQXTzol.bytes',
'0sLBvJeYl5UkiRzVlatI.asm',
'0sLBvJeYl5UkiRzVlatI.bytes',
'0sM2atqTo48GgnFwhQXu.asm',
'0sM2atqTo48GgnFwhQXu.bytes',
'0sM6DQlxcP3oAfd9ZVBT.asm',
'0sM6DQlxcP3oAfd9ZVBT.bytes',
'0sMxy5SNLu1YCkKJzaB3.bytes',
'0SUAYQWoTeZb8qmIPMK3.asm',
'0SUAYQWoTeZb8qmIPMK3.bytes',
'0sYDNzwVMaWE2TIZeA81.asm',
'0sYDNzwVMaWE2TIZeA81.bytes',
'0TBj7g1MRJScN9GmOzyi.asm',
'0TBj7g1MRJScN9GmOzyi.bytes',
'0tNShPVGYWjcmT2yXLdq.asm',
'0tNShPVGYWjcmT2yXLdq.bytes',
'0TpWNF7ULhd4Yk5bXQyP.asm',
'0TpWNF7ULhd4Yk5bXQyP.bytes',
'0TqZzRJOeno9NKpHIba5.asm',
'0TqZzRJOeno9NKpHIba5.bytes',
'0Tw8aQsVxM1ZgFDPUm21.asm',
'0Tw8aQsVxM1ZgFDPUm21.bytes',
'0tZy7jgOKCIJs1PRvrop.asm',
'0tZy7jgOKCIJs1PRvrop.bytes',
'0u3DIJTqnpkg187VSbOL.asm',
'0u3DIJTqnpkg187VSbOL.bytes',
'0u3fjhzJ2ZMyAgKYkVF9.asm',
'0u3fjhzJ2ZMyAgKYkVF9.bytes',
'0u4QcXSvNvYnMalOBCygH.asm',
'0u4QcXSvNvYnMalOBCygH.bytes',
'0UcDbq2yojiThfdLRNS4.asm',
'0UcDbq2yojiThfdLRNS4.bytes',
'0uDv25KpjABxvQ1rXgzT.asm',
'0uDv25KpjABxvQ1rXgzT.bytes',
'0ufTUZYESk6MepA7vzKJ.asm',
'0ufTUZYESk6MepA7vzKJ.bytes',
'0uLJEcpV1IWkTBeHZU4w.bytes',

'0uNkt6sirCnUWw175pj1.asm',
'0uNkt6sirCnUWw175pj1.bytes',
'0URC4HAzJrjWnDVMXF1x.asm',
'0URC4HAzJrjWnDVMXF1x.bytes',
'0UTxRyirCXl6tsSE5GAv.asm',
'0UTxRyirCXl6tsSE5GAv.bytes',
'0UwHzzuXy7GFStRKEgdx.asm',
'0UwHzzuXy7GFStRKEgdx.bytes',
'0UYATSktwajsio8LWpx2.asm',
'0UYATSktwajsio8LWpx2.bytes',
'0uZcAvaHLitbWmC1QMXf.asm',
'0uZcAvaHLitbWmC1QMXf.bytes',
'0V6UmdMg7wQGXs1AzORJ.asm',
'0V6UmdMg7wQGXs1AzORJ.bytes',
'0VfuK267xTdeBcmLaCp1.asm',
'0VfuK267xTdeBcmLaCp1.bytes',
'0Vi6RsKdlHBzySe5Zp4Y.asm',
'0Vi6RsKdlHBzySe5Zp4Y.bytes',
'0VLen4rGI1wNj0qzyZRb.asm',
'0VLen4rGI1wNj0qzyZRb.bytes',
'0VnTurd7B4YsH63wUGJ2.asm',
'0VnTurd7B4YsH63wUGJ2.bytes',
'0W1RChtwZvj4Qy78GYUJ.asm',
'0W1RChtwZvj4Qy78GYUJ.bytes',
'0w2L1G4eiBfby9JIMZgS.asm',
'0w2L1G4eiBfby9JIMZgS.bytes',
'0WdoYq78xDkFMcIwRpmJ.asm',
'0WdoYq78xDkFMcIwRpmJ.bytes',
'0wkFafuhqSr1Z5J7K6C9.bytes',
'0wmZBjQenOTsu5bpSkf6.asm',
'0wmZBjQenOTsu5bpSkf6.bytes',
'0wqQj9ngaAVxI4SLr5H7.asm',
'0wqQj9ngaAVxI4SLr5H7.bytes',
'0WQtf1pNPdRqUI7KJFAT.asm',
'0WQtf1pNPdRqUI7KJFAT.bytes',
'0wtIHxfBDmYRQ7GWSaTj.asm',
'0wtIHxfBDmYRQ7GWSaTj.bytes',
'0x47zZTBynC21leohIEs.asm',
'0x47zZTBynC21leohIEs.bytes',
'0X6U3SfcmbuZNodnH9qL.asm',
'0X6U3SfcmbuZNodnH9qL.bytes',
'0X8fojtR6QIB1ypvHlsM.asm',
'0X8fojtR6QIB1ypvHlsM.bytes',
'0XbS5KBuICayhs2p3NTU.asm',
'0XbS5KBuICayhs2p3NTU.bytes',
'0XEmjdgyY6hWFkulvCi.asm',
'0XEmjdgyY6hWFkulvCi.bytes',
'0xfjbVBPYX1pSt4rdWOe.asm',
'0xfjbVBPYX1pSt4rdWOe.bytes',
'0xkbZDBOAicQJK5wC7Pm.asm',
'0xkbZDBOAicQJK5wC7Pm.bytes',
'0xLRN4DuvoCGbjIFQKYB.asm',
'0xLRN4DuvoCGbjIFQKYB.bytes',
'0Xm5sqBgQASUyFui2kbH.asm',
'0Xm5sqBgQASUyFui2kbH.bytes',
'0XnPcouGZdhIKLiOFQHr.asm',
'0XnPcouGZdhIKLiOFQHr.bytes',
'0xpa8zDNdEnK9gB1Wmer.bytes',
'0XrH26wcU1ASSvK5ZqhD.asm',
'0XrH26wcU1ASSvK5ZqhD.bytes',
'0XSHfwZAgYm7eFVINuDT.asm',
'0XSHfwZAgYm7eFVINuDT.bytes',
'0XsnKFjqPQxzCf81DvWp.asm',
'0XsnKFjqPQxzCf81DvWp.bytes',
'0xC2vRzyYtqFa4VhLnG.asm',
'0xC2vRzyYtqFa4VhLnG.bytes',
'0XVbihU4clgRS9A6tsB2.asm',
'0XVbihU4clgRS9A6tsB2.bytes',
'0yBQ5dAc2gN4qx9YWuiS.asm',
'0yBQ5dAc2gN4qx9YWuiS.bytes',
'0YCvp8xd4rjiDMKfv9PA.asm',
'0YCvp8xd4rjiDMKfv9PA.bytes',
'0YIPlyX3fbNs2ntJRGgH.asm',
'0YIPlyX3fbNs2ntJRGgH.bytes',
'0yjkG4BHeCbTimIslgr7.asm',
'0yjkG4BHeCbTimIslgr7.bytes',
'0YJRabHQ1reDxN3yjUCd.asm',

'0YJRabHQ1reDxN3yjUCd.bytes',
'0yn4gphuOY1PsCz5JGFM.asm',
'0yn4gphuOY1PsCz5JGFM.bytes',
'0YRDpBU67WrqEugkXlnd.asm',
'0YRDpBU67WrqEugkXlnd.bytes',
'0yvBxsufKXzM6CwpcSYU.asm',
'0yvBxsufKXzM6CwpcSYU.bytes',
'0YWidX9hOD5sPrtTvc2M.asm',
'0YWidX9hOD5sPrtTvc2M.bytes',
'0Z3nQDPiU7LMEesmG4TB.bytes',
'0ZAJ1qmphO4fjnVB6bMz.asm',
'0ZAJ1qmphO4fjnVB6bMz.bytes',
'0Zb69USnM3qafJTWHCpL.asm',
'0Zb69USnM3qafJTWHCpL.bytes',
'0zfLRXlhm34ASQUtrdjV.asm',
'0zfLRXlhm34ASQUtrdjV.bytes',
'0ZHV6acpJ9KkAWPjEI71.asm',
'0ZHV6acpJ9KkAWPjEI71.bytes',
'0ZiAjFN81BOE7plhCnSm.asm',
'0ZiAjFN81BOE7plhCnSm.bytes',
'0ZiQmgtxzHe9v508Lf2k.asm',
'0ZiQmgtxzHe9v508Lf2k.bytes',
'0ZNejapCUtVsq41x5FDm.asm',
'0ZNejapCUtVsq41x5FDm.bytes',
'0zr6vIw74DgVSKXiRUY1.asm',
'0zr6vIw74DgVSKXiRUY1.bytes',
'0ZTEyLXaWReMK3rYVCjv.asm',
'0ZTEyLXaWReMK3rYVCjv.bytes',
'10Esk16DUSGugoWxvtYO.asm',
'10Esk16DUSGugoWxvtYO.bytes',
'10H9detwuiZByNSRMmxsf.asm',
'10H9detwuiZByNSRMmxsf.bytes',
'10oTLIHtSwzCN45DysmQ.asm',
'10oTLIHtSwzCN45DysmQ.bytes',
'125y4VsArzkCNOGZfu6o.asm',
'125y4VsArzkCNOGZfu6o.bytes',
'129LWuNOIS0RjQnq5om6.asm',
'129LWuNOIS0RjQnq5om6.bytes',
'12Jd4qpOzTtQC3E6PXDb.bytes',
'12pnF9KehUmEk4SZP3Lj.asm',
'12pnF9KehUmEk4SZP3Lj.bytes',
'12TEseaJ4jnbyORqKxhf.asm',
'12TEseaJ4jnbyORqKxhf.bytes',
'12tjh4qCkcHpObVBEeMr.asm',
'12tjh4qCkcHpObVBEeMr.bytes',
'13FZ9DOhywk78sJiVULu.asm',
'13FZ9DOhywk78sJiVULu.bytes',
'13HkPY29CL8Rzs0FOc5D.asm',
'13HkPY29CL8Rzs0FOc5D.bytes',
'13TchiQbYmGgO7aufZyE.asm',
'13TchiQbYmGgO7aufZyE.bytes',
'13YpdP5vTLOazSQFRgJn.asm',
'13YpdP5vTLOazSQFRgJn.bytes',
'148bgAqCj3NzaHwSP0Gu.asm',
'148bgAqCj3NzaHwSP0Gu.bytes',
'14DUIwIb5xrfkF30cThJ.asm',
'14DUIwIb5xrfkF30cThJ.bytes',
'141JRXiAchNqVFSOBLgb.asm',
'141JRXiAchNqVFSOBLgb.bytes',
'15ADIkydX7eNq90WfcB.asm',
'15ADIkydX7eNq90WfcB.bytes',
'15EAacDsS3Nj91QTfLPu.asm',
'15EAacDsS3Nj91QTfLPu.bytes',
'15loeAHtkJa8BuFi6Zry.asm',
'15loeAHtkJa8BuFi6Zry.bytes',
'15mfGKRS6aVevjAoYL8w.asm',
'15mfGKRS6aVevjAoYL8w.bytes',
'15sGnFeEvMIgpQ8acbUu.bytes',
'15tQNEdPb0mi63kfOWM8.asm',
'15tQNEdPb0mi63kfOWM8.bytes',
'15uCdAByiF48LpZ3YS6Q.asm',
'15uCdAByiF48LpZ3YS6Q.bytes',
'16cTMtKIjH5SbyovWuBq.asm',
'16cTMtKIjH5SbyovWuBq.bytes',
'16eaj17F2vQKhwi9dM4O.asm',
'16eaj17F2vQKhwi9dM4O.bytes',
'16gDhwBZTRjLMO9zI5cX.asm',

'16gDhwBZTRjLM09zI5cX.bytes',
'16KuXyrdDOCWGhc7P4kL.asm',
'16KuXyrdDOCWGhc7P4kL.bytes',
'175ZU8C1Osnifz03ybcX.asm',
'175ZU8C1Osnifz03ybcX.bytes',
'17aN50JFZXYrcAk2iCtG.asm',
'17aN50JFZXYrcAk2iCtG.bytes',
'17jNaPkYIgvUyw9HxC5D.asm',
'17jNaPkYIgvUyw9HxC5D.bytes',
'17kCoALmxaxpcNzd5KnE.asm',
'17kCoALmxaxpcNzd5KnE.bytes',
'17qvV08gWM3LtoZKS9d4.asm',
'17qvV08gWM3LtoZKS9d4.bytes',
'17zD8kITZON1BWYdfvr6.asm',
'17zD8kITZON1BWYdfvr6.bytes',
'18A9fzbeyripCSuLKo4.asm',
'18A9fzbeyripCSuLKo4.bytes',
'18AD51V4StiZfn0pmdGj.asm',
'18AD51V4StiZfn0pmdGj.bytes',
'18eZt9qWksQhoY3K60aE.bytes',
'18S9URJgMWqua2TZeVsh.asm',
'18S9URJgMWqua2TZeVsh.bytes',
'18SQsaz61MJIZCknYpgL.asm',
'18SQsaz61MJIZCknYpgL.bytes',
'18T9jLsd0DZQc2aRrSpb.asm',
'18T9jLsd0DZQc2aRrSpb.bytes',
'18te652nLyCi3IrZqMOJ.asm',
'18te652nLyCi3IrZqMOJ.bytes',
'194Iy5xv8QRz7XMTnmAk.asm',
'194Iy5xv8QRz7XMTnmAk.bytes',
'198PdXRpGVZAzYxCrF0e.asm',
'198PdXRpGVZAzYxCrF0e.bytes',
'19b2d1EtWXyeiqc67Bof.asm',
'19b2d1EtWXyeiqc67Bof.bytes',
'19B6wSHOLmjhgPoquTbW.asm',
'19B6wSHOLmjhgPoquTbW.bytes',
'19BtdpKYGOSyX3oZsirc.asm',
'19BtdpKYGOSyX3oZsirc.bytes',
'19Emlw3WgC2O4iJknVtP.asm',
'19Emlw3WgC2O4iJknVtP.bytes',
'19emSiJG8vrZgt2zXEIY.asm',
'19emSiJG8vrZgt2zXEIY.bytes',
'19hXQ20Dt6WvpEdxTGmo.asm',
'19hXQ20Dt6WvpEdxTGmo.bytes',
'191TEvHhVDprkGxgXLcf.asm',
'191TEvHhVDprkGxgXLcf.bytes',
'19UEqxPS8Go3MTNci4Hh.asm',
'19UEqxPS8Go3MTNci4Hh.bytes',
'19xM5LfBYviATcUkEOPm.bytes',
'19zYbuW3XONcEedv7xU1.asm',
'19zYbuW3XONcEedv7xU1.bytes',
'1A7UqoIMrxHgVuW3FS9c.asm',
'1A7UqoIMrxHgVuW3FS9c.bytes',
'1aAwe4J9VHrsq8uEoZhf.asm',
'1aAwe4J9VHrsq8uEoZhf.bytes',
'1aDvVlWtPApk6EuKLdCG.asm',
'1aDvVlWtPApk6EuKLdCG.bytes',
'1AFY9HiVpuymkQOeoDMU.asm',
'1AFY9HiVpuymkQOeoDMU.bytes',
'1AfYJKywcXSII6H80nqE.asm',
'1AfYJKywcXSII6H80nqE.bytes',
'1AhLoSNdHzDPfJrSXmQC.asm',
'1AhLoSNdHzDPfJrSXmQC.bytes',
'1AhQif5kMPb24RnCmjrg.asm',
'1AhQif5kMPb24RnCmjrg.bytes',
'1aJLj5T8sUYbM2trvX4n.asm',
'1aJLj5T8sUYbM2trvX4n.bytes',
'1ARrSUKCbhWeTFyjp4ix.asm',
'1ARrSUKCbhWeTFyjp4ix.bytes',
'1AvCEkV76rjuYp9aqPbK.asm',
'1AvCEkV76rjuYp9aqPbK.bytes',
'1aVxjv19Y5qNk8GbBsgE.asm',
'1aVxjv19Y5qNk8GbBsgE.bytes',
'1awKqmMpQUblukHhBy7n.asm',
'1awKqmMpQUblukHhBy7n.bytes',
'1awXjkTQiJpCeDV7dL8m.asm',
'1awXjkTQiJpCeDV7dL8m.bytes',

'1AXEYUDVt42T6ZFx8COo.bytes',
'1AXVt0D9mSU26fL4iFGc.asm',
'1AXVt0D9mSU26fL4iFGc.bytes',
'1BDCxE2AMQN9yH6sWdIo.asm',
'1BDCxE2AMQN9yH6sWdIo.bytes',
'1bgW3y0fFDHBjdps5Jhe.asm',
'1bgW3y0fFDHBjdps5Jhe.bytes',
'1bkniqgYzRMpZ2S5JBAQ.asm',
'1bkniqgYzRMpZ2S5JBAQ.bytes',
'1bL4yiwCUvSOg7tBJudf.asm',
'1bL4yiwCUvSOg7tBJudf.bytes',
'1BM418DCZmSpoYwE9Qcv.asm',
'1BM418DCZmSpoYwE9Qcv.bytes',
'1buWc2A0VeMrkNpohCPy.asm',
'1buWc2A0VeMrkNpohCPy.bytes',
'1BYZOXRoVfhtQGMJnkST.asm',
'1BYZOXRoVfhtQGMJnkST.bytes',
'1cB6EdrmebFjAX3QHkS9.asm',
'1cB6EdrmebFjAX3QHkS9.bytes',
'1CbXMgAIJtiNLR0So9ul.asm',
'1CbXMgAIJtiNLR0So9ul.bytes',
'1cbYnSz3otKJCKijjqAWV.asm',
'1cbYnSz3otKJCKijjqAWV.bytes',
'1cdtwSyhmiekLoH7OUEs.asm',
'1cdtwSyhmiekLoH7OUEs.bytes',
'1cEKCiXn0hLVeSk2dDyW.asm',
'1cEKCiXn0hLVeSk2dDyW.bytes',
'1CG2F537YPzwlraIoykc.asm',
'1CG2F537YPzwlraIoykc.bytes',
'1cgZSaQIy6UBTAqsdX9V.bytes',
'1Ch0ZFDqT9PzkGbSeyEV.asm',
'1Ch0ZFDqT9PzkGbSeyEV.bytes',
'1ckp5Ye83yzjIilWQCfm.asm',
'1ckp5Ye83yzjIilWQCfm.bytes',
'1cP8XiEvMd7tg3rJmnAj.asm',
'1cP8XiEvMd7tg3rJmnAj.bytes',
'1ctXu0GAOqqj23STKWV7s.asm',
'1ctXu0GAOqqj23STKWV7s.bytes',
'1Cub9LvOgMWkQqs3R2ef.asm',
'1Cub9LvOgMWkQqs3R2ef.bytes',
'1CWtrEfORZDj7psmnToA.asm',
'1CWtrEfORZDj7psmnToA.bytes',
'1cX2xrbYwzdaM4mWuoCl.asm',
'1cX2xrbYwzdaM4mWuoCl.bytes',
'1d485UYC6qhXKrAcDJPZ.asm',
'1d485UYC6qhXKrAcDJPZ.bytes',
'1D6IwN4EdfHta1YJuA95.asm',
'1D6IwN4EdfHta1YJuA95.bytes',
'1d79vVgNkt4fMjQELaw3.asm',
'1d79vVgNkt4fMjQELaw3.bytes',
'1dFAk8YzcKO4RVmsaS9Z.asm',
'1dFAk8YzcKO4RVmsaS9Z.bytes',
'1dhuKsA0t2Dnf9qNQLj1.asm',
'1dhuKsA0t2Dnf9qNQLj1.bytes',
'1dJWqVSQnsAKBXNkgy6j.asm',
'1dJWqVSQnsAKBXNkgy6j.bytes',
'1DmFGLjAlhPwn4yEVteq.asm',
'1DmFGLjAlhPwn4yEVteq.bytes',
'0pqtVjvLk89PdYr5MwZO.asm',
'0qPGt4cRVk9NoiJgubf2.asm',
'0RSW8EmilbznyYPIrgvD.asm',
'0sMxy5SNLu1YCkKJzaB3.asm',
'0uLJEcpVlIWkTBeHZU4w.asm',
'0wkFaFuhqSr1Z5J7K6C9.asm',
'0xpa8zDNdEnK9gB1Wmer.asm',
'0Z3nQDPiU7LMEesmG4TB.asm',
'12Jd4qpOzTtQC3E6PXDb.asm',
'15sGnFeEvMIgpQ8acbUu.asm',
'18eZt9qWksQhoY3K60aE.asm',
'19xM5LfBYviATcUkEOPm.asm',
'1AXEYUDVt42T6ZFx8COo.asm',
'1DNWriJEG56S03yvjTBu.bytes',
'1dxLT9no8t6fbykXEJiR.asm',
'1dxLT9no8t6fbykXEJiR.bytes',
'1e4vEcCMmyGJ3wrPXTq8.asm',
'1e4vEcCMmyGJ3wrPXTq8.bytes',
'1e8SrDV6mBaLUJOXOK2w.asm'

'1e8SrDV6mBgLUJQXOK2w.bytes',
'1E93CpP60RHFNiT5Qfvn.asm',
'1E93CpP60RHFNiT5Qfvn.bytes',
'1eAC0zqkIGHRBvm97P6i.asm',
'1eAC0zqkIGHRBvm97P6i.bytes',
'1eCQkhvgDSFNXYLJbAdj.asm',
'1eCQkhvgDSFNXYLJbAdj.bytes',
'1EDNRfXhHW17txAOdycU.asm',
'1EDNRfXhHW17txAOdycU.bytes',
'1EisIfzGOTX8hKaFH6uy.asm',
'1EisIfzGOTX8hKaFH6uy.bytes',
'1eJx3418pcAFvMuOwrjB.asm',
'1eJx3418pcAFvMuOwrjB.bytes',
'1el4qZPVwIXavrChcyGp.asm',
'1el4qZPVwIXavrChcyGp.bytes',
'1eOaAY4fpV38LIdhx195.asm',
'1eOaAY4fpV38LIdhx195.bytes',
'1ESMN0Gc6wRmC9BFPjWy.asm',
'1ESMN0Gc6wRmC9BFPjWy.bytes',
'1esMpo6SAPHiTUGEb3Bt.asm',
'1esMpo6SAPHiTUGEb3Bt.bytes',
'1EWLhwV47I6HYktzXsmS.asm',
'1EWLhwV47I6HYktzXsmS.bytes',
'1ezc5qf2alj08CsSw7Er.bytes',
'1f0427Tk3POoeFMpG9IR.asm',
'1f0427Tk3POoeFMpG9IR.bytes',
'1F7jeznmCoiBxNdIKsWT.asm',
'1F7jeznmCoiBxNdIKsWT.bytes',
'1FacC02JPfxSdXeD7MEw.asm',
'1FacC02JPfxSdXeD7MEw.bytes',
'1FAxSngsym35phwROG10.asm',
'1FAxSngsym35phwROG10.bytes',
'1Fe4Dk5vSltEWj3Kpg97.asm',
'1Fe4Dk5vSltEWj3Kpg97.bytes',
'1fJqS7HUCEa6PYx0ZLc9.asm',
'1fJqS7HUCEa6PYx0ZLc9.bytes',
'1fJVGN3On017kRohcF2E.asm',
'1fJVGN3On017kRohcF2E.bytes',
'1Flm7qajnOzdtwBZ29gL.asm',
'1Flm7qajnOzdtwBZ29gL.bytes',
'1fMQAqC26xFJIgKpkrwS.asm',
'1fMQAqC26xFJIgKpkrwS.bytes',
'1FsfvPlaKG3TpImqOejM.asm',
'1FsfvPlaKG3TpImqOejM.bytes',
'1fZPmkixLgY8KqnOQd5E.asm',
'1fZPmkixLgY8KqnOQd5E.bytes',
'1gaQlVW56mCk4NxrhPw.asm',
'1gaQlVW56mCk4NxrhPw.bytes',
'1GKhFnumAdIeykS3rZYi.asm',
'1GKhFnumAdIeykS3rZYi.bytes',
'1GmZULtOTHsYE4sNWDC.asm',
'1GmZULtOTHsYE4sNWDC.bytes',
'1GQ41zmoS3P96jkAKsRn.bytes',
'1gQDa419cNBbkO7n5IY0.asm',
'1gQDa419cNBbkO7n5IY0.bytes',
'1GqxsQhZ3FutTnbghPaf.asm',
'1GqxsQhZ3FutTnbghPaf.bytes',
'1GRsOmJXetEkqIraBuSd.asm',
'1GRsOmJXetEkqIraBuSd.bytes',
'1gtbK8Nhsx3ICnoD9BXY.asm',
'1gtbK8Nhsx3ICnoD9BXY.bytes',
'1GvV8Js09wLnje4kSyN.asm',
'1GvV8Js09wLnje4kSyN.bytes',
'1gx83bLB4PSsYIKCT1Zt.asm',
'1gx83bLB4PSsYIKCT1Zt.bytes',
'1gYyAWCIEbGBJaMHL0k2.asm',
'1gYyAWCIEbGBJaMHL0k2.bytes',
'1h0FDrPmV7gfCzs9Id43.asm',
'1h0FDrPmV7gfCzs9Id43.bytes',
'1h7nkscUw51DurHXPE20.asm',
'1h7nkscUw51DurHXPE20.bytes',
'1hAQV0ebMLSgy4ktpfoC.asm',
'1hAQV0ebMLSgy4ktpfoC.bytes',
'1hcfvnSGMXEHU94egR3F.asm',
'1hcfvnSGMXEHU94egR3F.bytes',
'1her6tuVWjBkqZPUf8KR.asm',
'1her6tuVWjBkqZPUf8KR.bytes'.

```

'1hex2VAaHSC61Pv3ILOw.asm',
'1hex2VAaHSC61Pv3ILOw.bytes',
'1hfkVAzp9jybt7YKEBrF.asm',
'1hfkVAzp9jybt7YKEBrF.bytes',
'1HFt83wrBnSQu7cP2WZy.bytes',
'1HJhGL4zxVvcI2se07in.asm',
'1HJhGL4zxVvcI2se07in.bytes',
'1hK3zVxBGXmtTP6J1Hwb.asm',
'1hK3zVxBGXmtTP6J1Hwb.bytes',
'1HRWOCzEX65rYmMiZ21v.asm',
'1HRWOCzEX65rYmMiZ21v.bytes',
'1hToAGU8YgHQxfVONLqJ.asm',
'1hToAGU8YgHQxfVONLqJ.bytes',
'1hvs5Z9VJaFmwSDnx16u.asm',
'1hvs5Z9VJaFmwSDnx16u.bytes',
'1hzZA7JPylptLbxKYOmI.asm',
'1hzZA7JPylptLbxKYOmI.bytes',
'1i2SThtkCGEplHeoJxqI.asm',
'1i2SThtkCGEplHeoJxqI.bytes',
'1I3Am5YxgsXpDh6BLHTz.asm',
'1I3Am5YxgsXpDh6BLHTz.bytes',
'1I3EGeyOpWHNtYjLvnk5.asm',
'1I3EGeyOpWHNtYjLvnk5.bytes',
'1I8aJRE7T0igPZHWCvVQ.asm',
'1I8aJRE7T0igPZHWCvVQ.bytes',
'1IbqA5sdXj3P0CMOGmiE.asm',
'1IbqA5sdXj3P0CMOGmiE.bytes',
'1iFYGHfzdnCJRXA5uby9.asm',
'1iFYGHfzdnCJRXA5uby9.bytes',
'1IGur13q9bJQ7pmtMD8O.asm',
'1IGur13q9bJQ7pmtMD8O.bytes',
'1IjdZ2HuarlMYEBFbTCK.asm',
'1IjdZ2HuarlMYEBFbTCK.bytes',
'1ik3ncdBTYIJAfmjpwz.bytes',
'1iPIVbHF1w5Mcyn8D2Nu.asm',
'1iPIVbHF1w5Mcyn8D2Nu.bytes',
'1IpWLz6eyhVxDafQMKEd.asm',
'1IpWLz6eyhVxDafQMKEd.bytes',
'1Iqjvx7dTxCu54M0ylBE.asm',
'1Iqjvx7dTxCu54M0ylBE.bytes',
'1iTgmuEOkVaUnhC1z2JL.asm',
'1iTgmuEOkVaUnhC1z2JL.bytes',
'1iWpGIDVnJEaTHR8QSPL.asm',
'1iWpGIDVnJEaTHR8QSPL.bytes',
'1iWTJCz1SZgempU2xf04.asm',
'1iWTJCz1SZgempU2xf04.bytes',
'1IyCi4PMVRA9ULJ8Z1Gd.asm',
'1IyCi4PMVRA9ULJ8Z1Gd.bytes',
'1Jd0PTkEWtKcyCSbBD6H.asm',
'1Jd0PTkEWtKcyCSbBD6H.bytes',
'1jedHW3Pu0SGBqv4ytO8.asm',
...
]

```

In [3]:

```

source = 'train'
destination = 'byteFiles'
# we will check if the folder 'byteFiles' exists if it not there we will create a folder with the same name
if not os.path.isdir(destination):
    os.makedirs(destination)

if os.path.isdir(source):
    os.rename(source, 'asmFiles')
    source='asmFiles'
    data_files = os.listdir(source)
    for file in data_files:
        if (file.endswith("bytes")):
            shutil.move('asmFiles/'+file,destination)

```

3.1. Distribution of malware classes in whole data set

In [14]:

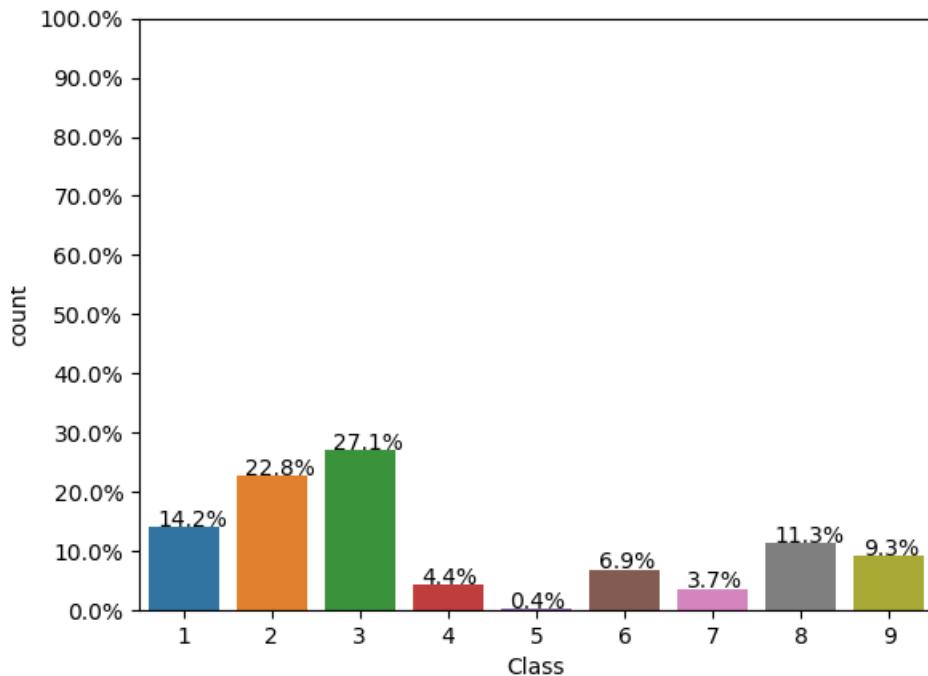
```

Y=pd.read_csv("trainLabels.csv")
total = len(Y)*1.
ax=sns.countplot(x="Class", data=Y)
for p in ax.patches:
    ax.annotate('{:.1f}%'.format(100*p.get_height()/total), (p.get_x()+0.1, p.get_height()+5))

#put 11 ticks (therefore 10 steps), from 0 to the total number of rows in the dataframe
ax.yaxis.set_ticks(np.linspace(0, total, 11))

#adjust the ticklabel to the desired format, without changing the position of the ticks.
ax.set_yticklabels(map('{:.1f}%'.format, 100*ax.yaxis.get_majorticklocs()/total))
plt.show()

```



3.2. Feature extraction

3.2.1 File size of byte files as a feature

In [21]:

```
Y=pd.read_csv("trainLabels.csv")
```

In [22]:

```

#file sizes of byte files

files=os.listdir('byteFiles')
filenames=Y['Id'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, st_nlink=1,
st_uid=0, st_gid=0,
    # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=1519638522)
    # read more about os.stat: here https://www.tutorialspoint.com/python/os_stat.htm
    statinfo=os.stat('byteFiles/'+file)
    # split the file name at '.' and take the first part of it i.e the file name
    file=file.split('.')[0]
    if any(file == filename for filename in filenames):

```

```

-- 
i=filenames.index(file)
class_bytes.append(class_y[i])
# converting into Mb's
sizebytes.append(statinfo.st_size/(1024.0*1024.0))
fnames.append(file)
data_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes})
print(data_size_byte.head())

```

	ID	size	Class
0	01azqd4InC7m9JpocGv5	5.012695	9
1	01IsoiSMh5gxyDYTl4CB	6.556152	2
2	01jsnpXSAlg6aPeDxrU	4.602051	9
3	01kCPWA9K2BoxQeS5Rju	0.679688	1
4	01SuzwMJEIXsK7A8dQb1	0.438965	8

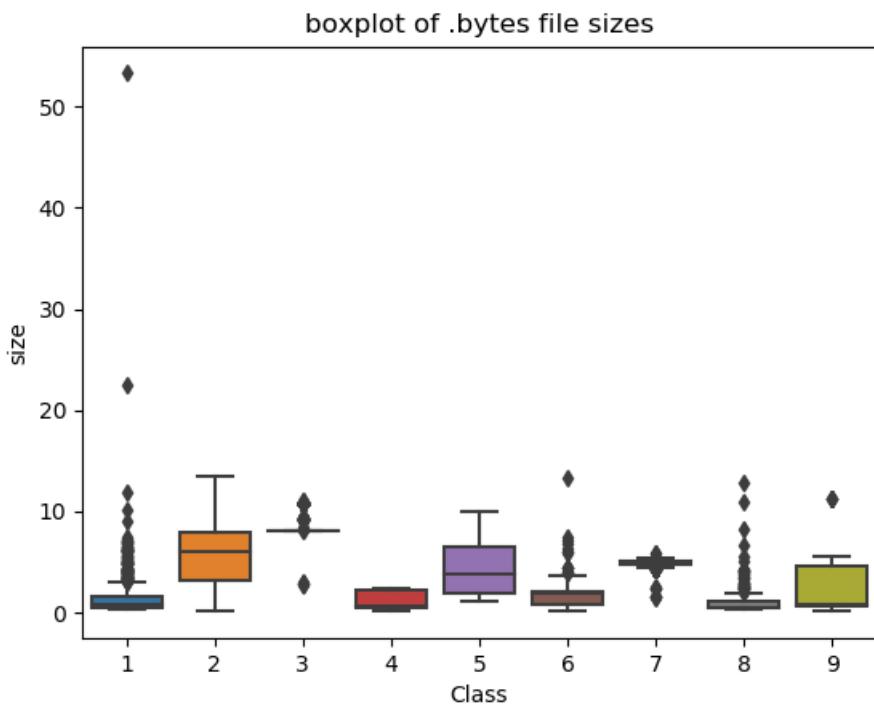
3.2.2 box plots of file size (.byte files) feature

In [16]:

```

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

```



3.2.3 feature extraction from byte files

In []:

```

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

files = os.listdir('byteFiles')
filenames=[]
array=[]
for file in files:
    if(file.endswith("bytes")):
        file=file.split('.')[0]

```

```

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

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

```

In [19]:

```

byte_features=pd.read_csv("result_with_size.csv")
print(byte_features.head())

```

	Unnamed: 0	ID	0	1	2	3	4	5	\			
0	0	01azqd4InC7m9JpocGv5	601905	3905	2816	3832	3345	3242				
1	1	01IsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844				
2	2	01jsnpXSAlg6aPeDxrU	93506	9542	2568	2438	8925	9330				
3	3	01kcPWA9K2B0xQeS5Rju	21091	1213	726	817	1257	625				
4	4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410				
	6	7	...	f9	fa	fb	fc	fd	fe	ff	??	\
0	3650	3201	...	3101	3211	3097	2758	3099	2759	5753	1824	
1	8420	7589	...	439	281	302	7639	518	17001	54902	8588	
2	9007	2342	...	2242	2885	2863	2471	2786	2680	49144	468	
3	550	523	...	485	462	516	1133	471	761	7998	13940	
4	262	249	...	350	209	239	653	221	242	2199	9008	
	size	Class										
0	4.234863	9										
1	5.538818	2										
2	3.887939	9										
3	0.574219	1										
4	0.370850	8										

[5 rows x 261 columns]

In [5]:

```

byte_features.columns

```

Out[5]:

```

Index(['Unnamed: 0', 'ID', '0', '1', '2', '3', '4', '5', '6', '7',
       ...
       'f9', 'fa', 'fb', 'fc', 'fd', 'fe', 'ff', '??', 'size', 'Class'],
      dtype='object', length=261)

```

In [23]:

```

result = pd.merge(byte_features, data_size_byte, on='ID', how='left')
result.head()

```

Out[23]:

	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	fb	fc	fd	fe	ff	'
0	0	01azqd4InC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	...	3097	2758	3099	2759	5753	18
1	1	01IsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	...	302	7639	518	17001	54902	85
2	2	01jsnpXSAlg6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	...	2863	2471	2786	2680	49144	4
3	3	01kcPWA9K2B0xQeS5Rju	21091	1213	726	817	1257	625	550	523	...	516	1133	471	761	7998	139

```
4 Unnamed:  
0 01SuzwMJEIXsK7A8dQbI 19760 710 302 433 554 415 262 240 ... 230 656 214 242 2198 90
```

5 rows × 263 columns

In [24]:

```
result = byte_features
```

In [25]:

```
# https://stackoverflow.com/a/29651514  
def normalize(df):  
    result1 = df.copy()  
    for feature_name in tqdm(df.columns):  
        if (str(feature_name) != str('ID') and str(feature_name) != str('Class')):  
            max_value = df[feature_name].max()  
            min_value = df[feature_name].min()  
            result1[feature_name] = (df[feature_name] - min_value) / (max_value - min_value)  
    return result1  
result = normalize(result)
```

100% |██████████| 261/261 [00:00<00:00, 298.87it/s]

In [26]:

```
data_y = result['Class']  
result.head()
```

Out[26]:

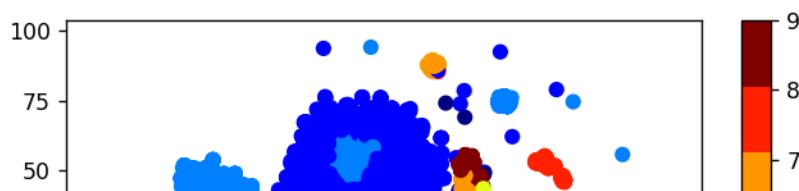
	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	f9
0	0.000000	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946	...	0.013560
1	0.000092	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984	...	0.001920
2	0.000184	01jsnpXSAlg6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155	...	0.009804
3	0.000276	01kcPWA9K2B0xQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481	...	0.002121
4	0.000368	01SuzwMJEIXsK7A8dQbI	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229	...	0.001530

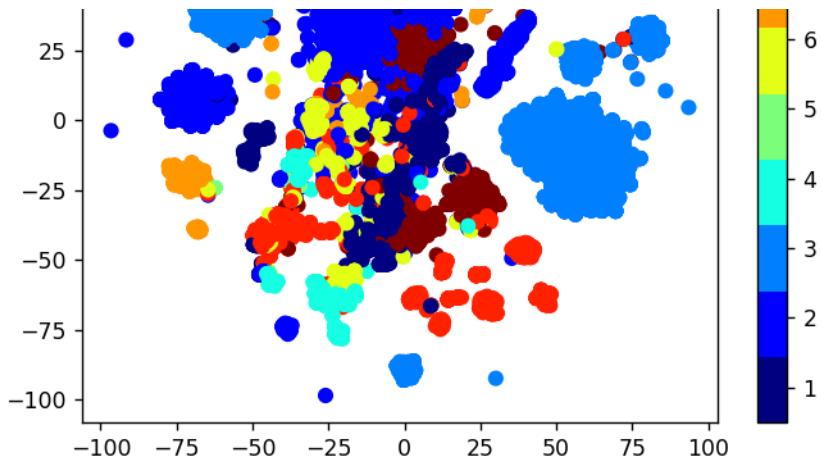
5 rows × 261 columns

3.2.4 Multivariate Analysis

In [0]:

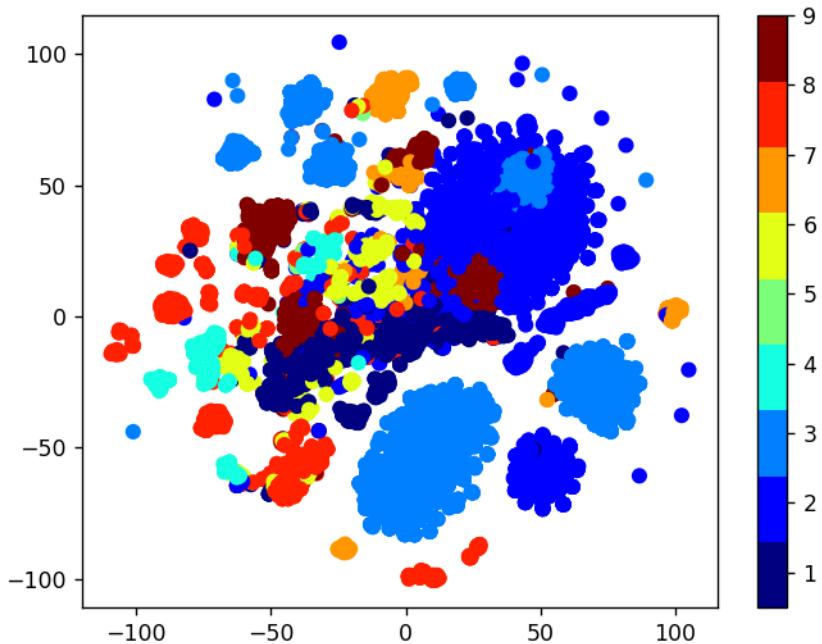
```
#multivariate analysis on byte files  
#this is with perplexity 50  
xtsne=TSNE(perplexity=50)  
results=xtsne.fit_transform(result.drop(['ID','Class'], axis=1))  
vis_x = results[:, 0]  
vis_y = results[:, 1]  
plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))  
plt.colorbar(ticks=range(10))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





In [0]:

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



Train Test split

In [29]:

```
#data_y = result['Class']
# split the data into test and train by maintaining same distribution of output variable 'y_true'
[stratify=y_true]
X_train, X_test, y_train, y_test = train_test_split(result.drop(['ID','Class'], axis=1), data_y,stratify=data_y,test_size=0.20)
# split the train data into train and cross validation by maintaining same distribution of output
variable 'y_train' [stratify=y_train]
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,stratify=y_train, test_size=0.20)
```

```
x_train, x_cv, y_train, y_cv = train_test_split(x_train, y_train, test_size=0.2)
```

In [30]:

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

```
Number of data points in train data: 6955
Number of data points in test data: 2174
Number of data points in cross validation data: 1739
```

In [0]:

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

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

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

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

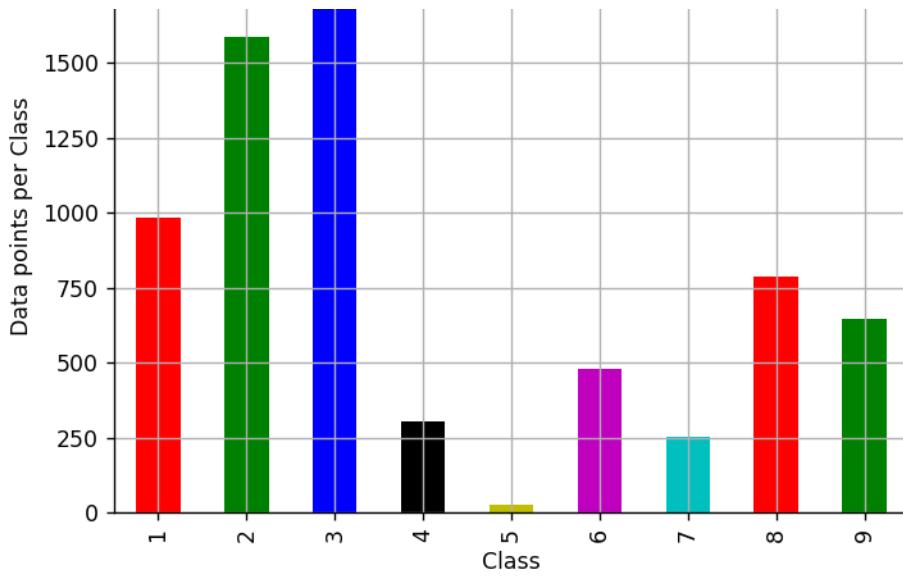
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(test_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

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

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

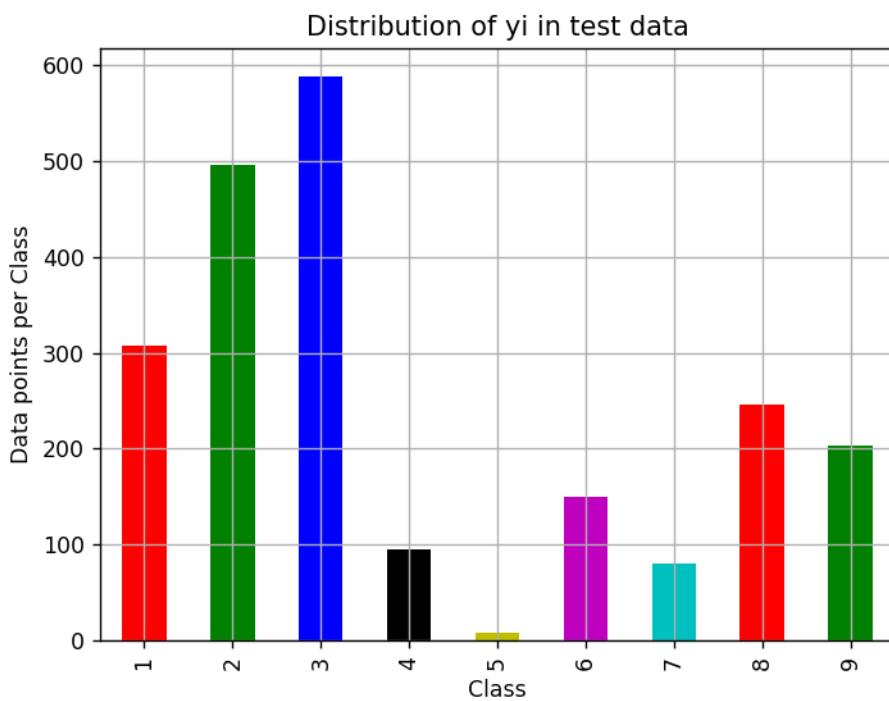
Distribution of yi in train data





```

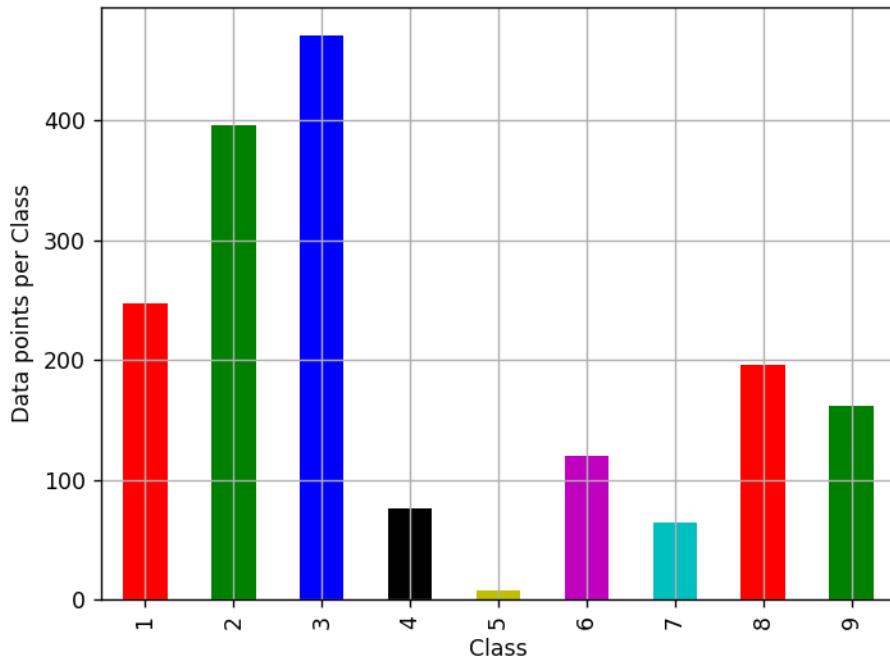
Number of data points in class 3 : 1883 ( 27.074 %)
Number of data points in class 2 : 1586 ( 22.804 %)
Number of data points in class 1 : 986 ( 14.177 %)
Number of data points in class 8 : 786 ( 11.301 %)
Number of data points in class 9 : 648 ( 9.317 %)
Number of data points in class 6 : 481 ( 6.916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



```

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 %)
-----
```

Distribution of yi in cross validation data



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

In [31]:

```
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    print("Number of misclassified points ",(len(test_y)-np.trace(C))/len(test_y)*100)
    # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j

    A =((C.T)/(C.sum(axis=1))).T
    #divid each element of the confusion matrix with the sum of elements in that column

    # C = [[1, 2],
    #       [3, 4]]
    # C.T = [[1, 3],
    #       [2, 4]]
    # C.sum(axis = 1) axis=0 corresonds to columns and axis=1 corresponds to rows in two
    #diamensional array
    # C.sum(axix =1) = [[3, 7]]
    # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
    #                           [2/3, 4/7]]

    # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
    #                               [3/7, 4/7]]
    # sum of row elements = 1

    B =(C/C.sum(axis=0))
    #divid each element of the confusion matrix with the sum of elements in that row
    # C = [[1, 2],
    #       [3, 4]]
    # C.sum(axis = 0) axis=0 corresonds to columns and axis=1 corresponds to rows in two
    #diamensional array
    # C.sum(axix =0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                      [3/4, 4/6]]

    labels = [1,2,3,4,5,6,7,8,9]
    cmap=sns.light_palette("green")
```

```

# representing A in heatmap format
print("-"*50, "Confusion matrix", "*"*50)
plt.figure(figsize=(10,5))
sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()

print("-"*50, "Precision matrix", "*"*50)
plt.figure(figsize=(10,5))
sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()
print("Sum of columns in precision matrix",B.sum(axis=0))

# representing B in heatmap format
print("-"*50, "Recall matrix" , "*"*50)
plt.figure(figsize=(10,5))
sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()
print("Sum of rows in precision matrix",A.sum(axis=1))

```

4. Machine Learning Models

4.1. Machine Learning Models on bytes files

4.1.1. Random Model

In [0]:

```

# we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to generate 9 numbers and divide each of the numbers by their sum
# ref: https://stackoverflow.com/a/18662466/4084039

test_data_len = X_test.shape[0]
cv_data_len = X_cv.shape[0]

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

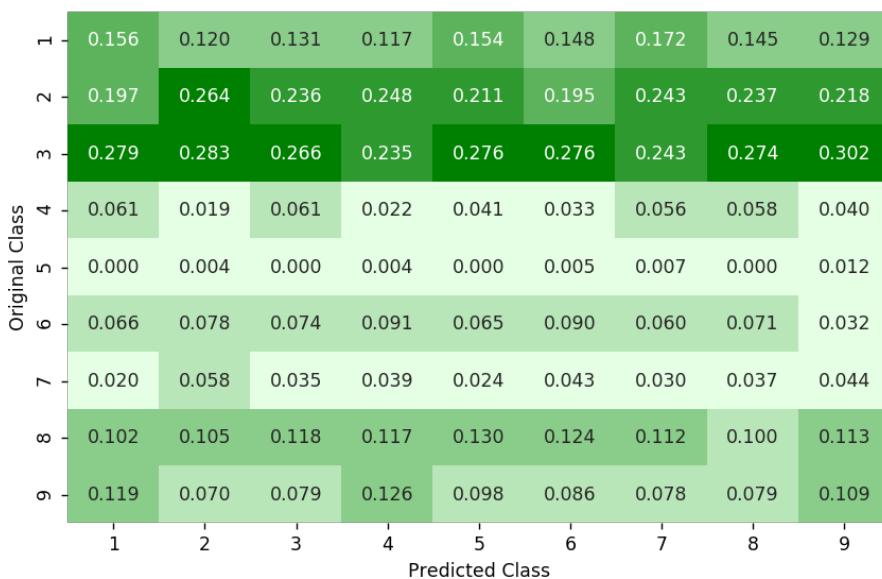
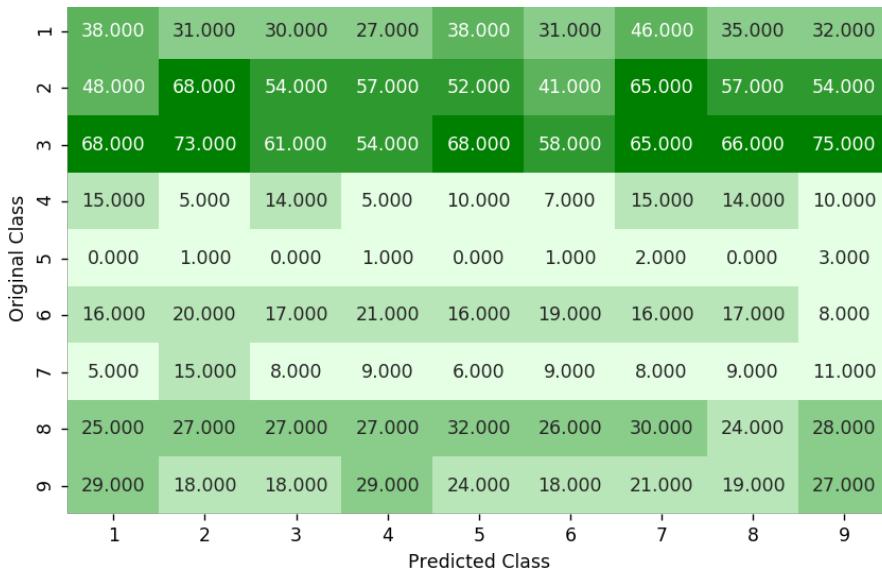
# Test-Set error.
#we create a output array that has exactly same as the test data
test_predicted_y = np.zeros((test_data_len,9))
for i in range(test_data_len):
    rand_probs = np.random.rand(1,9)
    test_predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random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

predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)

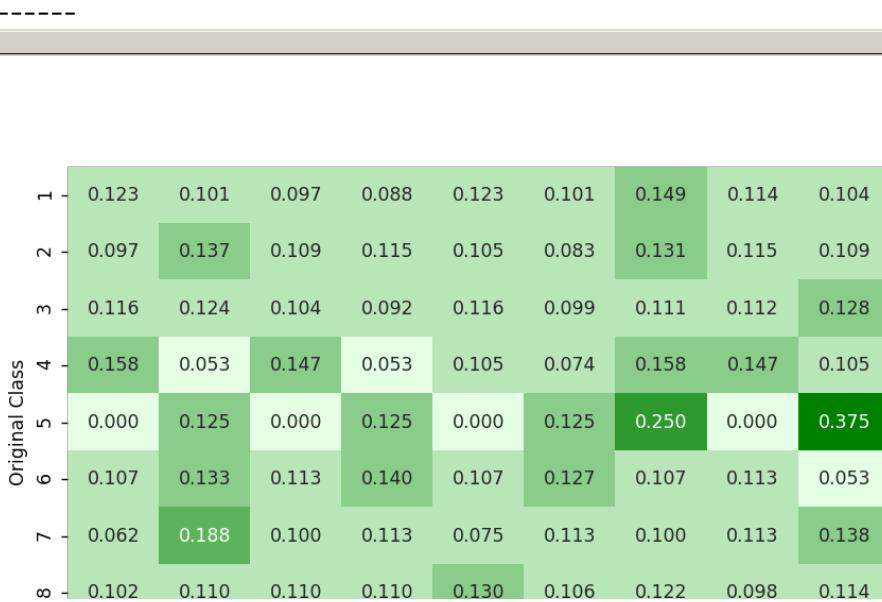
```

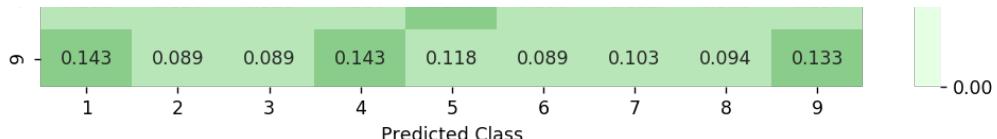
```

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



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





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

4.1.2. K Nearest Neighbour Classification

In [0]:

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

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


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


alpha = [x for x in range(1, 15, 2)]
cv_log_error_array=[]
for i in alpha:
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

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

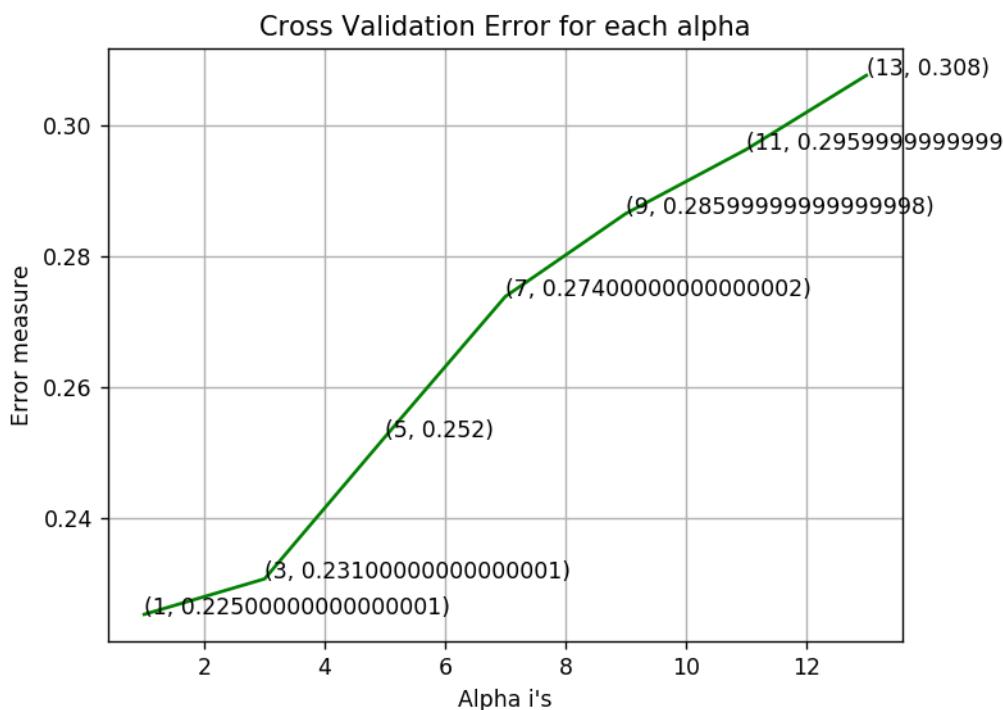
k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
k_cfl.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
sig_clf.fit(X_train, y_train)
```

```

predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
plot_confusion_matrix(y_test, sig_clf.predict(X_test))

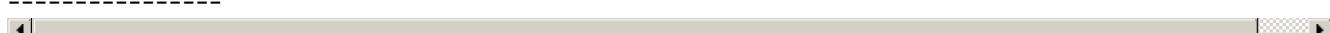
```

log_loss for k = 1 is 0.225386237304
 log_loss for k = 3 is 0.230795229168
 log_loss for k = 5 is 0.252421408646
 log_loss for k = 7 is 0.273827486888
 log_loss for k = 9 is 0.286469181555
 log_loss for k = 11 is 0.29623391147
 log_loss for k = 13 is 0.307551203154

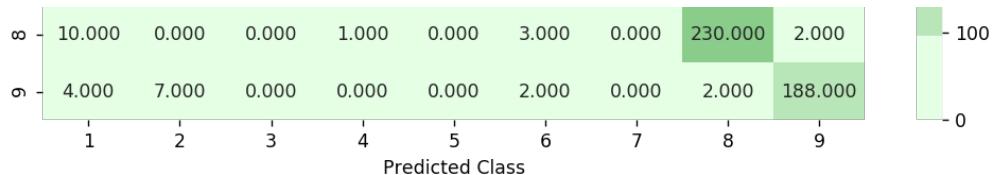


For values of best alpha = 1 The train log loss is: 0.0782947669247
 For values of best alpha = 1 The cross validation log loss is: 0.225386237304
 For values of best alpha = 1 The test log loss is: 0.241508604195
 Number of misclassified points 4.50781968721

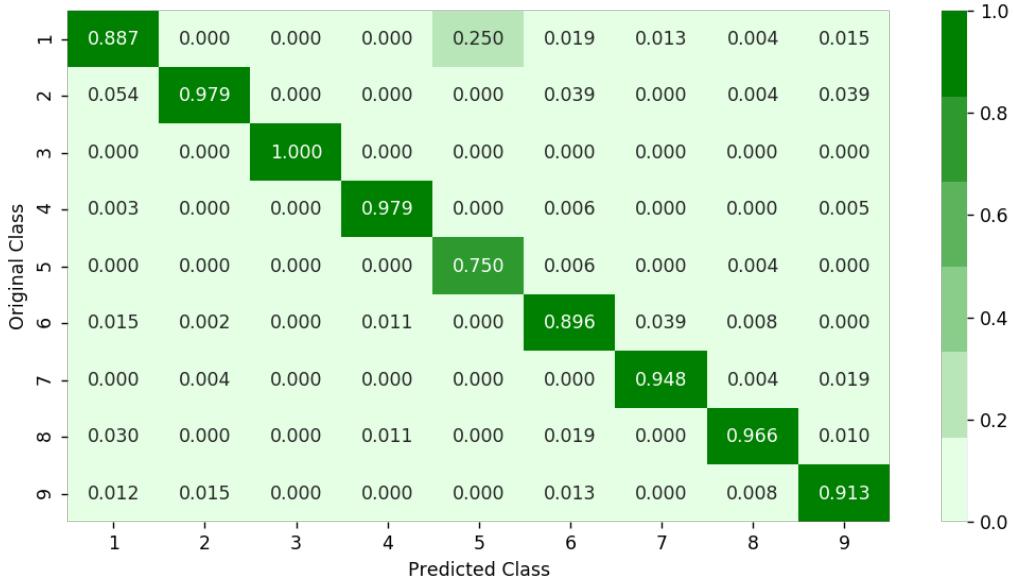
----- Confusion matrix -----



Original Class	1	2	3	4	5	6	7	8	9
1	298.000	0.000	0.000	0.000	2.000	3.000	1.000	1.000	3.000
2	18.000	463.000	0.000	0.000	0.000	6.000	0.000	1.000	8.000
3	0.000	0.000	588.000	0.000	0.000	0.000	0.000	0.000	0.000
4	1.000	0.000	0.000	92.000	0.000	1.000	0.000	0.000	1.000
5	0.000	0.000	0.000	0.000	6.000	1.000	0.000	1.000	0.000
6	5.000	1.000	0.000	1.000	0.000	138.000	3.000	2.000	0.000
7	0.000	2.000	0.000	0.000	0.000	0.000	73.000	1.000	4.000

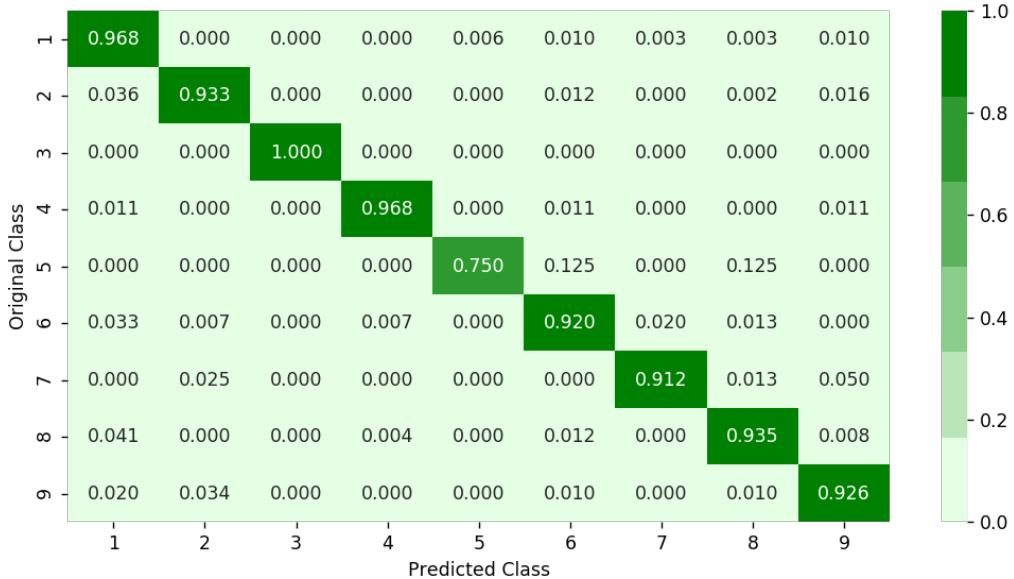


Precision matrix



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

Recall matrix



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

4.1.3. Logistic Regression

To 101:

```

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

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

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

alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=logisticR.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

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

logisticR=LogisticRegression(penalty='l2',C=alpha[best_alpha],class_weight='balanced')
logisticR.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
sig_clf.fit(X_train, y_train)
pred_y=sig_clf.predict(X_test)

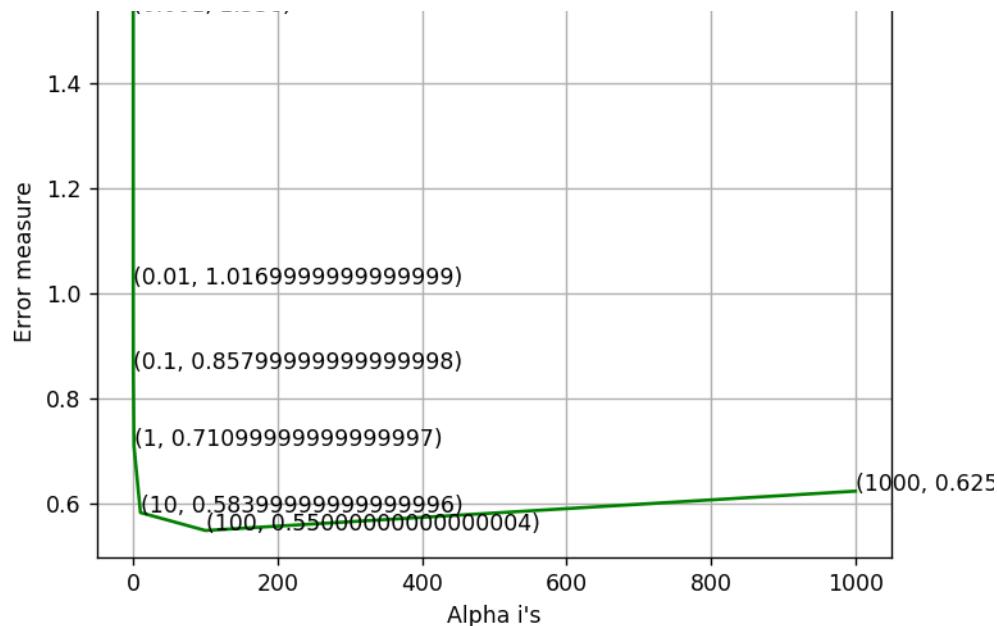
predict_y = sig_clf.predict_proba(X_train)
print ('log loss for train data',log_loss(y_train, predict_y, labels=logisticR.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_cv)
print ('log loss for cv data',log_loss(y_cv, predict_y, labels=logisticR.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print ('log loss for test data',log_loss(y_test, predict_y, labels=logisticR.classes_, eps=1e-15))
plot_confusion_matrix(y_test, sig_clf.predict(X_test))

log_loss for c =  1e-05 is 1.56916911178
log_loss for c =  0.0001 is 1.57336384417
log_loss for c =  0.001 is 1.53598598273
log_loss for c =  0.01 is 1.01720972418
log_loss for c =  0.1 is 0.857766083873
log_loss for c =  1 is 0.711154393309
log_loss for c =  10 is 0.583929522635
log_loss for c =  100 is 0.549929846589
log_loss for c =  1000 is 0.624746769121

```

Cross Validation Error for each alpha





```
log loss for train data 0.498923428696
log loss for cv data 0.549929846589
log loss for test data 0.528347316704
Number of misclassified points 12.3275068997
```

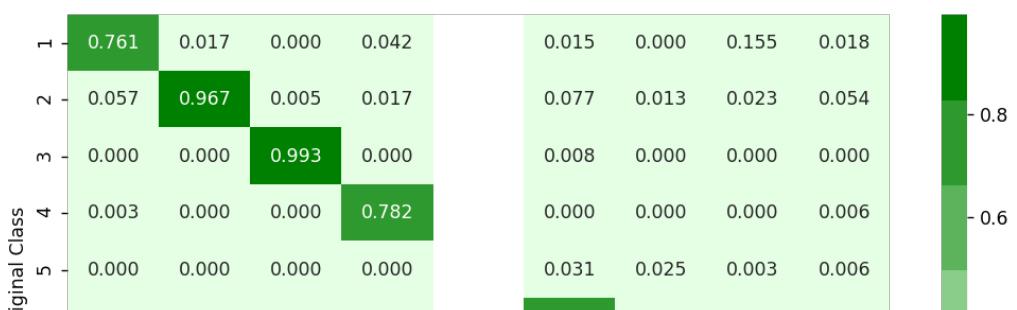
----- Confusion matrix -----

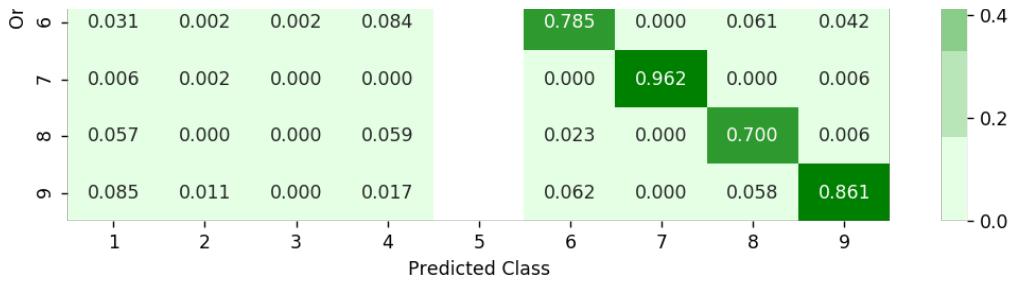
[◀] [▶]



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

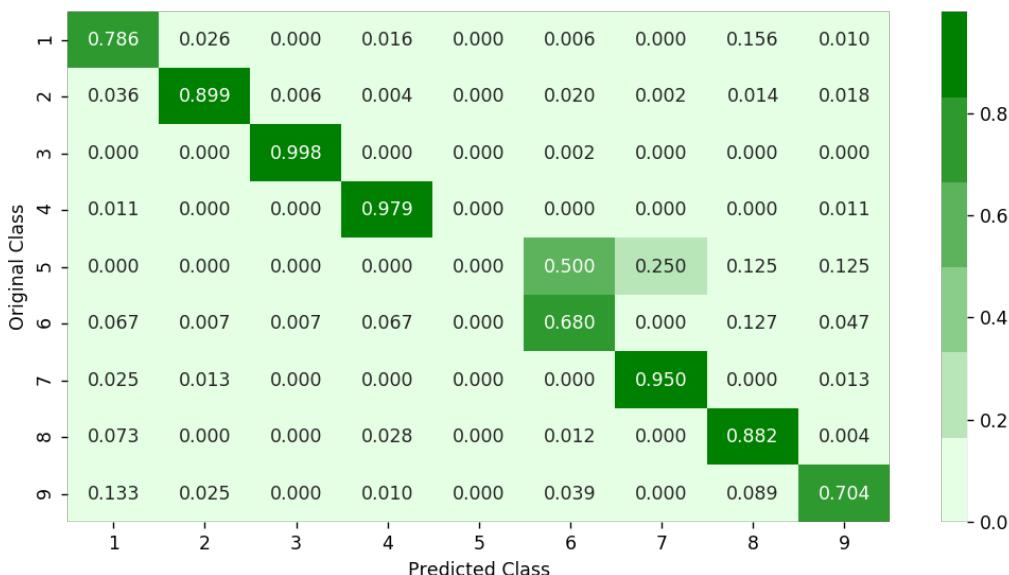
[◀] [▶]





```
Sum of columns in precision matrix [ 1. 1. 1. 1. nan 1. 1. 1. 1.]
```

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



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

4.1.4. Random Forest Classifier

In [0]:

```
# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None,
# verbose=0, warm_start=False,
# class_weight=None)

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

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

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

alpha=[10,50,100,500,1000,2000,3000]
```

```

cv_log_error_array=[]
train_log_error_array=[]
from sklearn.ensemble import RandomForestClassifier
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

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

r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
r_cfl.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train, y_train)

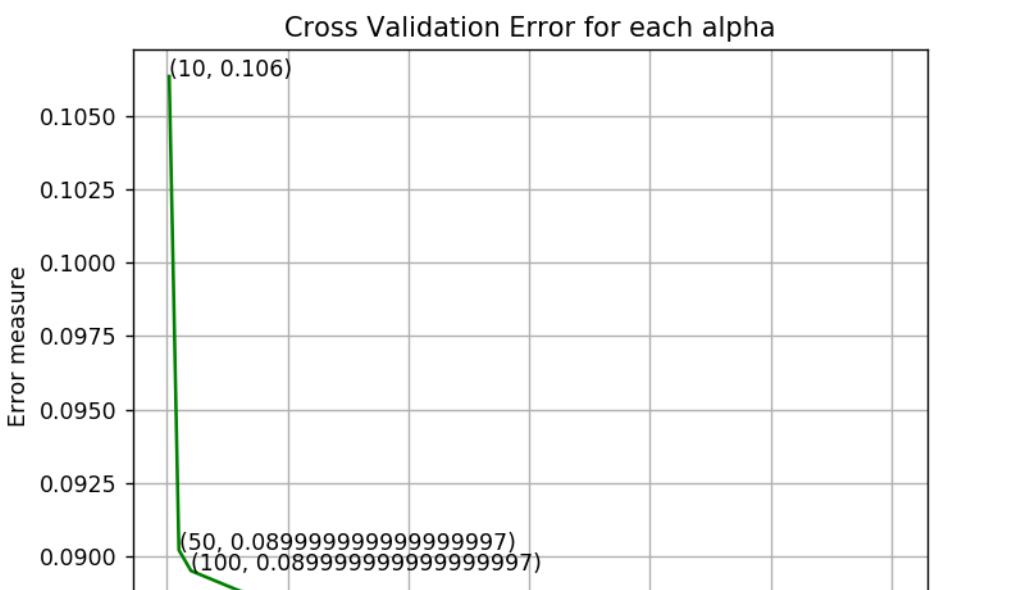
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,predict_y))
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
plot_confusion_matrix(y_test, sig_clf.predict(X_test))

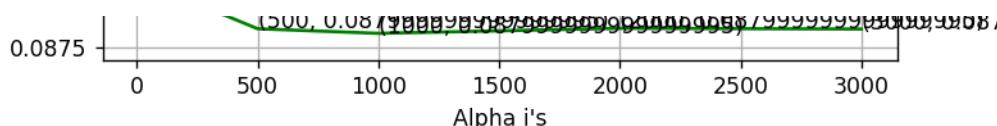
```

```

log_loss for c = 10 is 0.106357709164
log_loss for c = 50 is 0.0902124124145
log_loss for c = 100 is 0.0895043339776
log_loss for c = 500 is 0.0881420869288
log_loss for c = 1000 is 0.0879849524621
log_loss for c = 2000 is 0.0881566647295
log_loss for c = 3000 is 0.0881318948443

```





For values of best alpha = 1000 The train log loss is: 0.0266476291801
 For values of best alpha = 1000 The cross validation log loss is: 0.0879849524621
 For values of best alpha = 1000 The test log loss is: 0.0858346961407
 Number of misclassified points 2.02391904324

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



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



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

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





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

4.1.5. XgBoost Classification

In [0]:

```
# Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBClassifier function here
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# default parameters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
min_child_weight=1,
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
reg_lambda=1,
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# -----


alpha=[10,50,100,500,1000,2000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
    x_cfl.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

fig, ax = plt.subplots()
```

```

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

x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
x_cfl.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train, y_train)

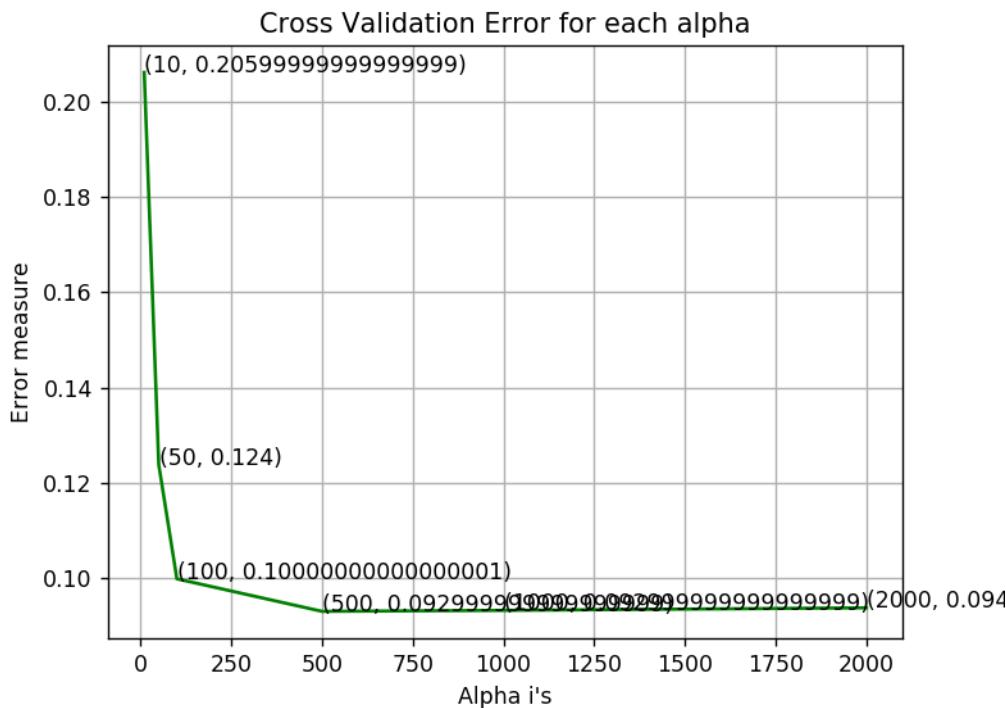
predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
plot_confusion_matrix(y_test, sig_clf.predict(X_test))

```

```

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

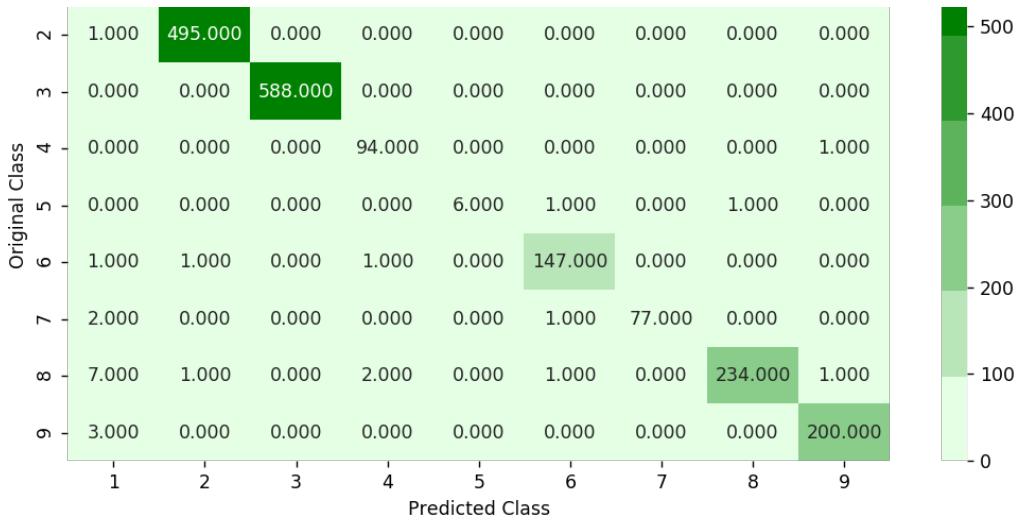
```



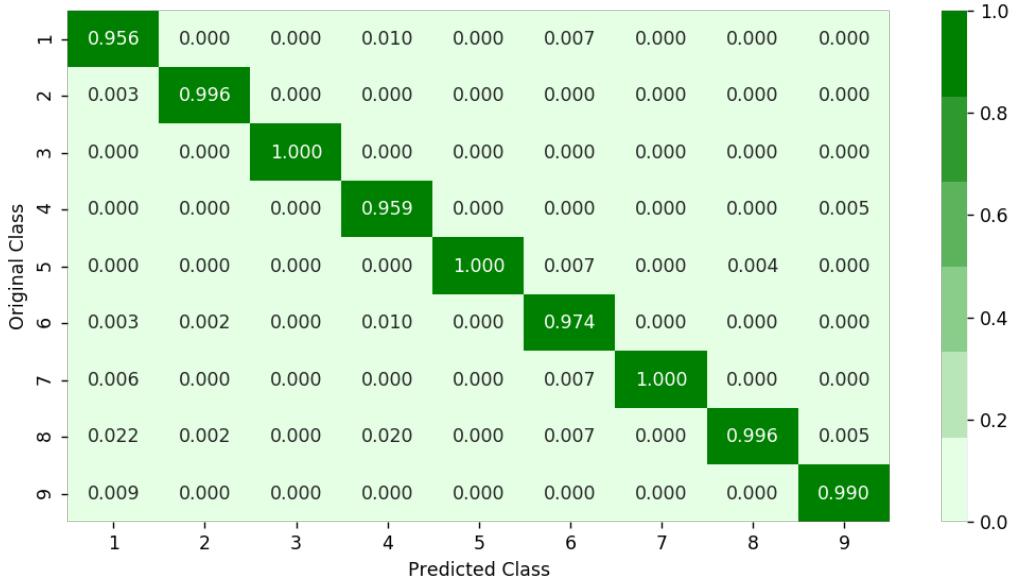
```

For values of best alpha = 500 The train log loss is: 0.0225231805824
For values of best alpha = 500 The cross validation log loss is: 0.0931035681289
For values of best alpha = 500 The test log loss is: 0.0792067651731
Number of misclassified points 1.24195032199
----- Confusion matrix -----
-----
```

1 - 306.000	0.000	0.000	1.000	0.000	1.000	0.000	0.000	0.000
-------------	-------	-------	-------	-------	-------	-------	-------	-------

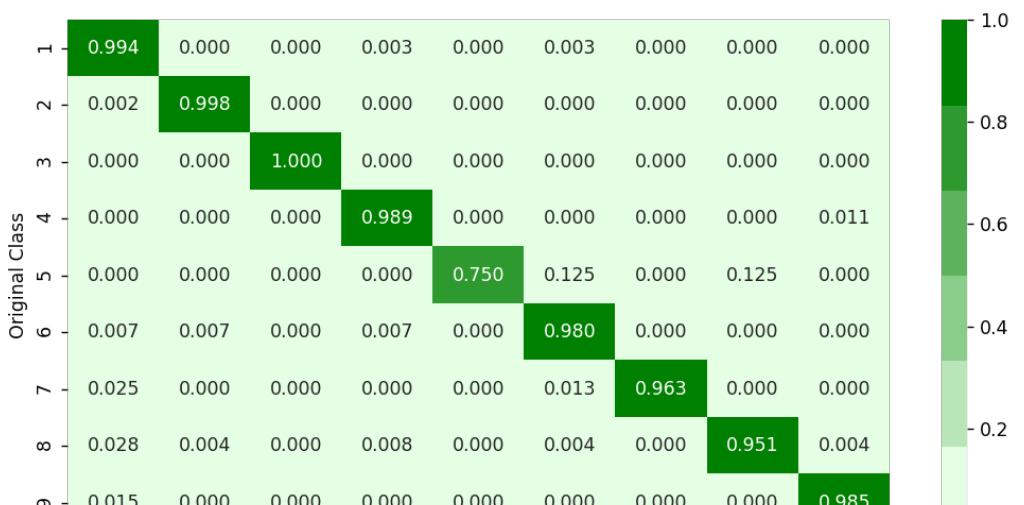


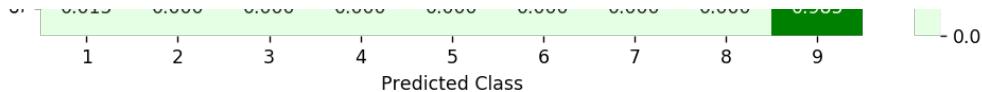
Precision matrix



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

Recall matrix





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

4.1.5. XgBoost Classification with best hyper parameters using RandomSearch

In [0]:

```
# https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/
x_cfl=XGBClassifier()

prams={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
    'n_estimators':[100,200,500,1000,2000],
    'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jobs=-1,)
random_cfl1.fit(X_train,y_train)
```

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

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

Out[0]:

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

In [0]:

```
print (random_cfl1.best_params_)
```

```
{'subsample': 1, 'n_estimators': 500, 'max_depth': 5, 'learning_rate': 0.05, 'colsample_bytree': 0
.5}
```

In [0]:

```
# Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBClassifier function here
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# default parameters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
# # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
# min_child_weight=1,
# # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
# reg_lambda=1,
# # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
```

```

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# -----


x_cfl=XGBClassifier(n_estimators=2000, learning_rate=0.05, colsample_bytree=1, max_depth=3)
x_cfl.fit(X_train,y_train)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train,y_train)

predict_y = c_cfl.predict_proba(X_train)
print ('train loss',log_loss(y_train, predict_y))
predict_y = c_cfl.predict_proba(X_cv)
print ('cv loss',log_loss(y_cv, predict_y))
predict_y = c_cfl.predict_proba(X_test)
print ('test loss',log_loss(y_test, predict_y))

train loss 0.022540976086
cv loss 0.0928710624158
test loss 0.0782688587098

```

4.2 Modeling with .asm files

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

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

We read the top solutions and handpicked the features from those papers/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 [0]:

```

#intially create five folders
#first
#second
#third
#fourth
#fifth
#this code tells us about random split of files into five folders
folder_1 ='first'
folder_2 ='second'

```

```

folder_1 = 'second'
folder_2 = 'third'
folder_3 = 'fourth'
folder_4 = 'fifth'
folder_5 = 'output'
for i in [folder_1, folder_2, folder_3, folder_4, folder_5, folder_6]:
    if not os.path.isdir(i):
        os.makedirs(i)

source='train/'
files = os.listdir('train')
ID=df['Id'].tolist()
data=range(0,10868)
r.shuffle(data)
count=0
for i in range(0,10868):
    if i % 5==0:
        shutil.move(source+files[data[i]],'first')
    elif i%5==1:
        shutil.move(source+files[data[i]],'second')
    elif i%5 ==2:
        shutil.move(source+files[data[i]],'third')
    elif i%5 ==3:
        shutil.move(source+files[data[i]],'fourth')
    elif i%5==4:
        shutil.move(source+files[data[i]],'fifth')

```

In [0]:

```

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

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

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

    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec',
'add','imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb','jz','rtn','lea','movzx']
    #best keywords that are taken from different blogs
    keywords = ['.dll','std::',':dword']
    #Below taken registers are general purpose registers and special registers
    #All the registers which are taken are best
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\asmssmallfile.txt","w+")
    files = os.listdir('first')
    for f in files:
        #filling the values with zeros into the arrays
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        # https://docs.python.org/3/library/codecs.html#codecs.ignore_errors
        # https://docs.python.org/3/library/codecs.html#codecs.Codec.encode
        with codecs.open('first/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                # https://www.tutorialspoint.com/python3/string_rstrip.htm
                line=lines.rstrip().split()
                l=line[0]
                #counting the prefixes in each and every line
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                #counting the opcodes in each and every line
                for i in range(len(opcodes)):
                    if opcodes[i] in line[0]:

```

```

        if any(opcodes[i]==li for li in line):
            features.append(opcodes[i])
            opcodescount[i]+=1
    #counting registers in the line
    for i in range(len(registers)):
        for li in line:
            # we will use registers only in 'text' and 'CODE' segments
            if registers[i] in li and ('text' in l or 'CODE' in l):
                registerscount[i]+=1
    #counting keywords in the line
    for i in range(len(keywords)):
        for li in line:
            if keywords[i] in li:
                keywordcount[i]+=1
#pushing the values into the file after reading whole file
for prefix in prefixescount:
    file1.write(str(prefix)+",")
for opcode in opcodescount:
    file1.write(str(opcode)+",")
for register in registerscount:
    file1.write(str(register)+",")
for key in keywordcount:
    file1.write(str(key)+",")
file1.write("\n")
file1.close()

#same as above
def secondprocess():
    prefixes = ['HEADER','.text','.Pav','.idata','.data','.bss','.rdata','.edata',
'.rsrc','.tls','.reloc','.BSS','.CODE']
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec',
'add','imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb','jz','rtn','lea','movz
x']
    keywords = ['.dll','std::',':dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\mediumasmfile.txt","w+")
    files = os.listdir('second')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('second/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                for i in range(len(registers)):
                    for li in line:
                        if registers[i] in li and ('text' in l or 'CODE' in l):
                            registerscount[i]+=1
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
            for prefix in prefixescount:
                file1.write(str(prefix)+",")
            for opcode in opcodescount:
                file1.write(str(opcode)+",")
            for register in registerscount:
                file1.write(str(register)+",")
            for key in keywordcount:
                file1.write(str(key)+",")
            file1.write("\n")
        file1.close()

```

```

# same as smallprocess() functions
def thirdprocess():
    prefixes = ['HEADER','.text','.Pav','.idata','.data','.bss','.rdata','.edata','.rsrc',
'.tls','.reloc','.BSS','.CODE']
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec',
'add','imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb','jz','rtn','lea','movz
x']
    keywords = ['.dll','std::':':dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\largeasmfile.txt","w+")
    files = os.listdir('thrid')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('thrid/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                for i in range(len(registers)):
                    for li in line:
                        if registers[i] in li and ('text' in l or 'CODE' in l):
                            registerscount[i]+=1
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
                for prefix in prefixescount:
                    file1.write(str(prefix)+",")
                for opcode in opcodescount:
                    file1.write(str(opcode)+",")
                for register in registerscount:
                    file1.write(str(register)+",")
                for key in keywordcount:
                    file1.write(str(key)+",")
        file1.write("\n")
    file1.close()

def fourthprocess():
    prefixes = ['HEADER','.text','.Pav','.idata','.data','.bss','.rdata','.edata','.rsrc',
'.tls','.reloc','.BSS','.CODE']
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec',
'add','imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb','jz','rtn','lea','movz
x']
    keywords = ['.dll','std::':':dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\hugeasmfile.txt","w+")
    files = os.listdir('fourth/')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('fourth/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1

```

```

        if prefixes[i] in line[0]:
            prefixescount[i]+=1
    line=line[1:]
    for i in range(len(opcodes)):
        if any(opcodes[i]==li for li in line):
            features.append(opcodes[i])
            opcodescount[i]+=1
    for i in range(len(registers)):
        for li in line:
            if registers[i] in li and ('text' in li or 'CODE' in li):
                registerscount[i]+=1
    for i in range(len(keywords)):
        for li in line:
            if keywords[i] in li:
                keywordcount[i]+=1
    for prefix in prefixescount:
        file1.write(str(prefix)+",")
    for opcode in opcodescount:
        file1.write(str(opcode)+",")
    for register in registerscount:
        file1.write(str(register)+",")
    for key in keywordcount:
        file1.write(str(key)+",")
    file1.write("\n")
file1.close()

def fifthprocess():
    prefixes = ['HEADER','.text','.Pav','.idata','.data','.bss','.rdata','.edata','.rsrc',
'.tls','.reloc','.BSS','.CODE']
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec',
'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb','jz','rtn','lea','movzx']
    keywords = ['.dll','std::':dword']
    registers=[edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
    file1=open("output\trainasmfile.txt","w+")
    files = os.listdir('fifth/')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('fifth/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                for i in range(len(registers)):
                    for li in line:
                        if registers[i] in li and ('text' in li or 'CODE' in li):
                            registerscount[i]+=1
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
        for prefix in prefixescount:
            file1.write(str(prefix)+",")
        for opcode in opcodescount:
            file1.write(str(opcode)+",")
        for register in registerscount:
            file1.write(str(register)+",")
        for key in keywordcount:
            file1.write(str(key)+",")
        file1.write("\n")
file1.close()

```

```

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

if __name__=="__main__":
    main()

```

In [10]:

```

# asmoutputfile.csv(output generated from the above two cells) will contain all the extracted features from .asm files
# this file will be uploaded in the drive, you can directly use this
dfasm=pd.read_csv("asmoutputfile.csv")
Y.columns = ['ID', 'Class']
result_asm = pd.merge(dfasm, Y,on='ID', how='left')
result_asm.head()

```

Out[10]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	edx	esi	eax	ebx	ecx	edi	eb
0	01kcPWA9K2B0xQeS5Rju		19	744	0	127	57	0	323	0	3	...	18	66	15	43	83	0 1
1	1E93CpP60RHFNiT5Qfvn		17	838	0	103	49	0	0	0	3	...	18	29	48	82	12	0 1
2	3ekVow2ajZHbTnBcsDfX		17	427	0	50	43	0	145	0	3	...	13	42	10	67	14	0 1
3	3X2nY7iQaPBIVDrAZqJe		17	227	0	43	19	0	0	0	3	...	6	8	14	7	2	0
4	46OZzdsSKDCFV8h7XWxf		17	402	0	59	170	0	0	0	3	...	12	9	18	29	5	0 1

5 rows × 53 columns

4.2.1.1 Files sizes of each .asm file

In [11]:

```

#file sizes of byte files

files=os.listdir('asmFiles')
filenames=Y['ID'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, st_nlink=1,
st_uid=0, st_gid=0,
    # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=1519638522)
    # read more about os.stat: here https://www.tutorialspoint.com/python/os_stat.htm
    statinfo=os.stat('asmFiles/'+file)
    # split the file name at '.' and take the first part of it i.e the file name
    file=file.split('.')[0]
    if any(file == filename for filename in filenames):

```

```

i=filenames.index(file)
class_bytes.append(class_y[i])
# converting into Mb's
sizebytes.append(statinfo.st_size/(1024.0*1024.0))
fnames.append(file)
asm_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes})
print(asm_size_byte.head())

```

	ID	size	Class
0	01azqd4InC7m9JpocGv5	56.229886	9
1	01IsoiSMh5gxyDYTl4CB	13.999378	2
2	01jsnpXSAlgw6aPeDxrU	8.507785	9
3	01kcPWA9K2B0xQeS5Rju	0.078190	1
4	01SuzwMJEIXsK7A8dQbl	0.996723	8

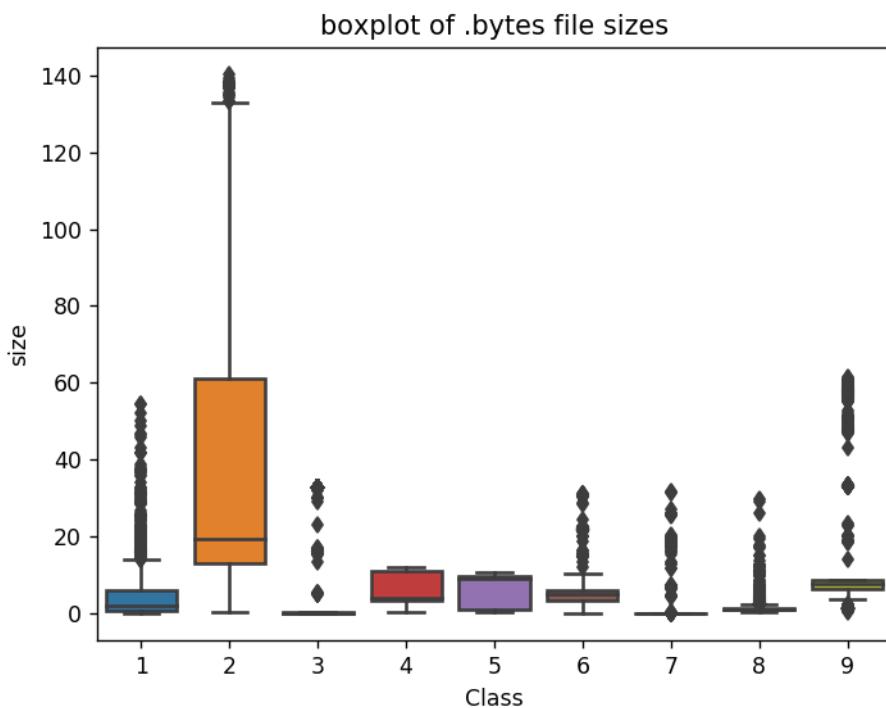
4.2.1.2 Distribution of .asm file sizes

In [0]:

```

#boxplot of asm files
ax = sns.boxplot(x="Class", y="size", data=asm_size_byte)
plt.title("boxplot of .bytes file sizes")
plt.show()

```



In [12]:

```

# add the file size feature to previous extracted features
print(result_asm.shape)
print(asm_size_byte.shape)
result_asm = pd.merge(result_asm, asm_size_byte.drop(['Class'], axis=1), on='ID', how='left')
result_asm.head()

```

(10868, 53)
(10868, 3)

Out[12]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	esi	eax	ebx	ecx	edi	ebp	es
0	01kcPWA9K2B0xQeS5Rju	19	744	0	127	57	0	323	0	3	...	66	15	43	83	0	17	4
1	1E93CpP60RHFNiT5Qfvn	17	838	0	103	49	0	0	0	3	...	29	48	82	12	0	14	

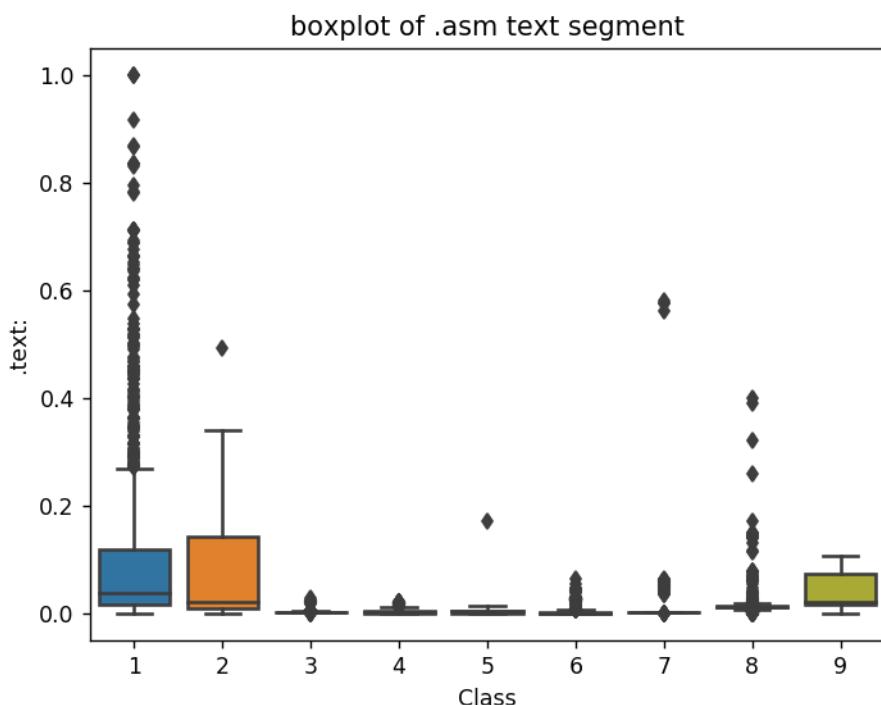
	ID	HEADER	.text	.Pav	.idata	.data	.bss	.rdata	.edata	.rsrc	...	esi	eax	ebx	ecx	edi	ebp	es
2	3ekVow2ajZHbTnBcsDfx	17	427	0	43	19	0	0	0	0	3	...	8	14	7	2	0	8
3	3X2nY7iQaPBIDrAZqJe	17	227	0	43	19	0	0	0	0	3	...	8	14	7	2	0	8
4	46OZzdsSKDCFV8h7XWxf	17	402	0	59	170	0	0	0	0	3	...	9	18	29	5	0	11

5 rows × 54 columns

4.2.2 Univariate analysis on asm file features

In [0]:

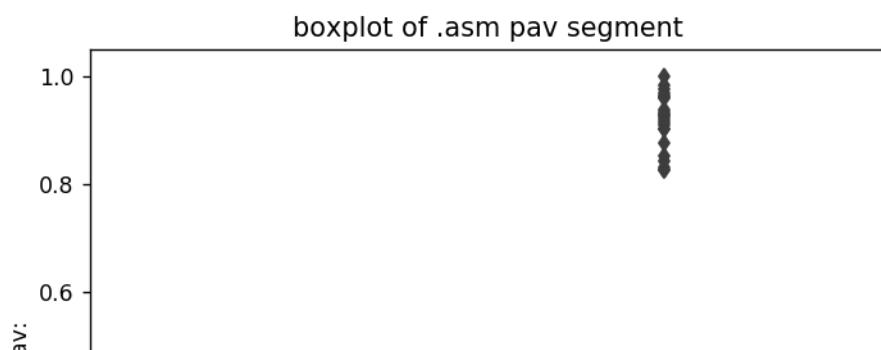
```
ax = sns.boxplot(x="Class", y=".text:", data=result_asm)
plt.title("boxplot of .asm text segment")
plt.show()
```

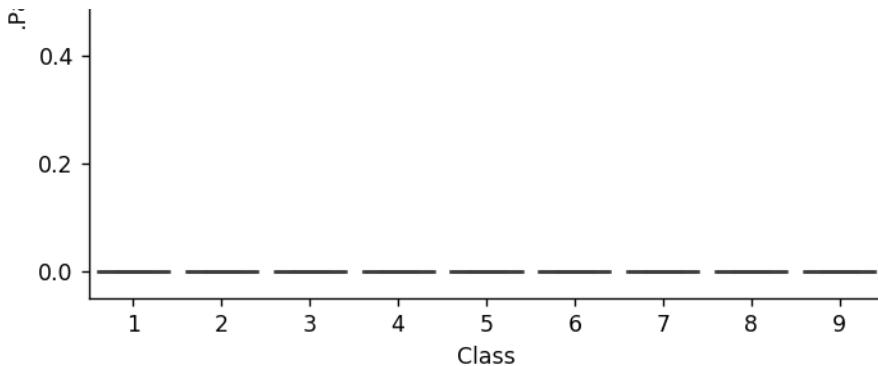


The plot is between Text and class
Class 1,2 and 9 can be easily separated

In [0]:

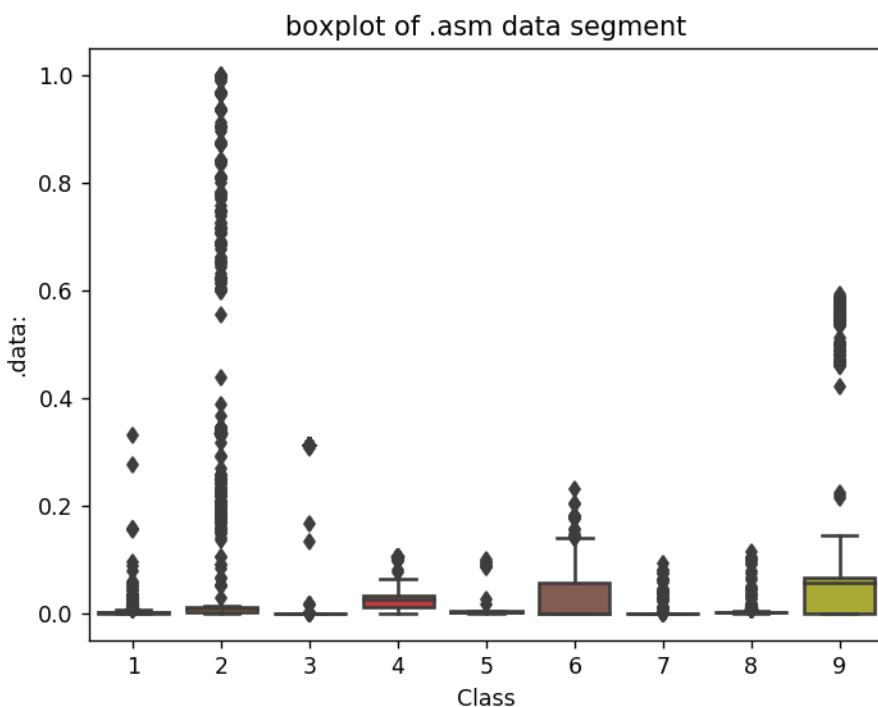
```
ax = sns.boxplot(x="Class", y=".Pav:", data=result_asm)
plt.title("boxplot of .asm pav segment")
plt.show()
```





In [0]:

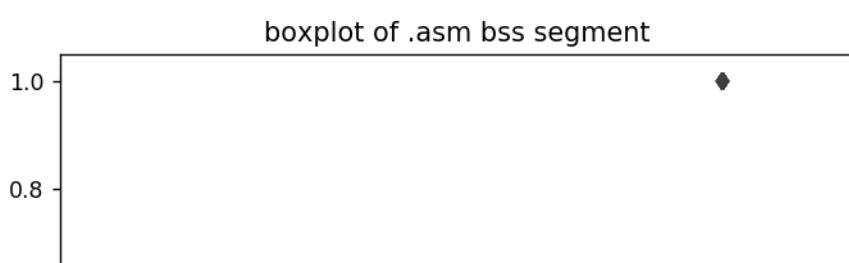
```
ax = sns.boxplot(x="Class", y=".data:", data=result_asm)
plt.title("boxplot of .asm data segment")
plt.show()
```

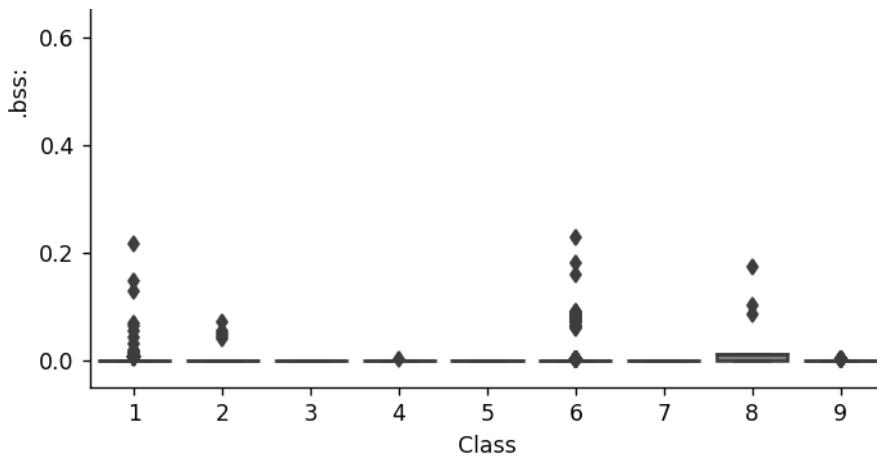


The plot is between data segment and class label
class 6 and class 9 can be easily separated from given points

In [0]:

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

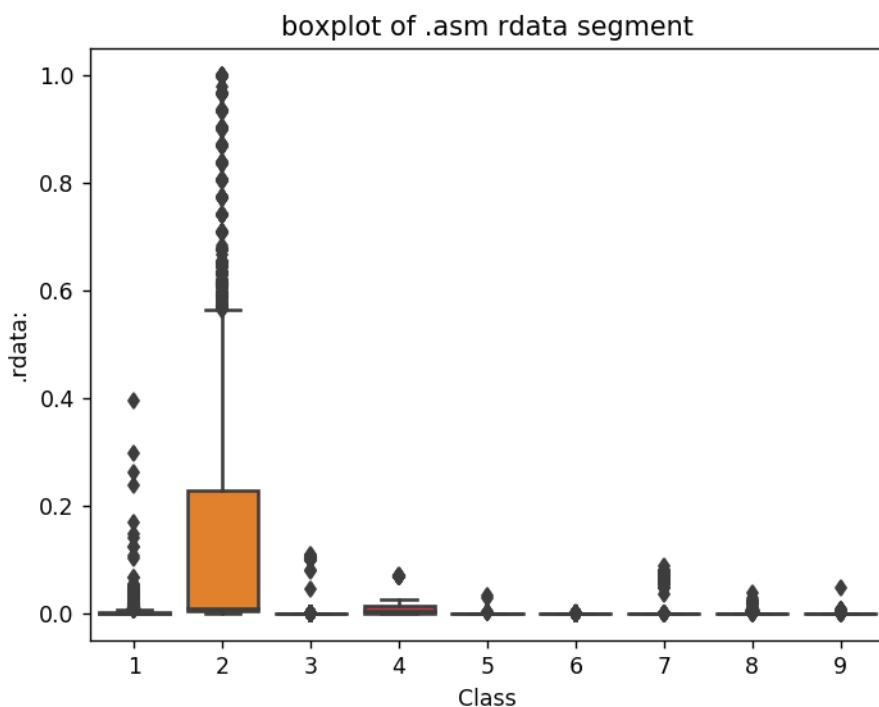




```
plot between bss segment and class label
very less number of files are having bss segment
```

In [0]:

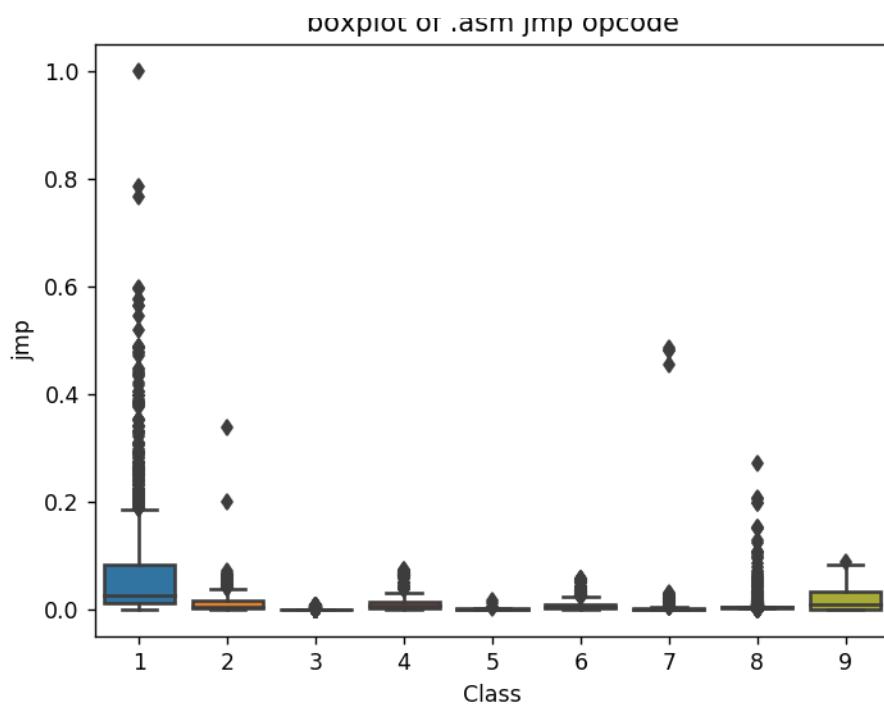
```
ax = sns.boxplot(x="Class", y=".rdata:", data=result_asm)
plt.title("boxplot of .asm rdata segment")
plt.show()
```



```
Plot between rdata segment and Class segment
Class 2 can be easily separated 75 percentile files are having 1M rdata lines
```

In [0]:

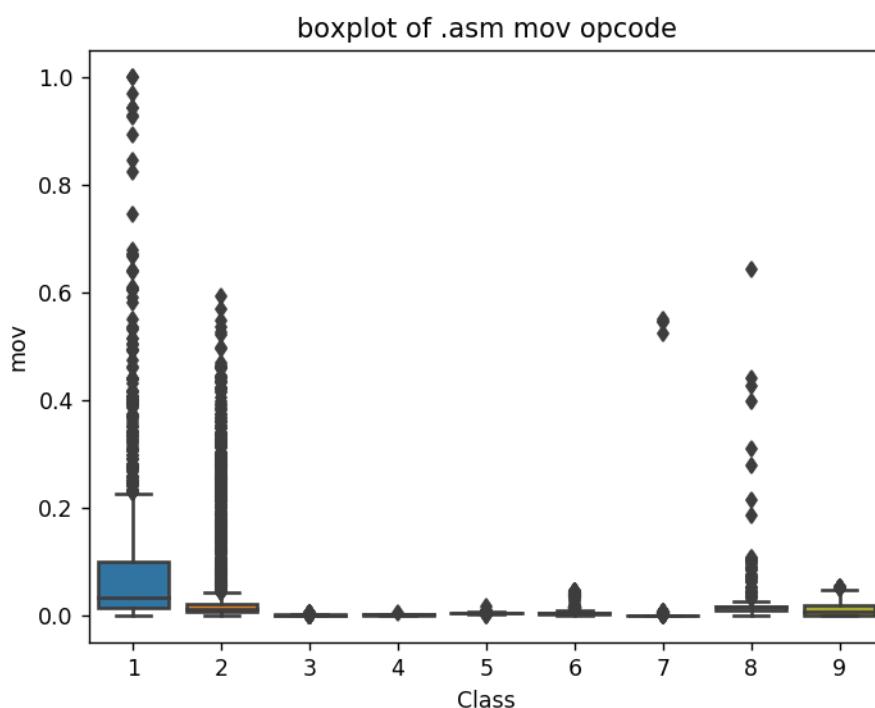
```
ax = sns.boxplot(x="Class", y="jmp", data=result_asm)
plt.title("boxplot of .asm jmp opcode")
plt.show()
```



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

In [0]:

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

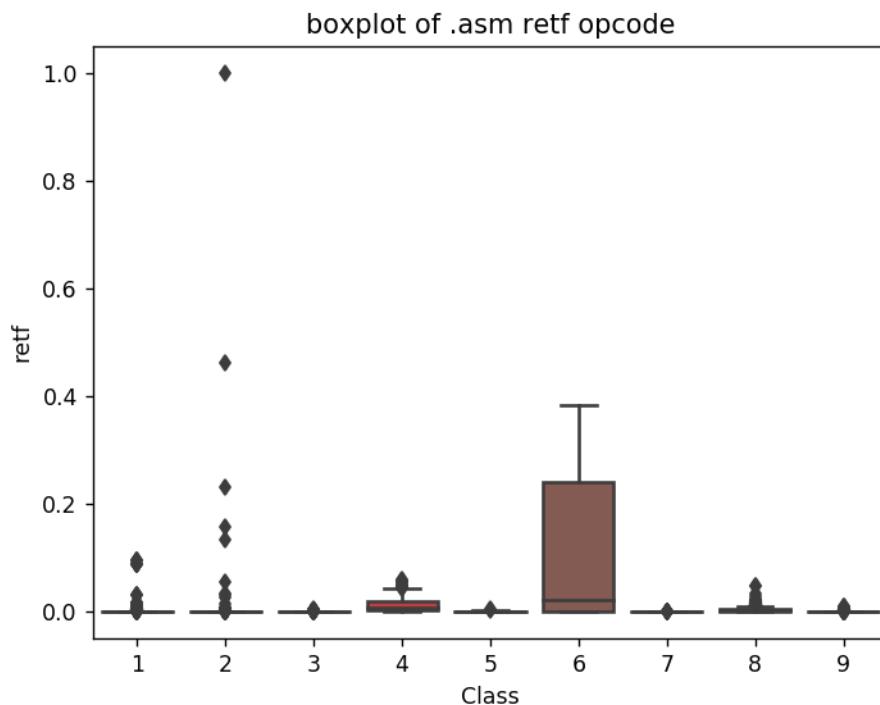


```
plot between Class label and mov opcode
Class 1 is having frequency of 2000 approx in 75 percentile of files
```

75 / 100

```
In [0]:
```

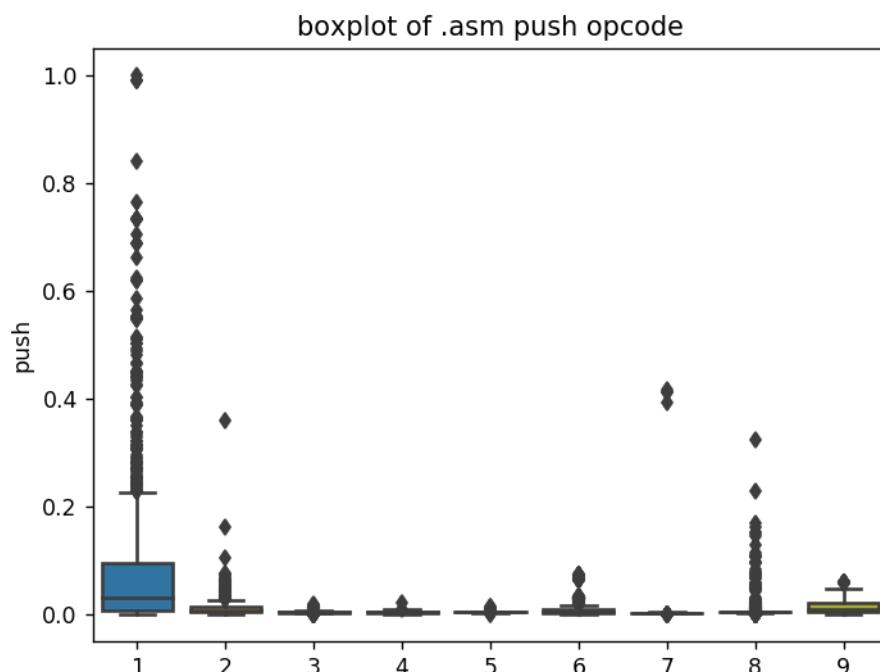
```
ax = sns.boxplot(x="Class", y="retf", data=result_asm)
plt.title("boxplot of .asm retf opcode")
plt.show()
```



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

```
In [0]:
```

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



Class

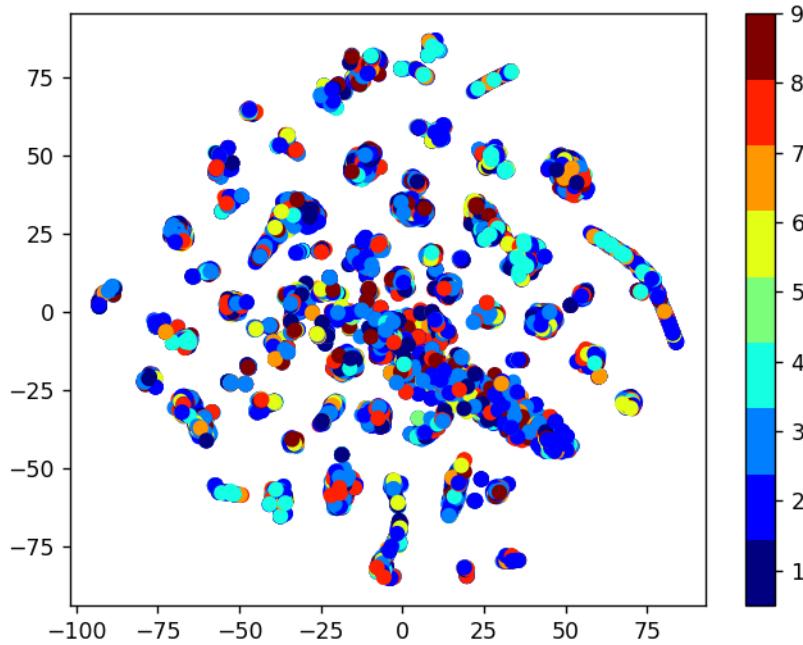
```
plot between push opcode and Class label
Class 1 is having 75 precentile files with push opcodes of frequency 1000
```

4.2.2 Multivariate Analysis on .asm file features

In [0]:

```
# check out the course content for more explantion on tsne algorithm
# https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/t-distributed-stochastic
-neighbourhood-embeddingt-sne-part-1/

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

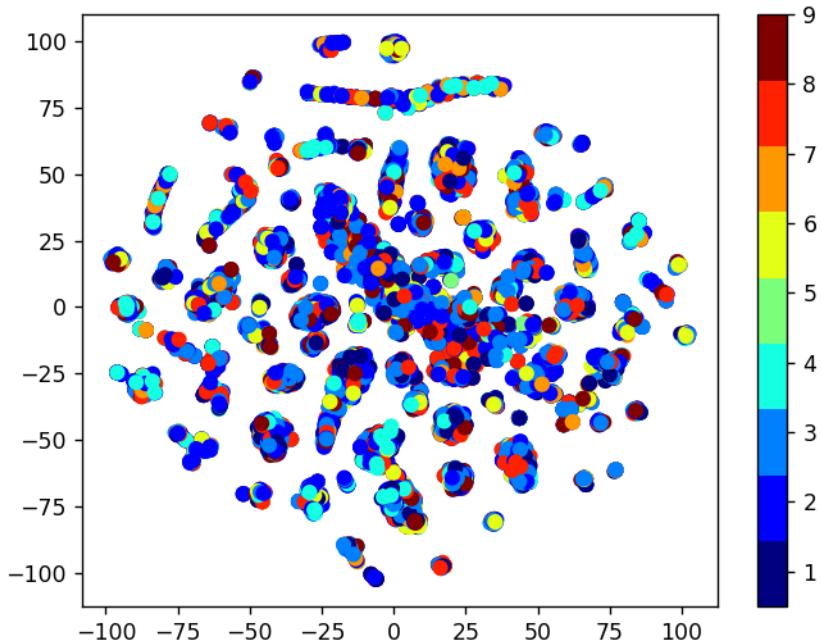


In [0]:

```
# by univariate analysis on the .asm file features we are getting very negligible information from
# 'rtn', '.BSS:' '.CODE' features, so heare we are trying multivariate analysis after removing tho
se features
# the plot looks very messy

xtsne=TSNE(perplexity=30)
results=xtsne.fit_transform(result_asm.drop(['ID','Class', 'rtn', '.BSS:', '.CODE','size'], axis=1))
vis_x = results[:, 0]
vis_y = results[:, 1]
plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
plt.colorbar(ticks=range(10))
plt.clim(0.5, 9)
```

```
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 [35]:

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

In [36]:

```
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x,asm_y ,stratify=asm_y,test_size=0.20)
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_asm,stratify=y_train_asm,test_size=0.20)
```

In [0]:

```
print( X_cv_asm.isnull().all())
```

```
HEADER:    False
.text:    False
.Pav:    False
.idata:    False
.data:    False
.bss:    False
```

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

4.4. Machine Learning models on features of .asm files

4.4.1 K-Nearest Neighbors

In [0]:

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

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


# find more about CalibratedClassifierCV here at http://scikit-
# learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default parameters
# "CalibratedClassifierCV(califier_type='base', estimator=None, method='sigmoid', n_
# classes=2, cv=5, verbose=0)"
```

```

# sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sigmoid', cv=5)
#
# some of the methods of CalibratedClassifierCV()
# fit(X, y[, sample_weight]) Fit the calibrated model
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict the target of new samples.
# predict_proba(X) Posterior probabilities of classification
#-----
# video link:
#-----

alpha = [x for x in range(1, 21, 2)]
cv_log_error_array=[]
for i in alpha:
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=k_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

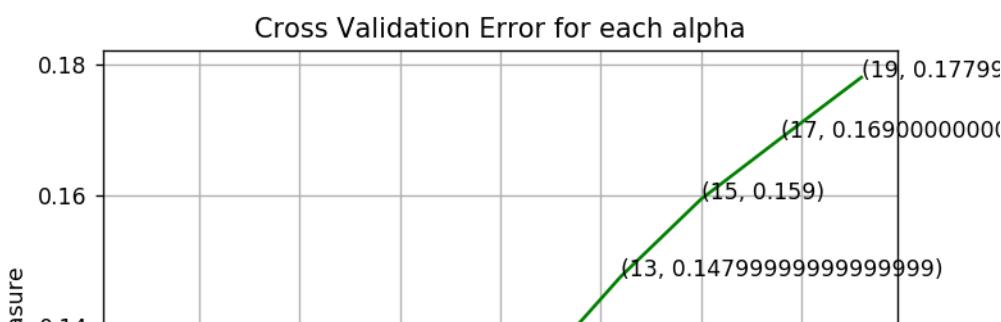
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

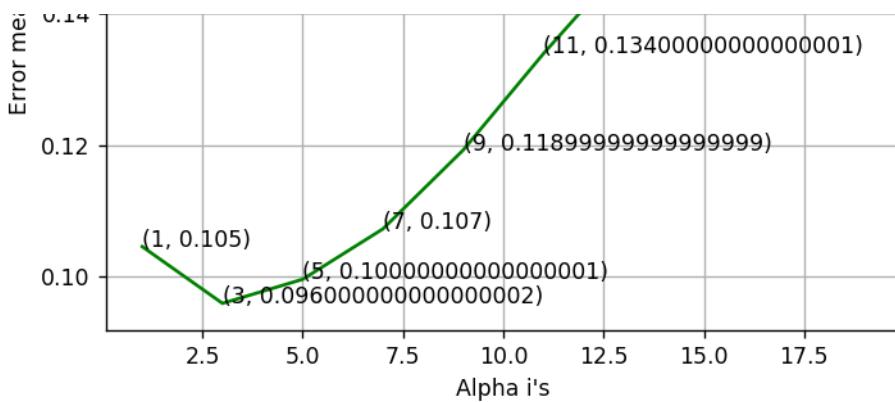
k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
k_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)
pred_y=sig_clf.predict(X_test_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',log_loss(y_train_asm, predict_y))
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',log_loss(y_cv_asm, predict_y))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',log_loss(y_test_asm, predict_y))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))

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

```





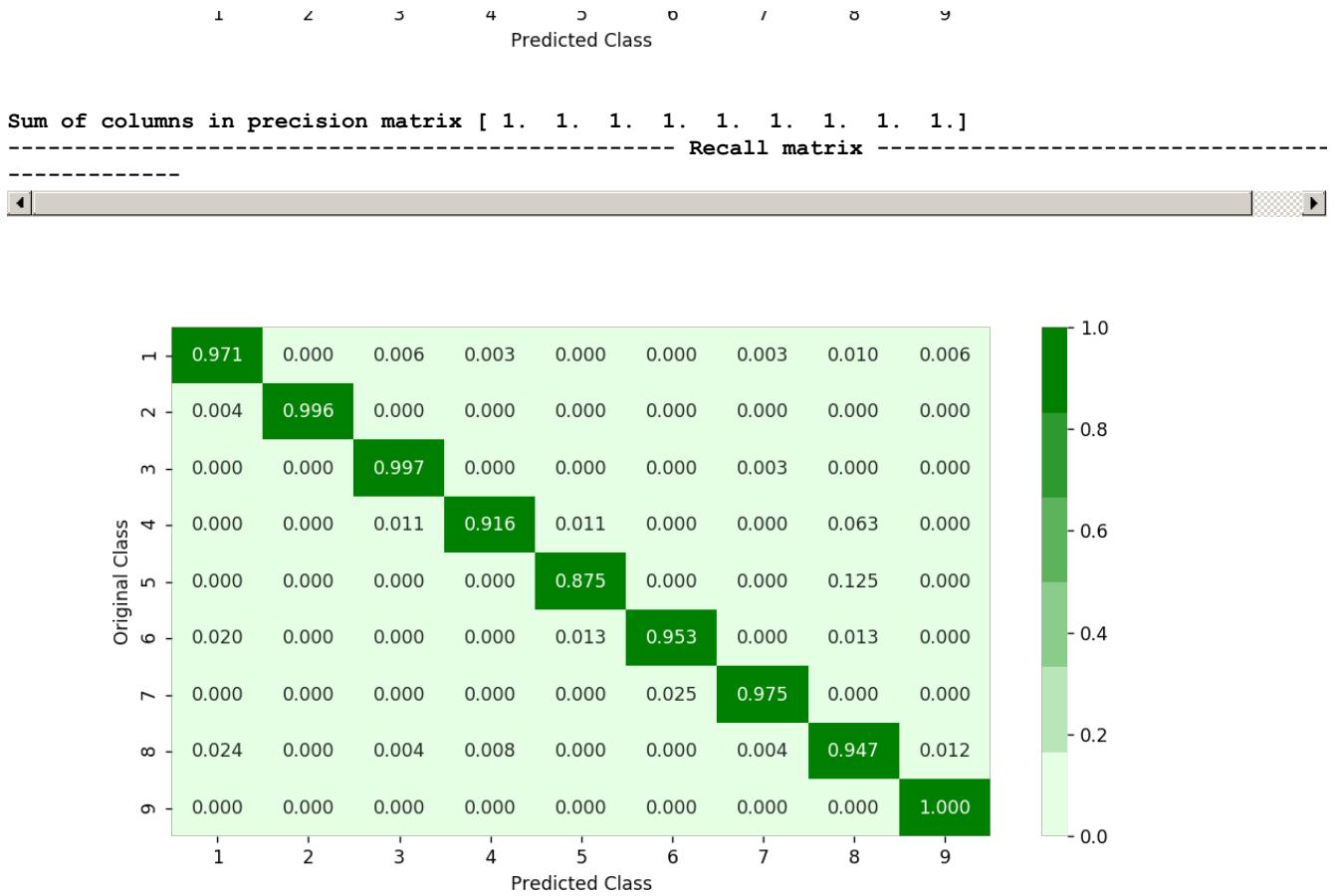
```
log loss for train data 0.0476773462198
log loss for cv data 0.0958800580948
log loss for test data 0.0894810720832
Number of misclassified points 2.02391904324
```

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



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





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

4.4.2 Logistic Regression

In [0]:

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

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

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


alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)
```

```

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

logisticR=LogisticRegression(penalty='l2',C=alpha[best_alpha],class_weight='balanced')
logisticR.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

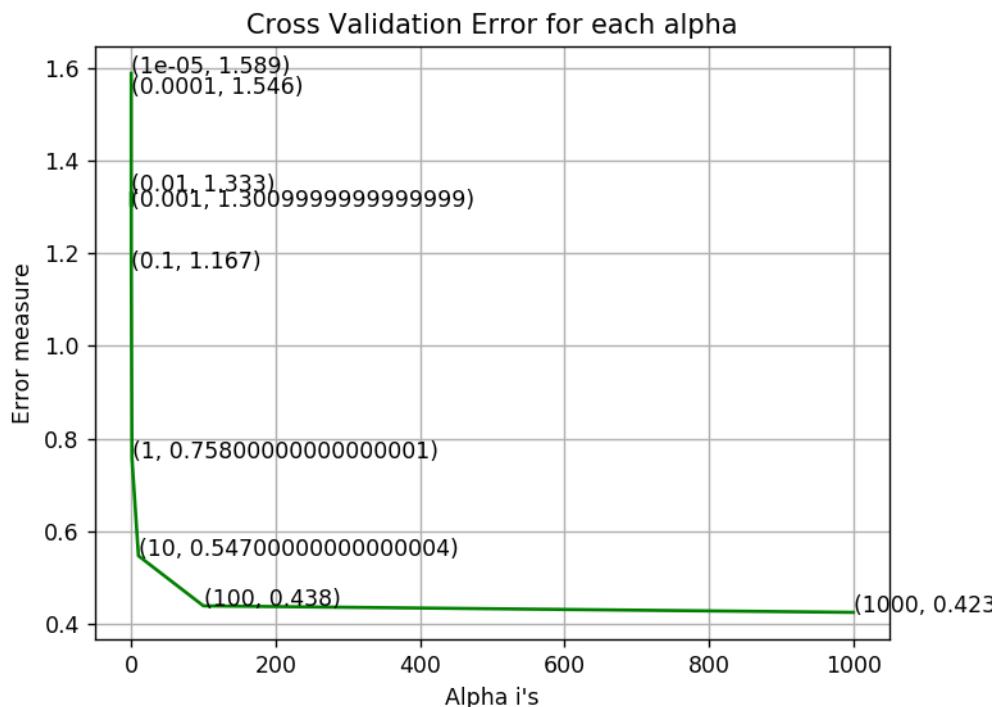
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=logisticR.classes_, eps=1e-15)))
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=logisticR.classes_, eps=1e-15)))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',(log_loss(y_test_asm, predict_y, labels=logisticR.classes_, eps=1e-15)))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))

```

```

log_loss for c =  1e-05 is 1.58867274165
log_loss for c =  0.0001 is 1.54560797884
log_loss for c =  0.001 is 1.30137786807
log_loss for c =  0.01 is 1.33317456931
log_loss for c =  0.1 is 1.16705751378
log_loss for c =  1 is 0.757667807779
log_loss for c =  10 is 0.546533939819
log_loss for c =  100 is 0.438414998062
log_loss for c =  1000 is 0.424423536526

```

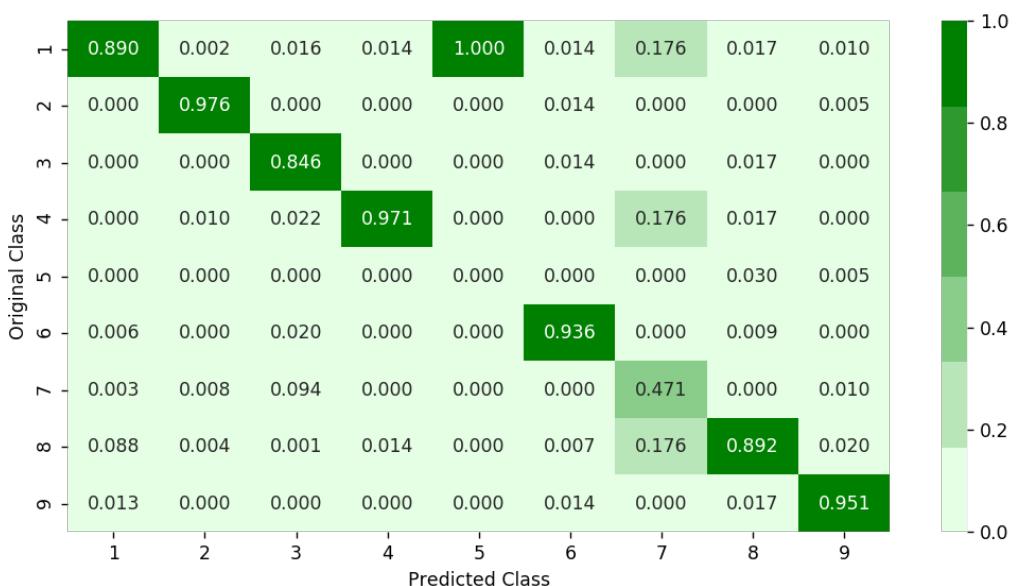


```

log loss for train data 0.396219394701
log loss for cv data 0.424423536526
log loss for test data 0.415685592517
Number of misclassified points  9.61361545538
----- Confusion matrix -----
-----
```

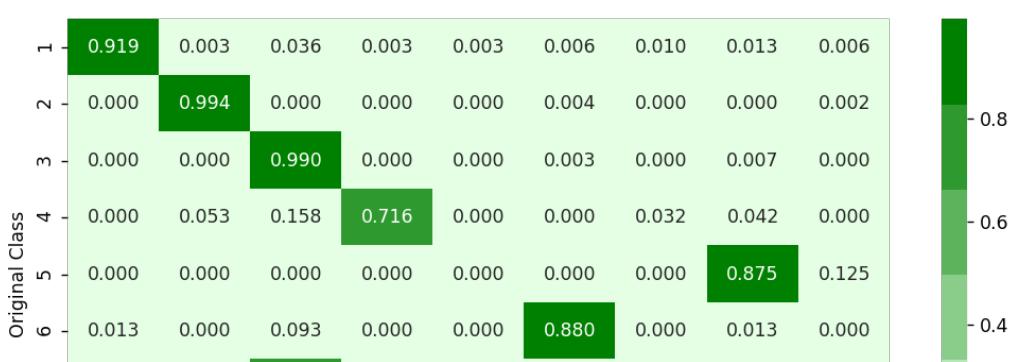


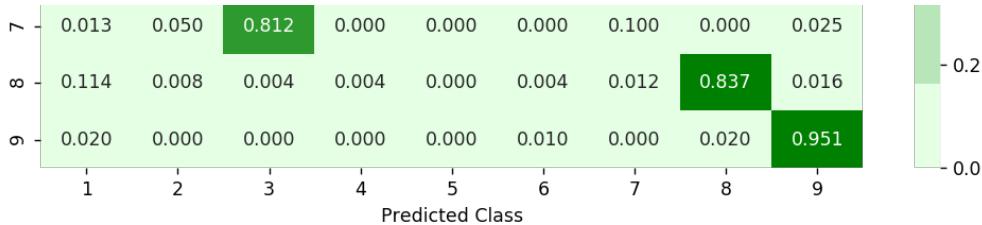
----- Precision matrix -----



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

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





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

4.4.3 Random Forest Classifier

In [0]:

```
# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None,
# verbose=0, warm_start=False,
# class_weight=None)

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

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

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


alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

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

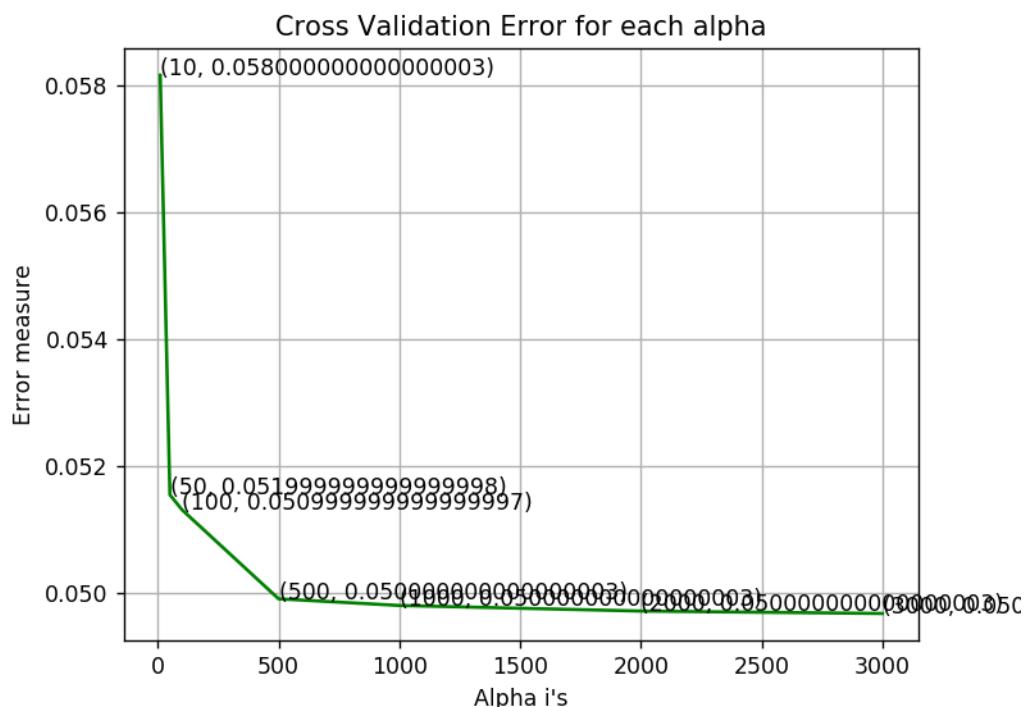
r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=sig_clf.classes_, eps=1e-15)))
predict_y = sig_clf.predict_proba(X_cv_asm)
```

```

print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=sig_clf.classes_, eps=1e-15)))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',(log_loss(y_test_asm, predict_y, labels=sig_clf.classes_, eps=1e-15)))
plot confusion matrix(v test asm.sig clf.predict(X test asm))

```

```
log_loss for c = 10 is 0.0581657906023
log_loss for c = 50 is 0.0515443148419
log_loss for c = 100 is 0.0513084973231
log_loss for c = 500 is 0.0499021761479
log_loss for c = 1000 is 0.0497972474298
log_loss for c = 2000 is 0.0497091690815
log_loss for c = 3000 is 0.0496706817633
```



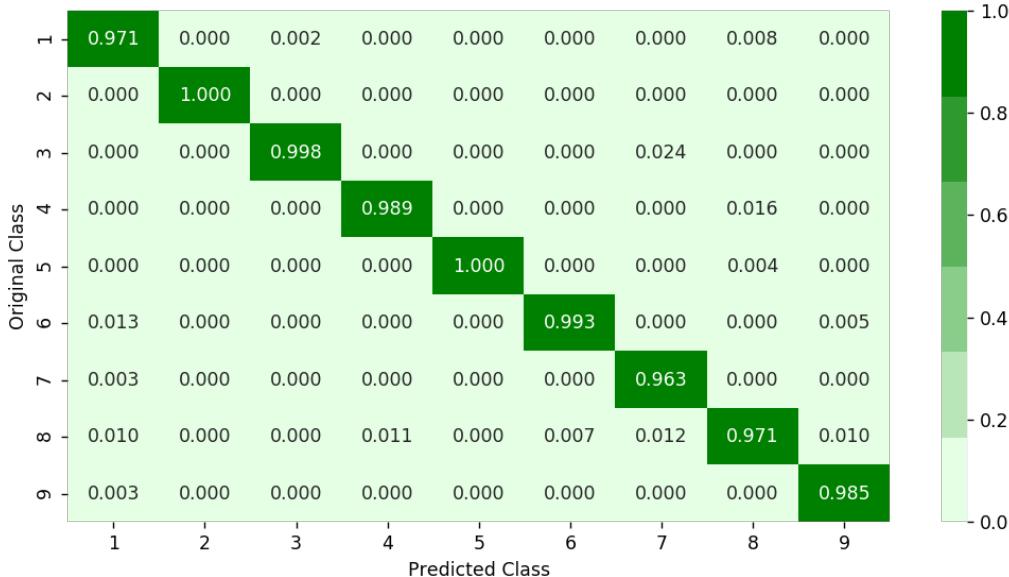
```
log loss for train data 0.0116517052676  
log loss for cv data 0.0496706817633  
log loss for test data 0.0571239496453  
Number of misclassified points 1.14995400184
```

Confusion matrix -----



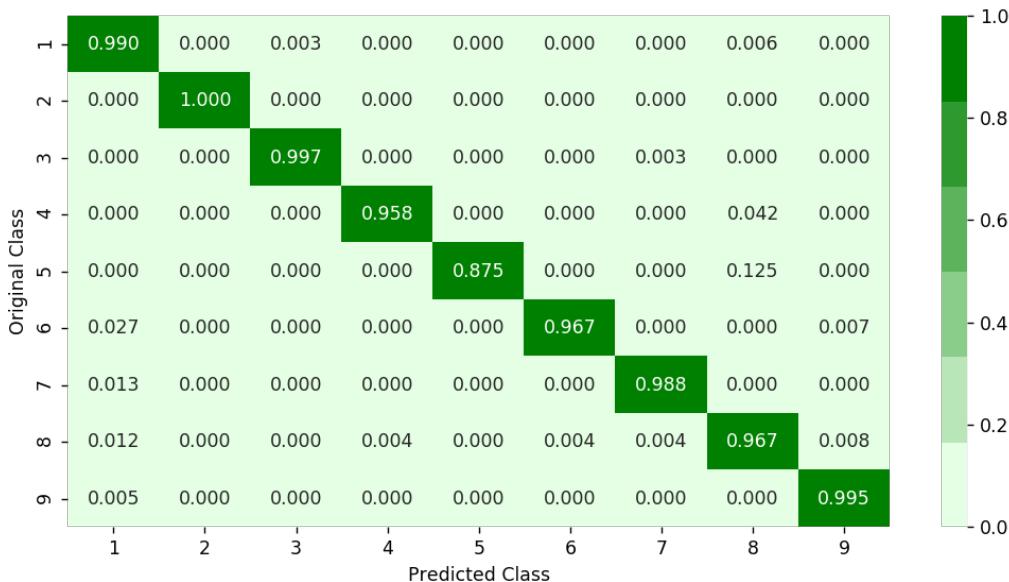
Predicted Class

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



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

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



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

4.4.4 XgBoost Classifier

In [0]:

```
# Training a hyper-parameter tuned Xg-Boost regressor on our train data  
# find more about XGBClассifier function here  
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier  
# -----
```

```

# default parameters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
min_child_weight=1,
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
reg_lambda=1,
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# ----

alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
    x_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=x_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

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

x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
x_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

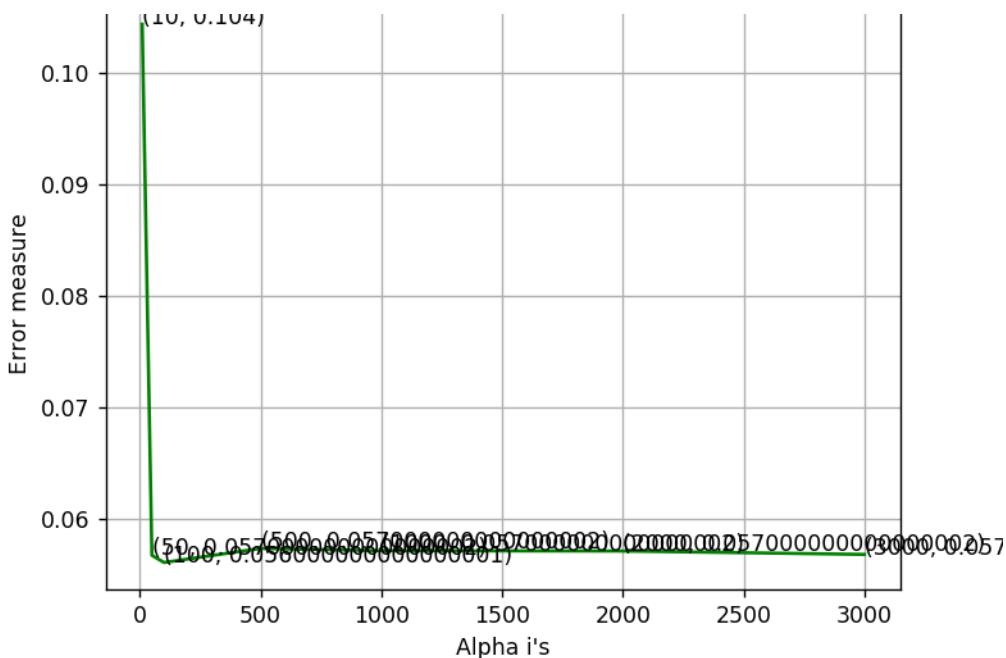
predict_y = sig_clf.predict_proba(X_train_asm)

print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train_asm, predict_y))
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv_asm, predict_y))
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test_asm, predict_y))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))

log_loss for c = 10 is 0.104344888454
log_loss for c = 50 is 0.0567190635611
log_loss for c = 100 is 0.056075038646
log_loss for c = 500 is 0.057336051683
log_loss for c = 1000 is 0.0571265109903
log_loss for c = 2000 is 0.057103406781
log_loss for c = 3000 is 0.0567993215778

```

Cross Validation Error for each alpha



For values of best alpha = 100 The train log loss is: 0.0117883742574
 For values of best alpha = 100 The cross validation log loss is: 0.056075038646
 For values of best alpha = 100 The test log loss is: 0.0491647763845
 Number of misclassified points 0.873965041398

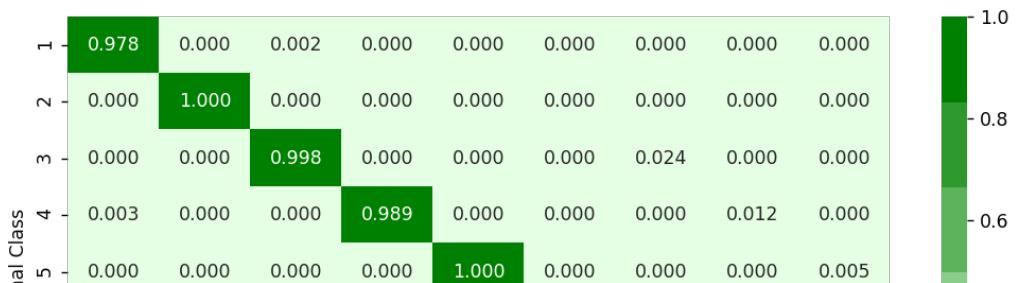
----- Confusion matrix -----

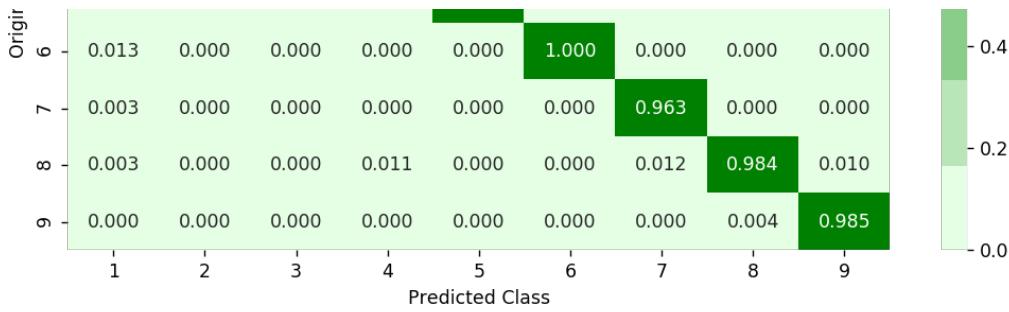
[◀ ▶]



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

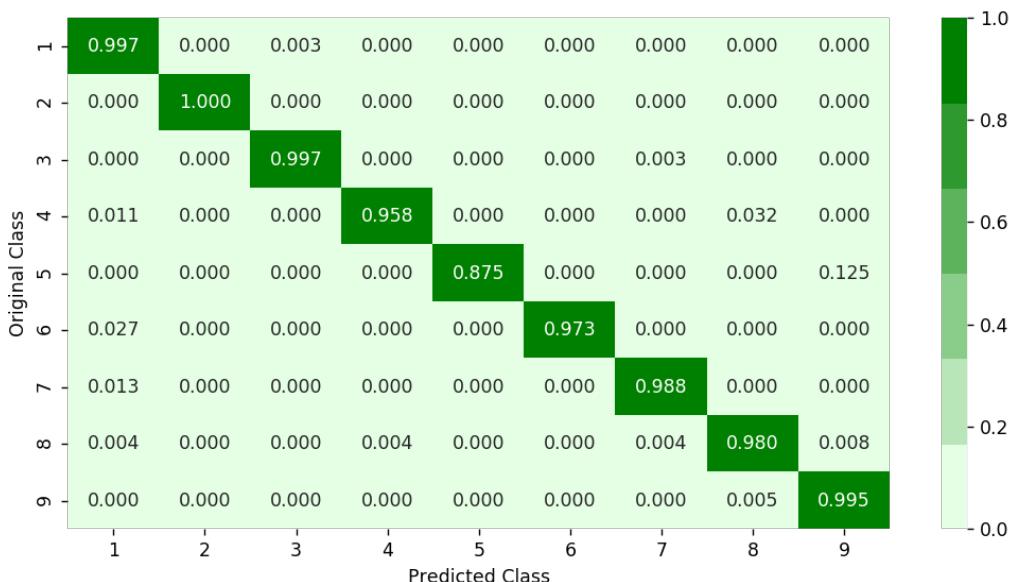
[◀ ▶]





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

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



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

4.4.5 Xgboost Classifier with best hyperparameters

In [0]:

```
x_cfl=XGBClassifier()

prams={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
    'n_estimators':[100,200,500,1000,2000],
    'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jobs=-1,)
random_cfl.fit(X_train_asm,y_train_asm)
```

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

```
[Parallel(n_jobs=-1)]: Done  2 tasks      | elapsed:   8.1s
[Parallel(n_jobs=-1)]: Done  9 tasks      | elapsed:  32.8s
[Parallel(n_jobs=-1)]: Done 19 out of 30 | elapsed: 1.1min remaining: 39.3s
[Parallel(n_jobs=-1)]: Done 23 out of 30 | elapsed: 1.3min remaining: 23.0s
[Parallel(n_jobs=-1)]: Done 27 out of 30 | elapsed: 1.4min remaining:  9.2s
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 2.3min finished
```

Out[0]:

```

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

```

In [0]:

```
print (random_cfl.best_params_)
```

```
{'subsample': 1, 'n_estimators': 200, 'max_depth': 5, 'learning_rate': 0.15, 'colsample_bytree': 0
.5}
```

In [0]:

```
# Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBClassifier function here
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# default parameters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
min_child_weight=1,
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
reg_lambda=1,
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbo
se=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is no
t thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-en
sembles/
# -----
```

```
x_cfl=XGBClassifier(n_estimators=200,subsample=0.5,learning_rate=0.15,colsample_bytree=0.5,max_dept
h=3)
```

```
x_cfl.fit(X_train_asm,y_train_asm)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train_asm,y_train_asm)
```

```
predict_y = c_cfl.predict_proba(X_train_asm)
print ('train loss',log_loss(y_train_asm, predict_y))
predict_y = c_cfl.predict_proba(X_cv_asm)
print ('cv loss',log_loss(y_cv_asm, predict_y))
predict_y = c_cfl.predict_proba(X_test_asm)
print ('test loss',log_loss(y_test_asm, predict_y))
```

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

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

4.5.1. Merging both asm and byte file features

In [13]:

```
result.head()
```

Out[13]:

	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	f9
0	0.000000	01azqd4InC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946	...	0.013560
1	0.000092	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984	...	0.001920
2	0.000184	01jsnpXSAlg6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155	...	0.009804
3	0.000276	01kcPWA9K2B0xQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481	...	0.002121
4	0.000368	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229	...	0.001530

5 rows × 261 columns

In [14]:

```
# we normalize the data each column
result_asm = normalize(result_asm)
```

100%|██████████| 54/54 [00:00<00:00, 803.10it/s]

In [15]:

```
result_asm.head()
```

Out[15]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	esi	eax
0	01kcPWA9K2B0xQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084	0.0	0.000072	...	0.000746	0.000301
1	1E93CpP60RHFNT5Qfvn	0.096045	0.001230	0.0	0.000617	0.000019	0.0	0.000000	0.0	0.000072	...	0.000328	0.000965
2	3ekVow2ajZHbTnBcsDfX	0.096045	0.000627	0.0	0.000300	0.000017	0.0	0.000038	0.0	0.000072	...	0.000475	0.000201
3	3X2nY7iQaPBIDrAZqJe	0.096045	0.000333	0.0	0.000258	0.000008	0.0	0.000000	0.0	0.000072	...	0.000090	0.000281
4	46OZzdsSKDCFV8h7XWxf	0.096045	0.000590	0.0	0.000353	0.000068	0.0	0.000000	0.0	0.000072	...	0.000102	0.000362

5 rows × 54 columns

In [16]:

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

(10868, 261)
(10868, 54)

In [17]:

```
result_asm.columns
```

Out[17]:

```
Index(['ID', 'HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:', '.edata:', '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE', 'jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'jz', ' rtn', 'lea', 'movzx', '.dll', 'std::', ':dword', 'edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip', 'Class', 'size'],
      dtype='object')
```

In [18]:

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

```
result_y = result_x['Class']
result_x = result_x.drop(['ID', 'rtn', '.BSS:', '.CODE', 'Class'], axis=1)
result_x.head()
```

Out[18]:

Unnamed:	0	1	2	3	4	5	6	7	8	...	edx	esi	...	
0	0.000000	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946	0.002638	...	0.015418	0.025875	0.0257
1	0.000092	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984	0.008267	...	0.004961	0.012316	0.0078
2	0.000184	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155	0.008104	...	0.000095	0.006181	0.0001
3	0.000276	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481	0.000959	...	0.000343	0.000746	0.0003
4	0.000368	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229	0.000376	...	0.000343	0.013875	0.0004

5 rows × 308 columns

In [19]:

```
result_y.head()
```

Out[19]:

```
0    9
1    2
2    9
3    1
4    8
Name: Class, dtype: int64
```

In [22]:

```
result_y.shape
```

Out[22]:

```
(10868,)
```

In [23]:

```
byte_features_with_size = pd.read_csv("result_with_size.csv")
asm_features_with_size = pd.read_csv("asmoutputfile.csv").drop(['rtn', '.BSS:', '.CODE'], axis=1)
```

In [24]:

```
normalized_bf = normalize(byte_features_with_size)
```

```
100%|██████████| 261/261 [00:00<00:00, 331.72it/s]
```

In [25]:

```
normalized_asm = normalize(asm_features_with_size)
```

```
100%|██████████| 49/49 [00:00<00:00, 779.84it/s]
```

In [122]:

```
normalized_asm.head()
```

Out[122]:

ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	:dword	edx	
0	01kcPWA9K2B0xQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084	0.0	0.000072	...	0.001028	0.000343
1	1E93CpP60RHFNiT5Qfvn	0.096045	0.001230	0.0	0.000617	0.000019	0.0	0.000000	0.0	0.000072	...	0.000975	0.000343

2	3ekVow2ajZHbTnBcsDfX	ID	0.096045	0.000627	Pav	0.0	0.000300	0.060077	date	.bss	0.000038	edata	0.0	0.000072000050	0.000248
3	3X2nY7iQaPBIWDrAZqJe		0.096045	0.000333	0.0	0.000258	0.000008	0.0	0.000000	0.0	0.0000720000188	0.000114			
4	46OZzdsSKDCFV8h7XWxf		0.096045	0.000590	0.0	0.000353	0.000068	0.0	0.000000	0.0	0.0000720000135	0.000229			

5 rows × 49 columns

In [120]:

```
normalized_bf = normalized_bf.drop('Class', axis=1)
```

In [123]:

```
normalized_bf.head()
```

Out[123]:

	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	f8
0	0.000000	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946	...	0.019969
1	0.000092	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984	...	0.035399
2	0.000184	01jsnpXSAlg6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155	...	0.012771
3	0.000276	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481	...	0.004728
4	0.000368	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229	...	0.005129

5 rows × 260 columns

In [27]:

```
#asm_features_with_size_normalized
normalized_asm.to_csv('normalized_asm.csv')
```

In [121]:

```
#byte_features_with_size_normalized
normalized_bf.to_csv('normalized_bf.csv')
```

In [26]:

```
result_y_df = result_y.to_frame()
```

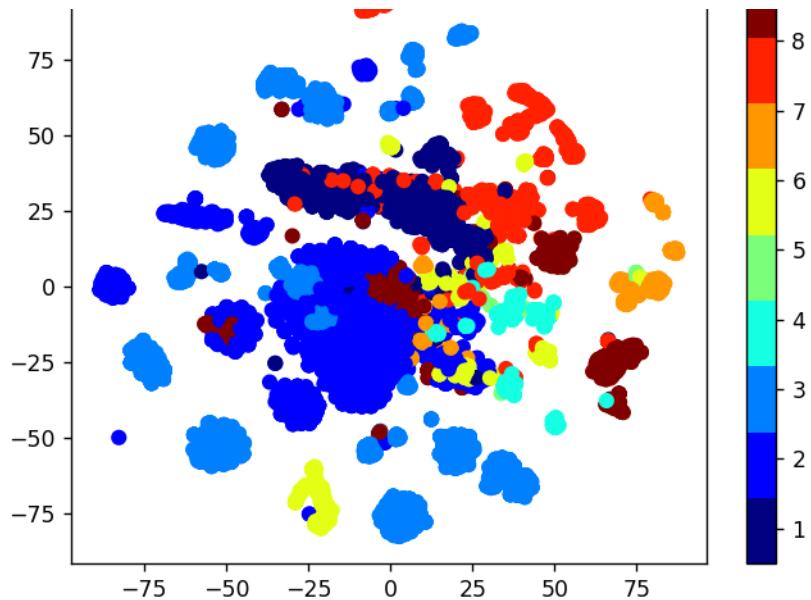
In [27]:

```
result_y_df.to_csv('result_y.csv')
```

4.5.2. Multivariate Analysis on final features

In [0]:

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



4.5.3. Train and Test split

In [0]:

```
X_train, X_test_merge, y_train, y_test_merge = train_test_split(result_x, result_y, stratify=result_y, test_size=0.20)
X_train_merge, X_cv_merge, y_train_merge, y_cv_merge = train_test_split(X_train, y_train, stratify=y_train, test_size=0.20)
```

4.5.4. Random Forest Classifier on final features

In [0]:

```
# -----
# default parameters
# sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None,
# verbose=0, warm_start=False,
# class_weight=None)

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

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

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


alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
from sklearn.ensemble import RandomForestClassifier
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_merge,y_train_merge)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_merge, y_train_merge)
    predict_y = sig_clf.predict_proba(X_cv_merge)
    cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=r_cfl.classes_, eps=1e-15))
```

```

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

best_alpha = np.argmin(cv_log_error_array)

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

r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
r_cfl.fit(X_train_merge,y_train_merge)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

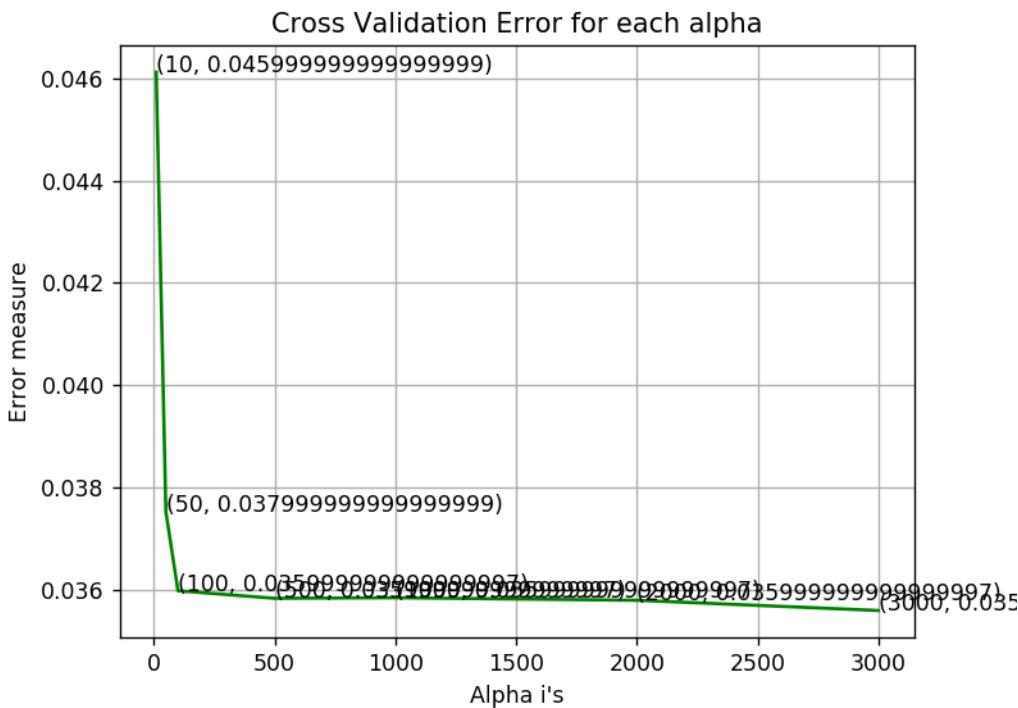
predict_y = sig_clf.predict_proba(X_train_merge)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss
is:",log_loss(y_train_merge, predict_y))
predict_y = sig_clf.predict_proba(X_cv_merge)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_lo
ss(y_cv_merge, predict_y))
predict_y = sig_clf.predict_proba(X_test_merge)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss
is:",log_loss(y_test_merge, predict_y))

```

```

log_loss for c =  10 is 0.0461221662017
log_loss for c =  50 is 0.0375229563452
log_loss for c =  100 is 0.0359765822455
log_loss for c =  500 is 0.0358291883873
log_loss for c =  1000 is 0.0358403093496
log_loss for c =  2000 is 0.0357908022178
log_loss for c =  3000 is 0.0355909487962

```



```

For values of best alpha =  3000 The train log loss is: 0.0166267614753
For values of best alpha =  3000 The cross validation log loss is: 0.0355909487962
For values of best alpha =  3000 The test log loss is: 0.0401141303589

```

4.5.5. XgBoost Classifier on final features

In [0]:

```
# Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBClассifier function here
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClассifier
# -----
# default parameters
# class xgboost.XGBClассifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
min_child_weight=1,
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
reg_lambda=1,
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)

# some of methods of RandomForrestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# -----


alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClассifier(n_estimators=i)
    x_cfl.fit(X_train_merge,y_train_merge)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(X_train_merge, y_train_merge)
    predict_y = sig_clf.predict_proba(X_cv_merge)
    cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=x_cfl.classes_, eps=1e-15))

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


best_alpha = np.argmin(cv_log_error_array)

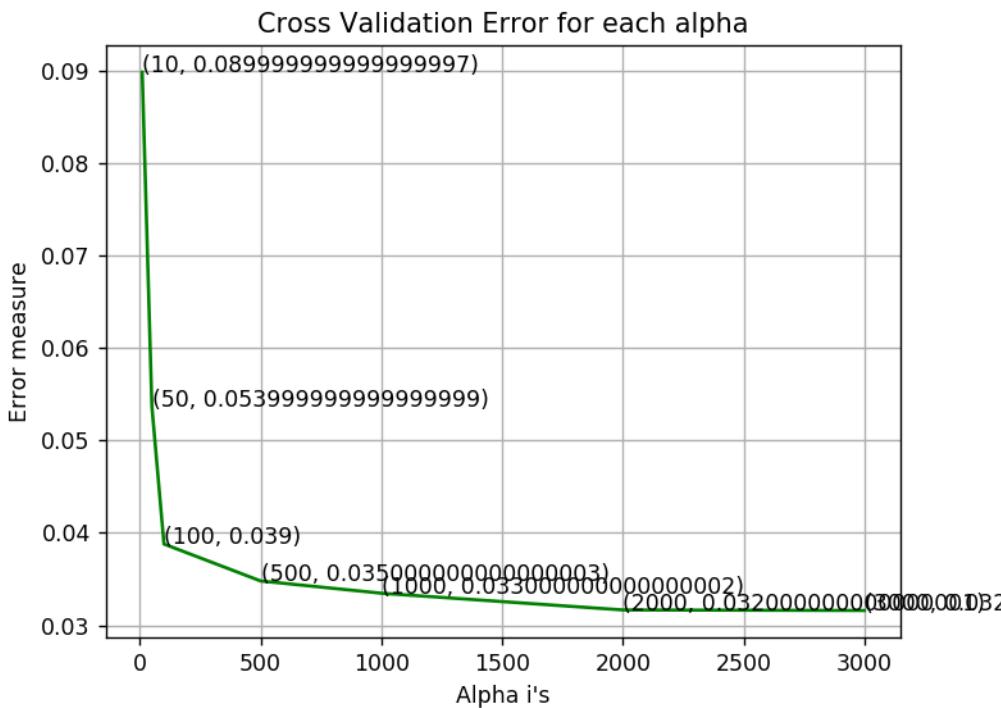
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

x_cfl=XGBClассifier(n_estimators=3000,nthread=-1)
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

predict_y = sig_clf.predict_proba(X_train_merge)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train_merge, predict_y))
predict_y = sig_clf.predict_proba(X_cv_merge)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv_merge, predict_y))
predict_y = sig_clf.predict_proba(X_test_merge)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test_merge, predict_y))

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

```
log_loss for c = 2000 is 0.0316569078846
log_loss for c = 3000 is 0.0315972694477
```



```
For values of best alpha = 3000 The train log loss is: 0.0111918809342
For values of best alpha = 3000 The cross validation log loss is: 0.0315972694477
For values of best alpha = 3000 The test log loss is: 0.0323978515915
```

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

In [0]:

```
x_cfl=XGBClassifier()

prams={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
    'n_estimators':[100,200,500,1000,2000],
    'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jobs=-1,)
random_cfl.fit(X_train_merge, y_train_merge)
```

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

```
[Parallel(n_jobs=-1)]: Done   2 tasks      | elapsed:  1.1min
[Parallel(n_jobs=-1)]: Done   9 tasks      | elapsed:  2.2min
[Parallel(n_jobs=-1)]: Done  19 out of 30 | elapsed:  4.5min remaining:  2.6min
[Parallel(n_jobs=-1)]: Done  23 out of 30 | elapsed:  5.8min remaining:  1.8min
[Parallel(n_jobs=-1)]: Done  27 out of 30 | elapsed:  6.7min remaining:  44.5s
[Parallel(n_jobs=-1)]: Done  30 out of 30 | elapsed:  7.4min finished
```

Out[0]:

```
RandomizedSearchCV(cv=None, error_score='raise',
                    estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, colsample_bytree=1,
                    gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
                    min_child_weight=1, missing=None, n_estimators=100, nthread=-1,
                    objective='binary:logistic', reg_alpha=0, reg_lambda=1,
                    scale_pos_weight=1, seed=0, silent=True, subsample=1),
                    fit_params=None, iid=True, n_iter=10, n_jobs=-1,
                    param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15, 0.2], 'n_estimators':
```

```
[100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'colsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},  
    pre_dispatch='2*n_jobs', random_state=None, refit=True,  
    return_train_score=True, scoring=None, verbose=10)
```

In [0]:

```
print (random_cfl.best_params_)
```

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

In [0]:

```
# find more about XGBClассifier function here  
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClассifier  
# -----  
# default parameters  
# class xgboost.XGBClассifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,  
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,  
min_child_weight=1,  
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,  
reg_lambda=1,  
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)  
  
# some of methods of RandomForestRegressor()  
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)  
# get_params([deep]) Get parameters for this estimator.  
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.  
# get_score(importance_type='weight') -> get the feature importance  
# -----  
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/  
# -----  
  
x_cfl=XGBClассifier(n_estimators=1000,max_depth=10,learning_rate=0.15,colsample_bytree=0.3,subsample=1,nthread=-1)  
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)  
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")  
sig_clf.fit(X_train_merge, y_train_merge)  
  
predict_y = sig_clf.predict_proba(X_train_merge)  
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss  
is:",log_loss(y_train_merge, predict_y))  
predict_y = sig_clf.predict_proba(X_cv_merge)  
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv_merge, predict_y))  
predict_y = sig_clf.predict_proba(X_test_merge)  
print('For values of best alpha = ', alpha[best_alpha], "The test log loss  
is:",log_loss(y_test_merge, predict_y))  
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_merge))
```

```
For values of best alpha = 3000 The train log loss is: 0.0121922832297  
For values of best alpha = 3000 The cross validation log loss is: 0.0344955487471  
For values of best alpha = 3000 The test log loss is: 0.0317041132442
```

5. Assignments

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 <=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

Bi-grams

```
In [95]:
```

```
byte_vocab =
"00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,17,18,19,1a,1b,1c,1d,1e,1f,20
22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,34,35,36,37,38,39,3a,3b,3c,3d,3e,3f,40,41,42,
4,45,46,47,48,49,4a,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,
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,
89,8a,8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,9a,9b,9c,9d,9e,9f,a0,a1,a2,a3,a4,a5,a6,a7,a8,a9,
b,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,c
,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
f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,ff,??"
```

```
In [96]:
```

```
byte_bigram_vocab = []
```

```
In [97]:
```

```
def byte_bigram():
    for i, v in enumerate(byte_vocab.split(',')):
        for j in range(0, len(byte_vocab.split(','))):
            byte_bigram_vocab.append(v + ' ' + byte_vocab.split(',')[j])
len(byte_bigram_vocab)
```

```
In [98]:
```

```
byte_bigram()
```

```
In [99]:
```

```
byte_bigram_vocab[:5]
```

```
Out[99]:
```

```
['00 00', '00 01', '00 02', '00 03', '00 04']
```

```
In [11]:
```

```
byte_trigram_vocab = []
```

```
In [12]:
```

```
def byte_trigram():
    byte_trigram_vocab = []
    for i, v in enumerate(byte_vocab.split(',')):
        for j in range(0, len(byte_vocab.split(','))):
            for k in range(0, len(byte_vocab.split(','))):
                byte_trigram_vocab.append(v + ' ' + byte_vocab.split(',')[j] + ' ' + byte_vocab.split(',')[k])
len(byte_trigram_vocab)
```

```
In [13]:
```

```
byte_trigram()
```

```
In [14]:
```

```
len(byte_bigram_vocab)
```

```
Out[14]:
```

```
66049
```

```
In [15]:
```

```
import pickle
filename = 'trigram'
outfile = open(filename, 'wb')
```

```
pickle.dump(byte_trigram_vocab,outfile)
outfile.close()
```

In [16]:

```
from tqdm import tqdm
from sklearn.feature_extraction.text import CountVectorizer
import scipy

vect = CountVectorizer(lowercase=False,ngram_range=(2,2), vocabulary=byte_bigram_vocab)
byte_bigram_vect = scipy.sparse.csr_matrix((10868, len(byte_bigram_vocab)))
for i, file in tqdm(enumerate(os.listdir('byteFiles'))):
    f = open('byteFiles/' + file)
    byte_bigram_vect[i,:] += scipy.sparse.csr_matrix(vect.fit_transform([f.read().replace('\n', ' ').lower()]))
    f.close()

10868it [15:30:56, 8.98s/it]
```

In [18]:

```
# import os.path
# os.path.exists('/media/ubuntu/E210D7AA10D783C7/train/01azqd4InC7m9JpocGv5.bytes')
```

In [17]:

```
scipy.sparse.save_npz('bytebigram.npz', byte_bigram_vect)
```

In [21]:

```
byte_bigram_vect = scipy.sparse.load_npz('bytebigram.npz')
```

In [22]:

```
byte_bigram_vect.shape
```

Out[22]:

```
(10868, 66049)
```

In [23]:

```
from sklearn.preprocessing import normalize

byte_bigram_vect_normalized = normalize(byte_bigram_vect, axis=0)
```

In [24]:

```
scipy.sparse.save_npz('byte_bigram_vect_normalized.npz', byte_bigram_vect_normalized)
```

N-Grams

In [7]:

```
#Ref https://www.edwardraff.com/publications/what_can_ngrams_learn.pdf
#Ref https://github.com/melanieihuei/Malware-Classification
```

```
opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'jz', 'rtn', 'lea', 'movzx']
```

In [30]:

```
asm_bigram = []
def asmopcodebigram():
    for i, v in tqdm(enumerate(opcodes)):
        for j in range(0, len(opcodes)):
```

```
        asm_bigram.append(v + ' ' + opcodes[j])
```

In [31]:

```
asmopcodebigram()
len(asm_bigram)
```

```
26it [00:00, 72701.27it/s]
```

Out[31]:

```
676
```

In [71]:

```
asm_trigram = []
def asmopcodetrigram():
    for i, v in enumerate(opcodes):
        for j in range(0, len(opcodes)):
            for k in range(0, len(opcodes)):
                asm_trigram.append(v + ' ' + opcodes[j] + ' ' + opcodes[k])
```

In [72]:

```
asmopcodetrigram()
len(asm_trigram)
```

Out[72]:

```
17576
```

In [8]:

```
asm_4gram = []
for i, v in enumerate(opcodes):
    for j in range(0, len(opcodes)):
        for k in range(0, len(opcodes)):
            for l in range(0, len(opcodes)):
                asm_4gram.append(v + ' ' + opcodes[j] + ' ' + opcodes[k] + ' ' + opcodes[l])
```

In [9]:

```
len(asm_4gram)
```

Out[9]:

```
456976
```

In [9]:

```
def opcode_collect():
    op_file = open("opcode_file.txt", "wt")
    for asmfile in os.listdir('asmFiles'):
        opcode_str = ""
        with codecs.open('asmFiles/' + asmfile, encoding='cp1252', errors ='replace') as fli:
            for lines in fli:
                line = lines.rstrip().split()
                for li in line:
                    if li in opcodes:
                        opcode_str += li + ' '
        op_file.write(opcode_str + "\n")
    op_file.close()
opcode_collect()
```

In [12]:

```
from sklearn.feature_extraction.text import CountVectorizer
import scipy
```

```
import numpy  
  
vect = CountVectorizer(ngram_range=(2, 2), vocabulary = asm_bigram)  
bigram_vect = scipy.sparse.csr_matrix((10868, len(asm_bigram)))  
raw_opcode = open('opcode_file.txt').read().split('\n')
```

In [13]:

```
for i in range(10868):  
    bigram_vect[i, :] += scipy.sparse.csr_matrix(vect.transform([raw_opcode[i]]))
```

In [15]:

```
bigram_vect
```

Out[15]:

```
<10868x676 sparse matrix of type '<class 'numpy.float64'>'  
with 1877099 stored elements in Compressed Sparse Row format>
```

In [16]:

```
scipy.sparse.save_npz('op_bigram.npz', bigram_vect)
```

In [17]:

```
vect = CountVectorizer(ngram_range=(3, 3), vocabulary = asm_trigram)  
trigram_vect = scipy.sparse.csr_matrix((10868, len(asm_trigram)))  
raw_opcode = open('opcode_file.txt').read().split('\n')
```

In [18]:

```
for i in range(10868):  
    trigram_vect[i, :] += scipy.sparse.csr_matrix(vect.transform([raw_opcode[i]]))
```

In [19]:

```
scipy.sparse.save_npz('op_trigram.npz', trigram_vect)
```

In [34]:

```
import scipy  
opcodebivect=scipy.sparse.load_npz('op_bigram.npz')
```

In [70]:

```
opcodedrivect=scipy.sparse.load_npz('op_trigram.npz')
```

In [24]:

```
asmopcodeitetragram = asm_4gram
```

In [25]:

```
vect = CountVectorizer(ngram_range=(4, 4), vocabulary = asmopcodeitetragram)  
opcodedetetravect = scipy.sparse.csr_matrix((10868, len(asopcodeitetragram)))  
  
for i in range(10868):  
    opcodedetetravect[i, :] += scipy.sparse.csr_matrix(vect.transform([raw_opcode[i]]))
```

In [26]:

```
opcodedetetravect
```

Out[26]:

```
<10868x456976 sparse matrix of type '<class 'numpy.float64'>'  
with 16604101 stored elements in Compressed Sparse Row format>
```

In [27]:

```
scipy.sparse.save_npz('opcodetetragram.npz', opcodetetravect)
```

In [3]:

```
opcodetetravect = scipy.sparse.load_npz('opcodetetragram.npz')
```

In [4]:

```
result_y = pd.read_csv('result_y.csv')  
result_y = result_y.drop('Unnamed: 0', axis=1)
```

In [60]:

```
import array  
from tqdm import tqdm  
def collect_img_asm():  
    #pix_file = open("../pixels.txt", "w+")  
    for asmfile in tqdm(os.listdir("./asmFiles")):  
        file_name = asmfile.split('.')[0]  
        file = codecs.open("./asmFiles/" + asmfile, 'rb')  
        file_len = os.path.getsize("./asmFiles/" + asmfile)  
        width = int(file_len ** 0.5)  
        rem = int(file_len / width)  
        arr = array.array('B')  
        arr.frombytes(file.read())  
        file.close()  
        reshaped = np.reshape(arr[:width * width], (width, width))  
        reshaped = np.uint8(reshaped)  
        scipy.misc.imsave('asm_image/' + file_name + '.png', reshaped)
```

In [61]:

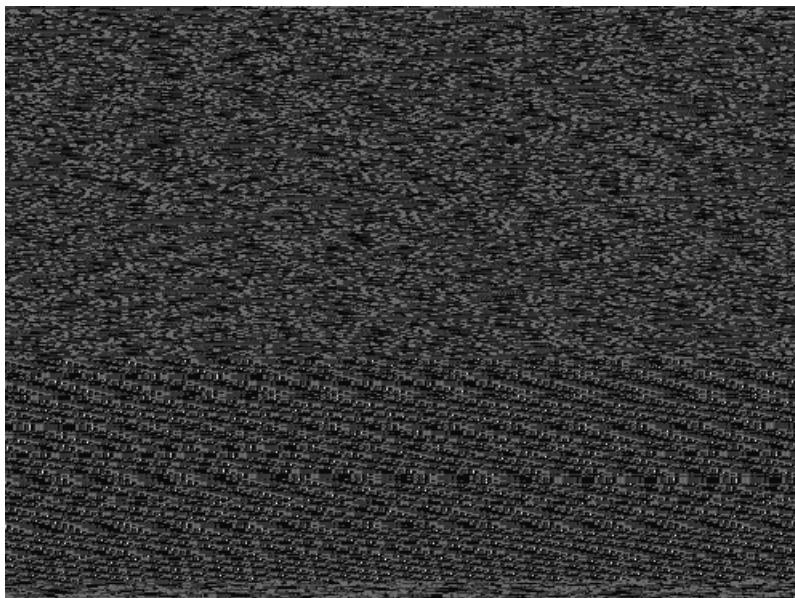
```
collect_img_asm()
```

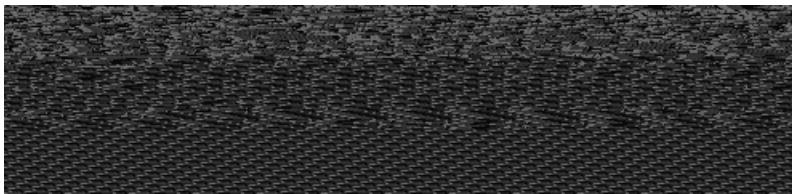
```
100%|██████████| 10868/10868 [1:21:06<00:00, 1.22s/it]
```

In [62]:

```
from IPython.display import Image  
Image(filename='asm_image/deTXH9Zau7qmM0yfYsRS.png')
```

Out[62]:





Features from First 500 image pixels

In [35]:

```
import cv2
image_features = np.zeros((10868, 500))
for i, asmfile in tqdm(enumerate(os.listdir("asmFiles"))):
    img = cv2.imread("asm_image/" + asmfile.split('.')[0] + '.png')
    img_arr = img.flatten()[:500]
    image_features[i, :] += img_arr
```

10868it [20:06, 7.24it/s]

In [38]:

```
image_features.shape
```

Out[38]:

(10868, 500)

In [39]:

```
import numpy as np
np.savetxt("image_features_500.txt", image_features, fmt="%s")
```

In [110]:

```
image_features = np.loadtxt("image_features.txt", dtype=float)
# image_features
```

In [40]:

```
from sklearn.preprocessing import normalize
img_feat = []
for i in tqdm(range(500)):
    img_feat.append('pix' + str(i))
img_final = pd.DataFrame(normalize(image_features, axis = 0), columns = img_feat)
```

100%|██████████| 500/500 [00:00<00:00, 1210826.79it/s]

In [41]:

```
img_final['ID'] = result.ID
```

In [42]:

```
img_final.head()
```

Out[42]:

	pix0	pix1	pix2	pix3	pix4	pix5	pix6	pix7	pix8	pix9	...	pix491	pix492	pix493
0	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320	0.007913	...	0.002942	0.010763	0.01076
1	0.006560	0.006560	0.006560	0.013504	0.013504	0.013504	0.012927	0.012927	0.012927	0.013963	...	0.019285	0.003312	0.00331
2	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320	0.007913	...	0.002942	0.010763	0.01076
3	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320	0.007913	...	0.002942	0.010763	0.01076

```
4 0.010268 0.010268 0.010268 0.008033 0.008033 0.008033 0.008033 0.008320 0.008320 0.008320 0.007919 ... 0.02942 0.010763 0.010763
```

5 rows x 501 columns

In [43]:

```
img_final.to_csv('img_final_500.csv')
```

500 most Important features in opcode bigram

In [44]:

```
from sklearn.preprocessing import normalize

normalized_opcodebivect = normalize(opcodebivect, axis = 0)
```

In [12]:

```
def imp_features(data, features, keep):
    rf = RandomForestClassifier(n_estimators = 100, n_jobs = -1)
    rf.fit(data, result_y)
    imp_feature_indx = np.argsort(rf.feature_importances_)[-1]
    imp_value = np.take(rf.feature_importances_, imp_feature_indx[:20])
    #print(imp_feature_indx)
    imp_feature_name = np.take(features, imp_feature_indx[:20])
    sns.set()
    plt.figure(figsize = (10, 5))
    ax = sns.barplot(x = imp_feature_name, y = imp_value)
    ax.set_xticklabels(labels = imp_feature_name, rotation = 45)
    sns.set_palette(reversed(sns.color_palette("husl", 10)), 10)
    plt.title('Important Features')
    plt.xlabel('Feature Names')
    plt.ylabel('Importance')
    plt.show()
    return imp_feature_indx[:keep]
```

In [55]:

```
len(asn_bigram)
```

Out[55]:

676

In [48]:

```
normalized_opcodebivect.shape
```

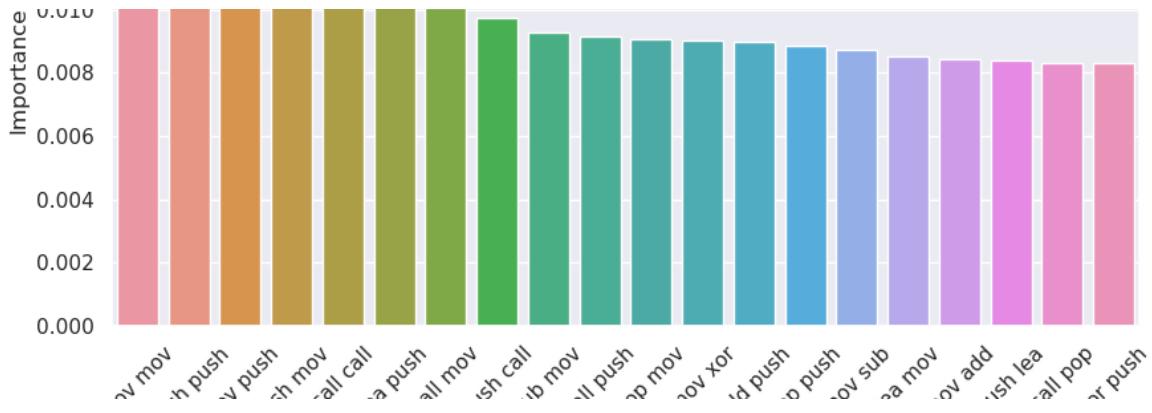
Out[48]:

(10868, 676)

In [56]:

```
opcode_bigram_indexes = imp_features(normalized_opcodebivect, asn_bigram, 500)
```





In [57]:

```
len(opcode_bigram_indexes)
```

Out[57]:

500

In [60]:

```
op_bi_df = pd.SparseDataFrame(normalized_opcodebivect, columns = asm_bigram)
```

In [62]:

```
for col in tqdm(op_bi_df.columns):
    if col not in np.take(asm_bigram, opcode_bigram_indexes):
        op_bi_df.drop(col, axis = 1, inplace = True)
```

100%|██████████| 676/676 [00:09<00:00, 70.65it/s]

In [63]:

```
op_bi_df['ID'] = result.ID
```

In [65]:

```
op_bi_df = op_bi_df.to_dense()
```

In [66]:

```
op_bi_df = op_bi_df.fillna(0)
```

In [76]:

```
op_bi_df.isnull().values.any()
```

Out[76]:

False

In [67]:

```
op_bi_df.head()
```

Out[67]:

	jmp	jmp	jmp mov	jmp retf	jmp push	jmp pop	jmp xor	jmp retn	jmp nop	jmp sub	jmp inc	...	movzx add	movzx imul	movzx or	movzx cmp
0	0.031815	0.003894	0.000000	0.00042	0.000000	0.002374	0.000000	0.0	0.00895	0.0	...	0.000000	0.0	0.000000	0.000471	
1	0.000000	0.000649	0.000000	0.00021	0.000374	0.000419	0.000703	0.0	0.00000	0.0	...	0.001429	0.0	0.004541	0.000471	

2	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	j _h _b	0.000000	j _h _b	...	0.001144	mov _{zx}	0.009839	0.005648
3	0.000000	0.000101	0.000000	0.000007	0.000000	0.000000	0.000000	0.000000	0.000000	n _o _B	0.000000	sub	...	0.000000	imul	0.000000	0.000000
4	0.000362	0.001156	0.001467	0.00028	0.000374	0.000140	0.000000	0.0	0.000000	0.0	0.000000	0.0	...	0.000000	0.0	0.001514	0.005177

5 rows x 501 columns

In [68]:

```
op_bi_df.to_csv('op_bigram_500.csv')
```

1000 Important Feature Among Trigram

In [73]:

```
from sklearn.preprocessing import normalize
normalize_opcodetrivect = normalize(opcodetrivect, axis=0)
```

In [77]:

```
len(asn_trigram)
```

Out[77]:

17576

In [80]:

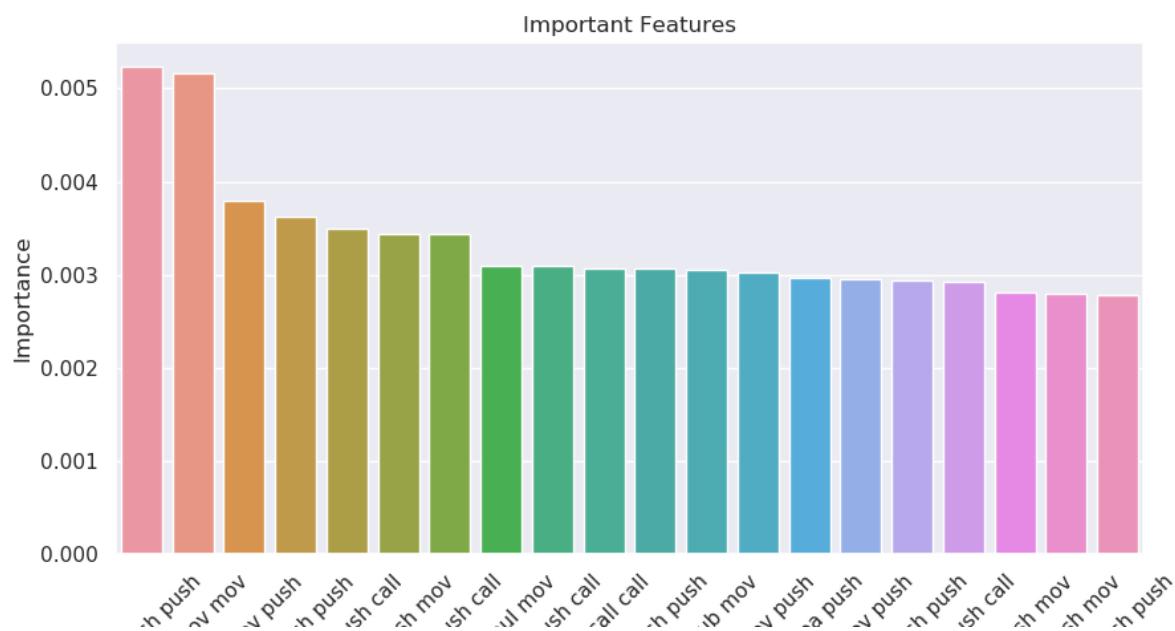
```
normalize_opcodetrivect.shape
```

Out[80]:

(10868, 17576)

In [81]:

```
opcode_trigram_indexes = imp_features(normalize_opcodetrivect, asn_trigram, 1000)
```



In [83]:

```
len(opcode_trigram_indexes)
```

Out[83]:

1000

In [82]:

```
op_tri_df = pd.SparseDataFrame(normalize_opcodetrivect, columns = asm_trigram)
```

In [84]:

```
op_tri_df = op_tri_df.loc[:, np.intersect1d(op_tri_df.columns, np.take(asm_trigram, opcode_trigram_indexes))]
```

In [85]:

```
op_tri_df['ID'] = result.ID
```

In [86]:

```
op_tri_df.head()
```

Out[86]:

	add add add	add add add add cmp	add add jmp	add add lea	add add mov	add add or	add add pop	add add push	add add sub	add add xor	xor sub cmp	xor sub mov	xor sub push	xor a			
0	NaN	NaN	0.002522	NaN	0.005422	0.00267		NaN	0.004902	0.000662	0.00138	...	0.004508	0.018631	0.005498	N		
1	0.0002	0.002123		NaN	NaN	0.003425	0.00534		NaN	0.004902	0.000662	NaN	...	NaN	0.003388	0.005498	N	
2	0.0001		NaN	NaN	NaN	NaN	NaN		NaN	0.003268		NaN	NaN	...	NaN	NaN	0.021993	0.0058
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN		NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	N
4	0.0001		NaN	0.002522	NaN	0.000571		NaN	0.001891		NaN	NaN	NaN	...	NaN	NaN	NaN	N

5 rows × 1001 columns

In [88]:

```
op_tri_df = op_tri_df.to_dense()
```

In [89]:

```
op_tri_df = op_tri_df.fillna(0)
```

In [90]:

```
op_tri_df.head()
```

Out[90]:

	add add add	add add add add cmp	add add jmp	add add lea	add add mov	add add or	add add pop	add add push	add add sub	add add xor	xor sub cmp	xor sub mov	xor sub push	xor a
0	0.0000	0.000000	0.002522	0.0	0.005422	0.00267	0.000000	0.004902	0.000662	0.00138	...	0.004508	0.018631	0.005498	0.0000
1	0.0002	0.002123	0.000000	0.0	0.003425	0.00534	0.000000	0.004902	0.000662	0.000000	...	0.000000	0.003388	0.005498	0.0000
2	0.0001	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.003268	0.000000	0.000000	...	0.000000	0.000000	0.021993	0.0058
3	0.0000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000	0.0000
4	0.0001	0.000000	0.002522	0.0	0.000571	0.000000	0.001891	0.000000	0.000000	0.000000	...	0.000000	0.000000	0.000000	0.000000

5 rows × 1001 columns

In [91]:

```
op_tri_df.to_csv('op_trigram_1000.csv')
```

1000 Important Feature Among tetramgram

In [5]:

```
from sklearn.preprocessing import normalize  
  
normalized_opcodetetravect = normalize(opcodetetravect , axis = 0)
```

In [10]:

```
len(asn_4gram)
```

Out[10]:

456976

In [13]:

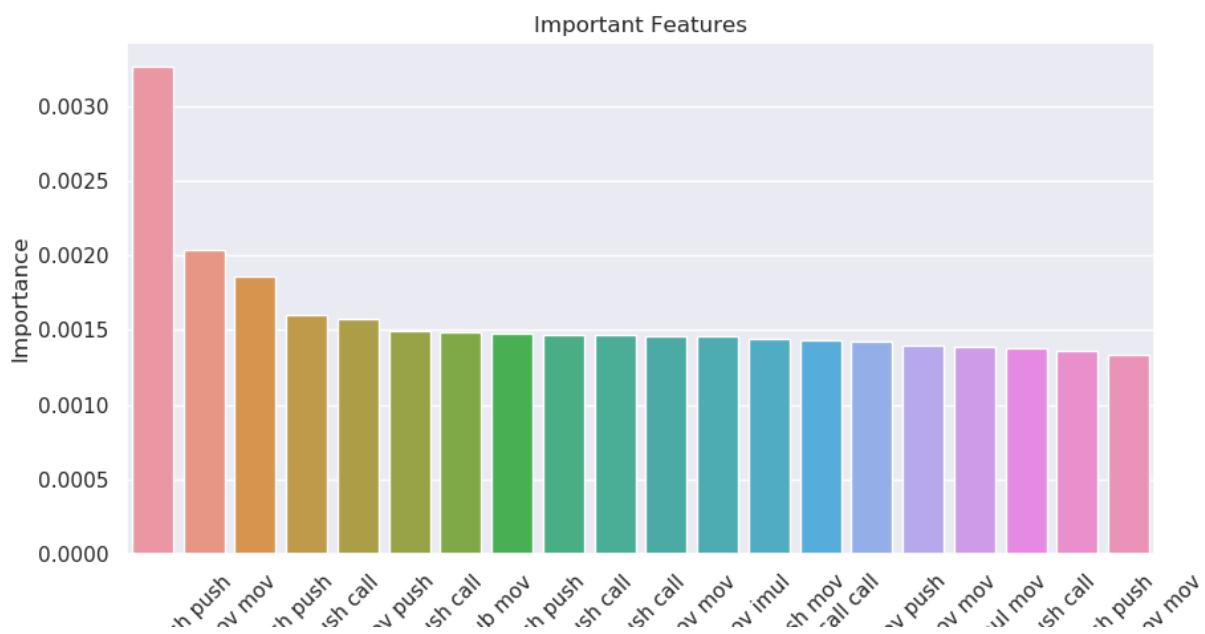
```
normalized_opcodetetravect
```

Out[13]:

```
<10868x456976 sparse matrix of type '<class 'numpy.float64'>'  
with 16604101 stored elements in Compressed Sparse Column format>
```

In [15]:

```
op_tetra_indexes = imp_features(normalized_opcodetetravect, asn_4gram, 1000)
```



In [34]:

```
scipy.sparse.save_npz('normalized_opcodetetravect.npz' , normalized_opcodetetravect)
```

In [16]:

```
op_tetra_df = pd.SparseDataFrame(normalized_opcodetetravect, columns = asn_4gram)
```

In [17]:

```
op_tetra_df = op_tetra_df.loc[:, np.intersect1d(op_tetra_df.columns, np.take(asn_4gram, op_tetra_in
```

```
axes))]
```

In [28]:

```
op_tetra_df['ID'] = result.ID
```

In [29]:

```
op_tetra_df.head()
```

Out[29]:

	add	add	add	add	add	add add	add add	add	add	add	add add	add add	add	add add	xor	xor	xor	xor	xor	xor
0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	push	push	push	push	push	ret
1	NaN	NaN	NaN	NaN	0.003854	NaN	NaN	NaN	NaN	NaN	0.005671	...	NaN	NaN	push	push	push	push	push	ret
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	0.003677	NaN	NaN	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	0.011386	NaN	NaN	NaN	NaN	NaN	NaN	NaN	0.000567	...	0.001965	NaN	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 1001 columns

In [30]:

```
op_tetra_df = op_tetra_df.to_dense()
```

In [31]:

```
op_tetra_df = op_tetra_df.fillna(0)
```

In [32]:

```
#checking if there are any NaN values
op_tetra_df.isnull().values.any()
```

Out[32]:

False

In [33]:

```
op_tetra_df.head()
```

Out[33]:

	add	add	add	add	add	add add	add add	add	add	add	add add	add add	add	add add	xor	xor	xor	xor	xor	xor
0	0.0	0.0	0.000000	0.0	0.000000	0.0	0.0	0.004399	0.0	0.006238	...	0.000000	0.000525	0.0	0.002512	0.0	0.003927			
1	0.0	0.0	0.000000	0.0	0.003854	0.0	0.0	0.000000	0.0	0.005671	...	0.000000	0.000000	0.0	0.007537	0.0	0.000000			
2	0.0	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.000000	...	0.000000	0.003677	0.0	0.000000	0.0	0.000000			
3	0.0	0.0	0.000000	0.0	0.000000	0.0	0.0	0.000000	0.0	0.000000	...	0.000000	0.000000	0.0	0.000000	0.0	0.000000			
4	0.0	0.0	0.011386	0.0	0.000000	0.0	0.0	0.000000	0.0	0.000567	...	0.001965	0.000000	0.0	0.000000	0.0	0.000000			

5 rows × 1001 columns

In [34]:

```
op_tetra_df.to_csv('op_tetragram_1000.csv')
```

```
In [20]:
```

```
op_tetra_df = pd.read_csv('op_tetragram.csv')
```

Extracting top 1000 features from byte_bigram

```
In [92]:
```

```
byte_bigram_vect_normalized = scipy.sparse.load_npz('byte_bigram_vect_normalized.npz')
```

```
In [93]:
```

```
byte_bigram_vect_normalized.shape
```

```
Out[93]:
```

```
(10868, 66049)
```

```
In [100]:
```

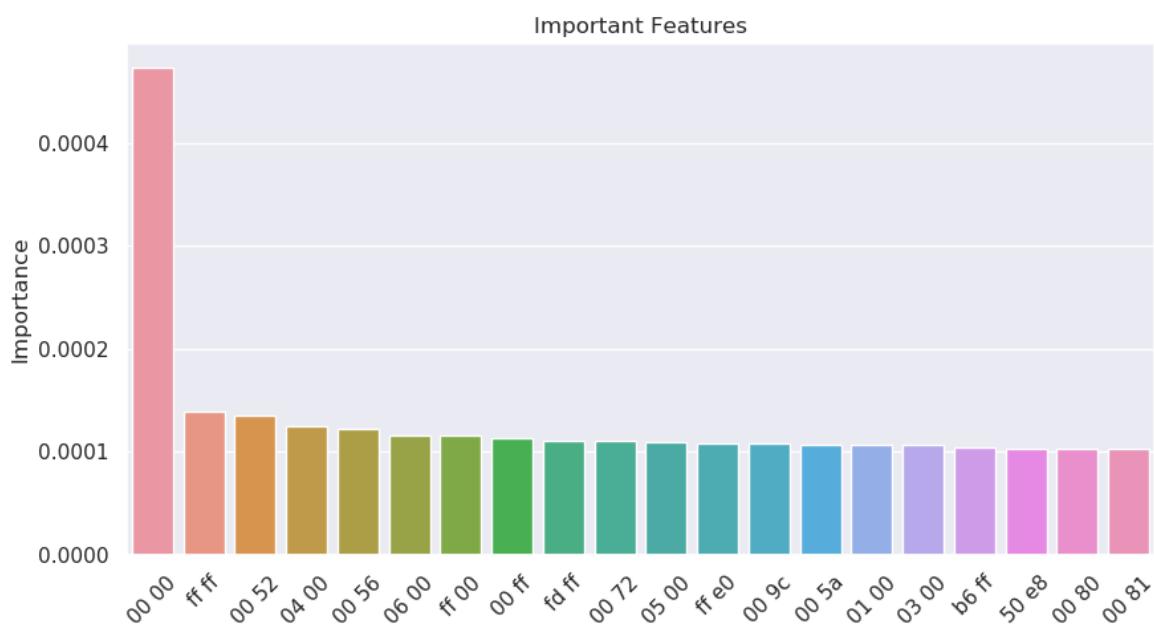
```
len(byte_bigram_vocab)
```

```
Out[100]:
```

```
66049
```

```
In [101]:
```

```
byte_bi_indexes = imp_features(byte_bigram_vect_normalized, byte_bigram_vocab, 1000)
```



```
In [65]:
```

```
np.save('byte_bi_indexes', byte_bi_indexes)
```

```
In [103]:
```

```
len(byte_bi_indexes)
```

```
Out[103]:
```

```
1000
```

In [104]:

```
top_byte_bigrams = np.zeros((10868, 0))
for i in tqdm(byte_bi_idxes):
    sliced = byte_bigram_vect_normalized[:, i].todense()
    top_byte_bigrams = np.hstack([top_byte_bigrams, sliced])
```

100%|██████████| 1000/1000 [00:16<00:00, 62.25it/s]

In [105]:

```
top_byte_bigrams.shape
```

Out[105]:

(10868, 1000)

In [106]:

```
len(byte_bigram_vocab)
```

Out[106]:

66049

In [107]:

```
np.take(byte_bigram_vocab, byte_bi_idxes).shape
```

Out[107]:

(1000,)

In [108]:

```
byte_bi_df = pd.SparseDataFrame(top_byte_bigrams, columns = np.take(byte_bigram_vocab,
byte_bi_idxes))
```

In [109]:

```
byte_bi_df['ID'] = result.ID
```

In [110]:

```
byte_bi_df.head()
```

Out[110]:

	00 00	ff ff	00 52	04 00	00 56	06 00	ff 00	00 ff	fd ff	00 72	...	ed ba	3f 3f	50 2
0	0.045798	0.000767	0.011985	0.014964	0.021535	0.022229	0.008520	0.007967	0.000379	0.008840	...	0.007672	0.002898	0.00063
1	0.003330	0.004335	0.000192	0.004179	0.000544	0.000292	0.011985	0.004612	0.004950	0.035056	...	0.000000	0.000000	0.00031
2	0.002689	0.002120	0.001545	0.006752	0.010777	0.008455	0.016569	0.015702	0.003687	0.001539	...	0.005480	0.000131	0.00268
3	0.001661	0.001221	0.000229	0.001415	0.000881	0.000523	0.003607	0.003308	0.000935	0.001338	...	0.002192	0.001436	0.00000
4	0.002564	0.000068	0.000823	0.000472	0.000112	0.000184	0.000136	0.000289	0.000152	0.000044	...	0.000000	0.000052	0.00086

5 rows × 1001 columns

In [112]:

```
byte_bi_df.isnull().values.any()
```

Out[112]:

```
False
```

```
In [113]:
```

```
byte_bi_df.to_dense().to_csv('byte_bigram_df_1000.csv')
```

```
In [100]:
```

```
result_x['ID'] = result.ID
```

```
In [101]:
```

```
result_x.head()
```

```
Out[101]:
```

	Unnamed: 0	0	1	2	3	4	5	6	7	8	...	esi	eax	...
0	0.000000	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946	0.002638	...	0.025875	0.025744	0.0049
1	0.000092	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984	0.008267	...	0.012316	0.007858	0.0075
2	0.000184	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155	0.008104	...	0.006181	0.000100	0.0037
3	0.000276	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481	0.000959	...	0.000746	0.000301	0.0003
4	0.000368	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229	0.000376	...	0.013875	0.000482	0.0129

5 rows × 309 columns

```
In [106]:
```

```
result_x.to_csv('result_x.csv')
```

```
In [126]:
```

```
normalized_asm.head()
```

```
Out[126]:
```

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	:dword	edx
0	01kcPWA9K2B0xQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084	0.0	0.000072	...	0.001028	0.000343
1	1E93CpP60RHFNiT5Qfvn	0.096045	0.001230	0.0	0.000617	0.000019	0.0	0.000000	0.0	0.000072	...	0.000975	0.000343
2	3ekVow2ajZHbTnBcsDfX	0.096045	0.000627	0.0	0.000300	0.000017	0.0	0.000038	0.0	0.000072	...	0.000630	0.000248
3	3X2nY7iQaPBIWDrAZqJe	0.096045	0.000333	0.0	0.000258	0.000008	0.0	0.000000	0.0	0.000072	...	0.000188	0.000114
4	46OZzdsSKDCFV8h7XWxf	0.096045	0.000590	0.0	0.000353	0.000068	0.0	0.000000	0.0	0.000072	...	0.000135	0.000229

5 rows × 49 columns

Extracting some features has taken too long. Especially Bytebigrams, about 15 hours

My kernel was hanging at times when i was trying to run the models with additional features.
(bytebigrams+opcode+image features)

So all the csv files of the dataframes were saved and the machine learning models are in another notebook.

Please open ms_malware_models.ipynb to check how the models have performed on the features