

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> (<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> (<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>
(<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>)
(<https://www.reviversoft.com/file-extensions/asm>)
 2. .bytes file (the raw data contains the hexadecimal representation of the file's binary content, without the PE header)
- Total train dataset consist of 200GB data out of which 50Gb of data is .bytes files and 150GB of data is .asm files:
- **Lots of Data for a single-box/computer.**
- There are total 10,868 .bytes files and 10,868 asm files total 21,736 files
- There are 9 types of malwares (9 classes) in our give data
- Types of Malware:
 1. Ramnit
 2. Lollipop
 3. Kelihos_ver3
 4. Vundo
 5. Simda
 6. Tracur
 7. Kelihos_ver1
 8. Obfuscator.ACY
 9. Gatak

2.1.2. Example Data Point

.asm file

```

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

```

.bytes file

```

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

```

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

2.2.1. Type of Machine Learning Problem

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

2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss

- Confusion matrix

2.2.3. Machine Learning Objectives and Constraints

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

Constraints:

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

2.3. Train and Test Dataset

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

2.4. Useful blogs, videos and reference papers

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

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

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

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

https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu_plB6ua?dl=0

(https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu_plB6ua?dl=0)

" Cross validation is more trustworthy than domain knowledge."

3. Exploratory Data Analysis

```
In [5]: 1 import warnings
2 warnings.filterwarnings("ignore")
3 import shutil
4 import os
5 import pandas as pd
6 import matplotlib
7 matplotlib.use(u'nbAgg')
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import numpy as np
11 import pickle
12 from sklearn.manifold import TSNE
13 from sklearn import preprocessing
14 import pandas as pd
15 from multiprocessing import Process# this is used for multithreading
16 import multiprocessing
17 import codecs# this is used for file operations
18 import random as r
19 from xgboost import XGBClassifier
20 from sklearn.model_selection import RandomizedSearchCV
21 from sklearn.tree import DecisionTreeClassifier
22 from sklearn.calibration import CalibratedClassifierCV
23 from sklearn.neighbors import KNeighborsClassifier
24 import dask.dataframe as dd
25 import pandas as pd
26 import scipy.sparse as sp
27 from sklearn.metrics import log_loss
28 from sklearn.metrics import confusion_matrix
29 from sklearn.model_selection import train_test_split
30 from sklearn.linear_model import LogisticRegression
31 from sklearn.ensemble import RandomForestClassifier
32 from tqdm import tqdm
33 import scipy.sparse as sp
34 import scipy.sparse
35 import gc
36 import pickle as pkl
37 from datetime import datetime as dt
38
39 from IPython.core.display import display, HTML
40 display(HTML("<style>.container { width:100% !important; }</style>"))
```

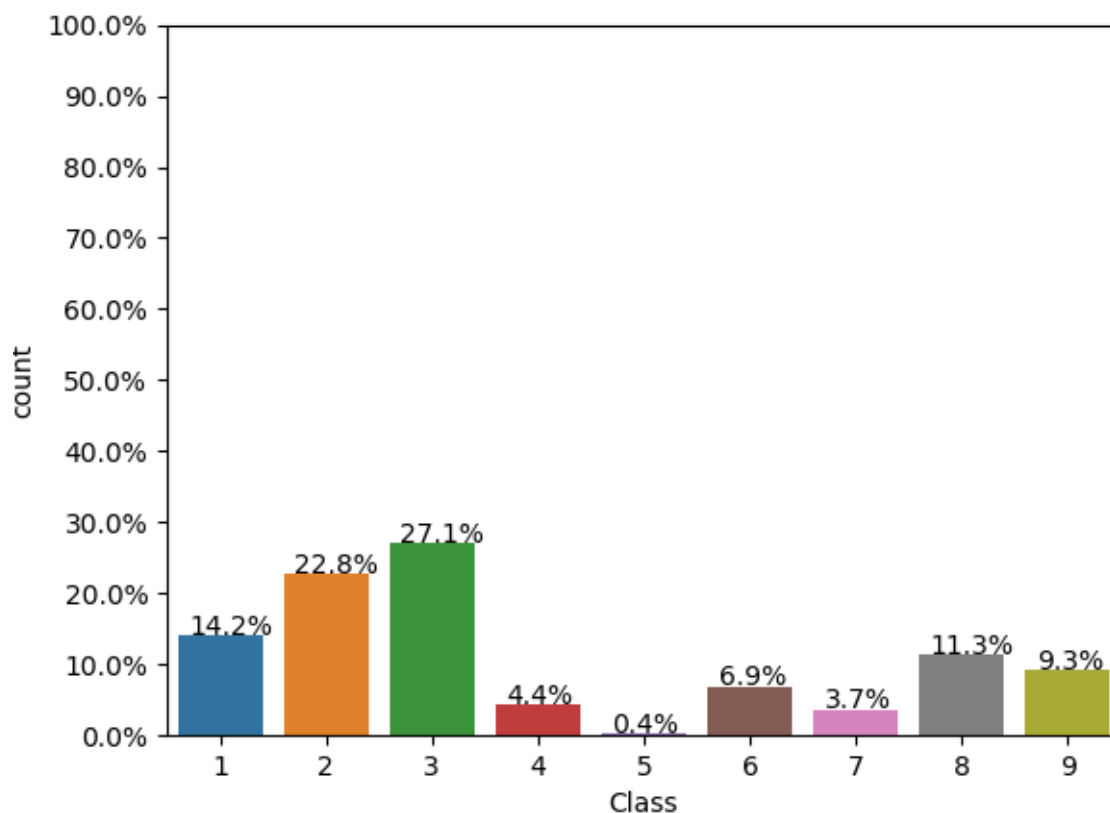
```
In [2]: 1 source = 'train'
2 destination = 'byteFiles'
3
4 #We will check if the folder 'byteFiles' exists if it not there we will creat
5 if not os.path.isdir(destination):
6     os.makedirs(destination)
7
8 # if we have folder called 'train' (train folder contains both .asm files and
9 # for every file that we have in our 'asmFiles' directory we check if it is e
10 # 'byteFiles' folder
11
12 # so by the end of this snippet we will separate all the .byte files and .asm
13 if os.path.isdir(source):
14     os.rename(source, 'asmFiles')
15     source = 'asmFiles'
16     asm_files = os.listdir(source)
17     for file in asm_files:
18         if (file.endswith("bytes")):
19             shutil.move(source+"/"+file, destination)
20
21 print("All byte files moved to 'byteFiles' folder and all asm files moved to
```

All byte files moved to 'byteFiles' folder and all asm files moved to 'asmFiles' folder..

3.1. Distribution of malware classes in whole data set

```
In [43]: 1 Y=pd.read_csv("trainLabels.csv")
2 total = len(Y)*1.
3 ax=sns.countplot(x="Class", data=Y)
4 for p in ax.patches:
5     ax.annotate('{:.1f}%'.format(100*p.get_height()/total), (p.get_x()+0.
6
7     #put 11 ticks (therefore 10 steps), from 0 to the total number of rows in the
8     ax.yaxis.set_ticks(np.linspace(0, total, 11))
9
10    #adjust the ticklabel to the desired format, without changing the position of
11    ax.set_yticklabels(map('{:.1f}%'.format, 100*ax.yaxis.get_majorticklocs()/tot
12    plt.show()
```

<IPython.core.display.Javascript object>



Here we should keep in mind that the dataset is highly imbalanced. We can see that class 1,2,3 are

the ones which has maximum number of datapoints. These 3 classes are closely followed by classes 6,8,9. Class 5 has the lowest number of files. Class 4,6 and 7 also has considerably less number of data points. So this is basically a multi-class problem with an imbalanced dataset.

3.2. Feature extraction

3.2.1 File size of byte files as a feature

In [23]:

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

Out[23]:

	ID	size	Class
0	01azqd4InC7m9JpocGv5	4.148438	9
1	01lsoiSMh5gxyDYTI4CB	5.425781	2
2	01jsnpXSAIgw6aPeDxrU	3.808594	9
3	01kcPWA9K2BOxQeS5Rju	0.562500	1
4	01SuzwMJEIXsK7A8dQbl	0.363281	8

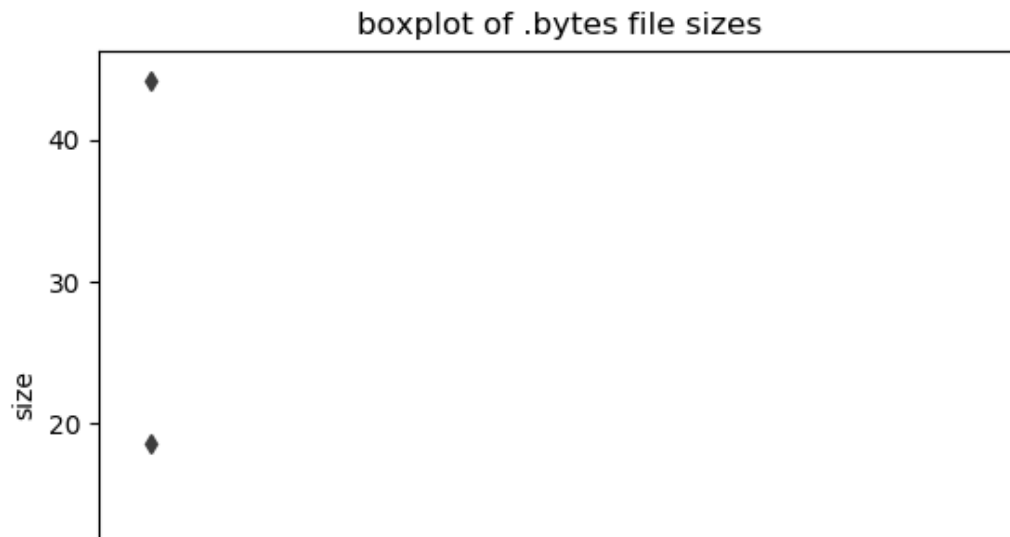
In [24]:

```
1 #Save byte file size dataframe
2 if not os.path.isdir("features"):
3     os.makedirs("features")
4
5 data_size_byte.to_csv("features/data_size_byte.csv")
```

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

```
In [25]: 1 #boxplot of byte files
2 ax = sns.boxplot(x="Class", y="size", data=data_size_byte)
3 plt.title("boxplot of .bytes file sizes")
4 plt.show()
```

<IPython.core.display.Javascript object>



Looking at the box plot carefully, we can see that the size feature is somewhat useful in determining some of the class labels if not all. In the above box plot, class 2 can be clearly separated from other classes (1,3,4,6,7,8) by just using the size feature. Class 3 has one of the lowest file sizes and they can be easily separated by classes. IF not fully, file sizes partially helps us to differentiate between few classes. So file sizes are indeed useful features. We will keep them and make use of them while doing our analysis and also while building our models.

3.2.3 Feature extraction from byte files

```

In [0]: 1 #removal of addres from byte files
2 # contents of .byte files
3 # -----
4 #00401000 56 8D 44 24 08 50 8B F1 E8 1C 1B 00 00 C7 06 08
5 #-----
6 #we remove the starting address 00401000
7
8 files = os.listdir('byteFiles')
9 filenames=[]
10 array=[]
11 for file in files:
12     if(file.endswith("bytes")):
13         file=file.split('.')[0]
14         text_file = open('byteFiles/'+file+".txt", 'w+')
15         with open('byteFiles/'+file+".bytes", "r") as fp:
16             lines=""
17             for line in fp:
18                 a=line.rstrip().split(" ")[1:]
19                 b=' '.join(a)
20                 b=b+"\n"
21                 text_file.write(b)
22             fp.close()
23             os.remove('byteFiles/'+file+".bytes")
24         text_file.close()
25
26 files = os.listdir('byteFiles')
27 filenames2=[]
28 feature_matrix = np.zeros((len(files),257),dtype=int)
29 k=0
30
31
32 #program to convert into bag of words of bytefiles
33 #this is custom-built bag of words this is unigram bag of words
34 byte_feature_file=open('result.csv','w+')
35 byte_feature_file.write("ID,0,1,2,3,4,5,6,7,8,9,0a,0b,0c,0d,0e,0f,10,11,12,13
36 byte_feature_file.write("\n")
37 for file in files:
38     filenames2.append(file)
39     byte_feature_file.write(file+",")
40     if(file.endswith("txt")):
41         with open('byteFiles/'+file,"r") as byte_flie:
42             for lines in byte_flie:
43                 line=line.rstrip().split(" ")
44                 for hex_code in line:
45                     if hex_code=='?':
46                         feature_matrix[k][256]+=1
47                     else:
48                         feature_matrix[k][int(hex_code,16)]+=1
49             byte_flie.close()
50     for i, row in enumerate(feature_matrix[k]):
51         if i!=len(feature_matrix[k])-1:
52             byte_feature_file.write(str(row)+",")
53         else:
54             byte_feature_file.write(str(row))
55     byte_feature_file.write("\n")
56

```

```

57     k += 1
58
59 byte_feature_file.close()

```

```

In [26]: 1 byte_features=pd.read_csv("result.csv").drop(columns=["Unnamed: 0"])
          2 byte_features['ID'] = byte_features['ID'].str.split('.').str[0]
          3 byte_features["size"]=data_size_byte["size"]
          4 byte_features.head()

```

```

Out[26]:

```

		ID	0	1	2	3	4	5	6	7	8	...	f9	
0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...	3101	32	
1	01lsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...	439	2	
2	01jsnpXSAIgw6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	9107	...	2242	28	
3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	1078	...	485	4	
4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	422	...	350	2	

5 rows × 260 columns

```

In [27]: 1 byte_features_with_size = byte_features
          2 byte_features_with_size.to_csv("features/byte_features_with_size.csv",index=N

```

```

In [71]: 1 byte_features_with_size=pd.read_csv("features/byte_features_with_size.csv")
          2 byte_features_with_size.head()

```

```

Out[71]:

```

		ID	0	1	2	3	4	5	6	7	8	...	f9	
0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...	3101	32	
1	01lsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...	439	2	
2	01jsnpXSAIgw6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	9107	...	2242	28	
3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	1078	...	485	4	
4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	422	...	350	2	

5 rows × 260 columns

```

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

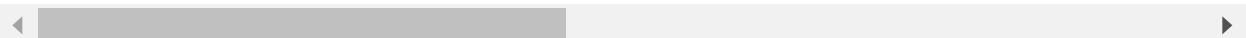
```

```
In [59]: 1 result = normalize(byte_features_with_size)
          2 result.head(5)
```

```
Out[59]:
```

	ID	0	1	2	3	4	5	6
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747
2	01jsnpXSAlgW6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148

5 rows × 260 columns

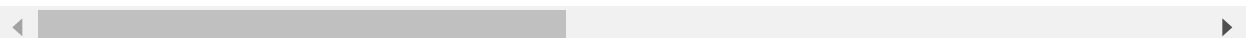


```
In [101]: 1 data_y = result['Class']
           2 result.head()
```

```
Out[101]:
```

	ID	0	1	2	3	4	5	6
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747
2	01jsnpXSAlgW6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148

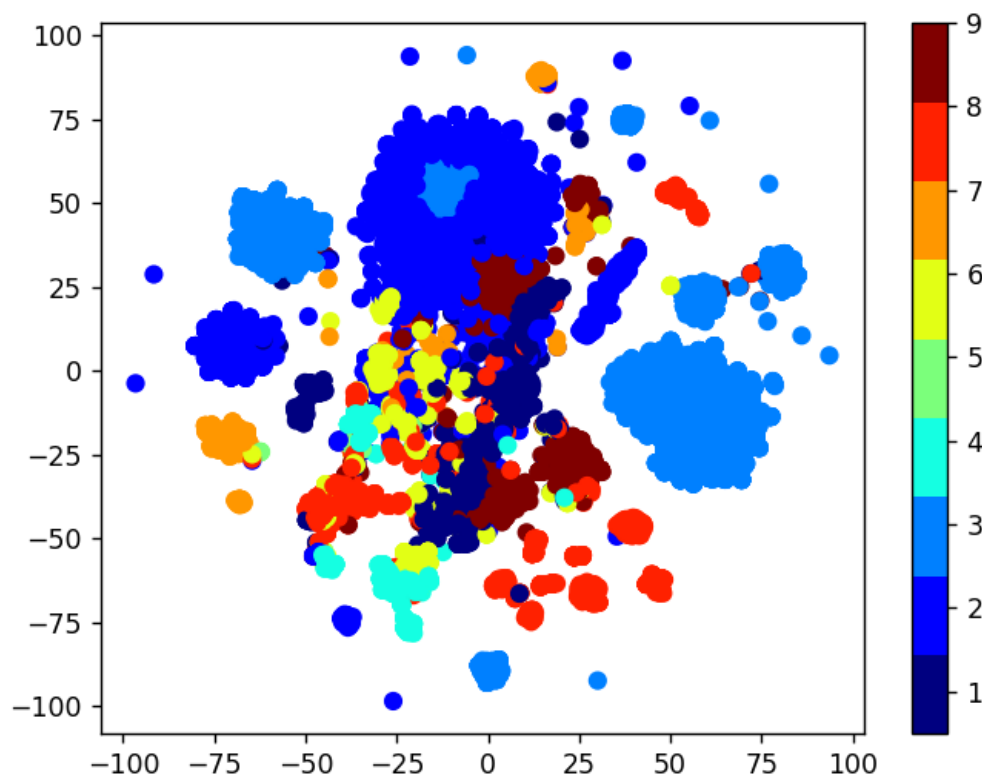
5 rows × 260 columns



3.2.4 Multivariate Analysis Plot TSNE

```
In [0]: 1 #Multivariate analysis on byte files features extracted using unigrams.  
2 def draw_tsne_byte(p):  
3     xtsne=TSNE(perplexity=p)  
4     results=xtsne.fit_transform(result.drop(['ID','Class'], axis=1))  
5     vis_x = results[:, 0]  
6     vis_y = results[:, 1]  
7     plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))  
8     plt.colorbar(ticks=range(10))  
9     plt.clim(0.5, 9)  
10    plt.show()  
11  
12 draw_tsne_byte(30)
```

<IPython.core.display.Javascript object>



At various values of perplexity we can see that the clusters are partially separable. There is some overlapping but in general classes 1,2,3 are well separated from the rest. The red points belonging to cluster 8 are scattered over a larger area, but inspite of this, there is no significant overlapping of class 8 with other clusters. This tells us that the unigrams features extracted from using a custom bag of words approach is pretty useful in separating the data.

Train Test split

```
In [102]: 1 data_y = result['Class']
          2 # split the data into test and train by maintaining same distribution of outp
          3 X_train, X_test, y_train, y_test = train_test_split(result.drop(['ID', 'Class'
          4 # split the train data into train and cross validation by maintaining same di
          5 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_t
```

```
In [103]: 1 print('Number of data points in train data:', X_train.shape[0])
          2 print('Number of data points in test data:', X_test.shape[0])
          3 print('Number of data points in cross validation data:', X_cv.shape[0])
```

Number of data points in train data: 6955

Number of data points in test data: 2174

Number of data points in cross validation data: 1739

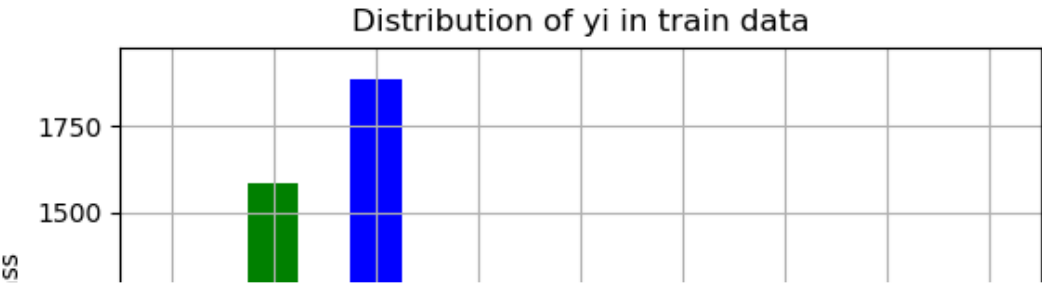
In [107]:

```

1  # it returns a dict, keys as class labels and values as the number of data po
2  train_class_distribution = y_train.value_counts().sortlevel()
3  test_class_distribution = y_test.value_counts().sortlevel()
4  cv_class_distribution = y_cv.value_counts().sortlevel()
5
6  my_colors = ["r","g","b","orange","y","m","c","g","black"]
7  train_class_distribution.plot(kind='bar', color=my_colors)
8  plt.xlabel('Class')
9  plt.ylabel('Data points per Class')
10 plt.title('Distribution of yi in train data')
11 plt.grid()
12 plt.show()
13
14 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.arg
15 # -(train_class_distribution.values): the minus sign will give us in decreasi
16 sorted_yi = np.argsort(-train_class_distribution.values)
17 for i in sorted_yi:
18     print('Number of data points in class', i+1, ':', train_class_distribution
19
20 print('-'*80)
21 test_class_distribution.plot(kind='bar', color=my_colors)
22 plt.xlabel('Class')
23 plt.ylabel('Data points per Class')
24 plt.title('Distribution of yi in test data')
25 plt.grid()
26 plt.show()
27
28 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.arg
29 # -(train_class_distribution.values): the minus sign will give us in decreasi
30 sorted_yi = np.argsort(-test_class_distribution.values)
31 for i in sorted_yi:
32     print('Number of data points in class', i+1, ':', test_class_distribution.
33
34 cv_class_distribution.plot(kind='bar', color=my_colors)
35 plt.xlabel('Class')
36 plt.ylabel('Data points per Class')
37 plt.title('Distribution of yi in cross validation data')
38 plt.grid()
39 plt.show()
40
41 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.arg
42 # -(train_class_distribution.values): the minus sign will give us in decreasi
43 sorted_yi = np.argsort(-train_class_distribution.values)
44 for i in sorted_yi:
45     print('Number of data points in class', i+1, ':', cv_class_distribution.va
46

```

<IPython.core.display.Javascript object>



```

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

```

4. Machine Learning Models

4.1. Machine Learning Models on bytes files

4.1.1. Random Model

```

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

```

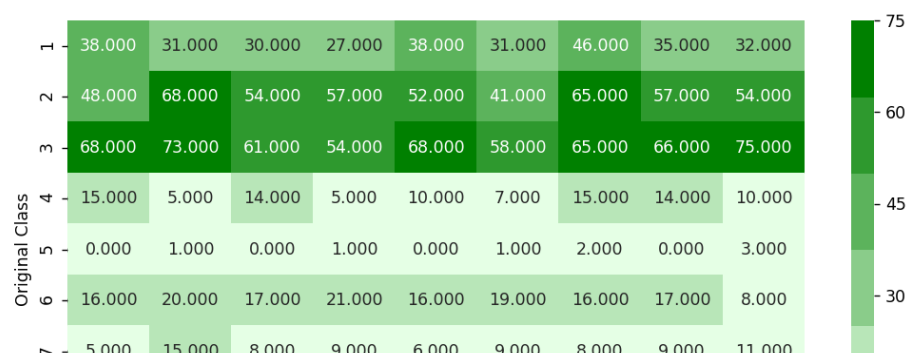
Log loss on Cross Validation Data using Random Model 2.45615644965

Log loss on Test Data using Random Model 2.48503905509

Number of misclassified points 88.5004599816

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

<IPython.core.display.Javascript object>



4.1.2. K Nearest Neighbour Classification

```

In [0]: 1 # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/
2 # -----
3 # default parameter
4 # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', le
5 # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)
6
7 # methods of
8 # fit(X, y) : Fit the model using X as training data and y as target values
9 # predict(X):Predict the class labels for the provided data
10 # predict_proba(X):Return probability estimates for the test data X.
11 #-----
12 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online
13 #-----
14
15
16 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stab
17 # -----
18 # default paramters
19 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sig
20 #
21 # some of the methods of CalibratedClassifierCV()
22 # fit(X, y[, sample_weight]) Fit the calibrated model
23 # get_params([deep]) Get parameters for this estimator.
24 # predict(X) Predict the target of new samples.
25 # predict_proba(X) Posterior probabilities of classification
26 #-----
27 # video link:
28 #-----
29
30 alpha = [x for x in range(1, 15, 2)]
31 cv_log_error_array=[]
32 for i in alpha:
33     k_cfl=KNeighborsClassifier(n_neighbors=i)
34     k_cfl.fit(X_train,y_train)
35     sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
36     sig_clf.fit(X_train, y_train)
37     predict_y = sig_clf.predict_proba(X_cv)
38     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cfl.classes_
39
40 for i in range(len(cv_log_error_array)):
41     print ('log_loss for k = ',alpha[i],'is',cv_log_error_array[i])
42
43 best_alpha = np.argmin(cv_log_error_array)
44
45 fig, ax = plt.subplots()
46 ax.plot(alpha, cv_log_error_array,c='g')
47 for i, txt in enumerate(np.round(cv_log_error_array,3)):
48     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
49 plt.grid()
50 plt.title("Cross Validation Error for each alpha")
51 plt.xlabel("Alpha i's")
52 plt.ylabel("Error measure")
53 plt.show()
54
55 k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
56 k_cfl.fit(X_train,y_train)

```

```

57 sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
58 sig_clf.fit(X_train, y_train)
59
60 predict_y = sig_clf.predict_proba(X_train)
61 print('For values of best alpha = ', alpha[best_alpha], "The train log loss")
62 predict_y = sig_clf.predict_proba(X_cv)
63 print('For values of best alpha = ', alpha[best_alpha], "The cross validation")
64 predict_y = sig_clf.predict_proba(X_test)
65 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is")
66 plot_confusion_matrix(y_test, sig_clf.predict(X_test))

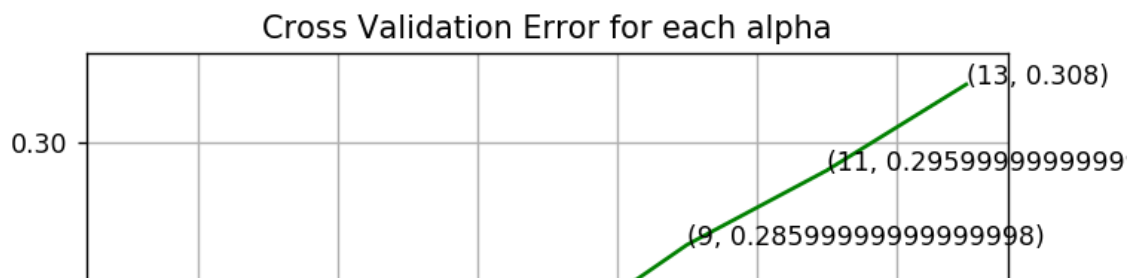
```

```

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

```

<IPython.core.display.Javascript object>



4.1.3. Logistic Regression

```

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

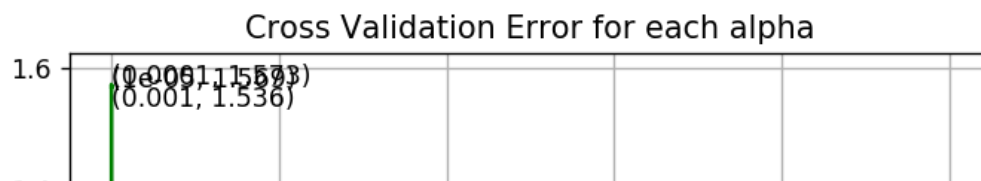
```

log_loss for c = 1e-05 is 1.56916911178

log_loss for c = 0.0001 is 1.57336384417

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

<IPython.core.display.Javascript object>



4.1.4. Random Forest Classifier


```

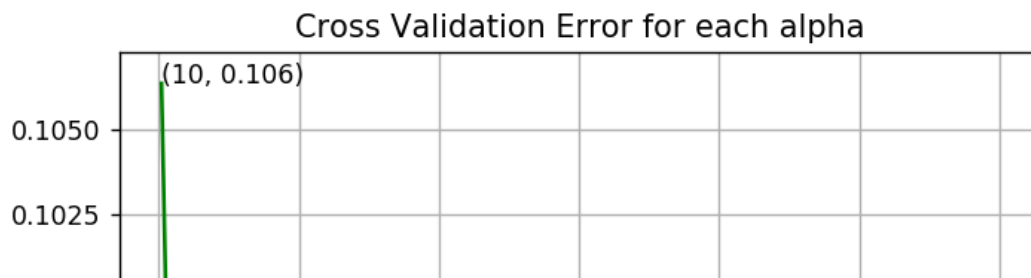
In [0]: 1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini',
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given trai
10 # predict(X) Perform classification on samples in X.
11 # predict_proba (X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online
19 # -----
20
21 alpha=[10,50,100,500,1000,2000,3000]
22 cv_log_error_array=[]
23 train_log_error_array=[]
24 from sklearn.ensemble import RandomForestClassifier
25 for i in alpha:
26     r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
27     r_cfl.fit(X_train,y_train)
28     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
29     sig_clf.fit(X_train, y_train)
30     predict_y = sig_clf.predict_proba(X_cv)
31     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_
32
33 for i in range(len(cv_log_error_array)):
34     print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
35
36
37 best_alpha = np.argmin(cv_log_error_array)
38
39 fig, ax = plt.subplots()
40 ax.plot(alpha, cv_log_error_array,c='g')
41 for i, txt in enumerate(np.round(cv_log_error_array,3)):
42     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
43 plt.grid()
44 plt.title("Cross Validation Error for each alpha")
45 plt.xlabel("Alpha i's")
46 plt.ylabel("Error measure")
47 plt.show()
48
49
50 r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n
51 r_cfl.fit(X_train,y_train)
52 sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
53 sig_clf.fit(X_train, y_train)
54
55 predict_y = sig_clf.predict_proba(X_train)
56 print('For values of best alpha = ', alpha[best_alpha], "The train log loss i

```

```
57 predict_y = sig_clf.predict_proba(X_cv)
58 print('For values of best alpha = ', alpha[best_alpha], "The cross validation
59 predict_y = sig_clf.predict_proba(X_test)
60 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
61 plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```

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

<IPython.core.display.Javascript object>



4.1.5. XgBoost Classification

```

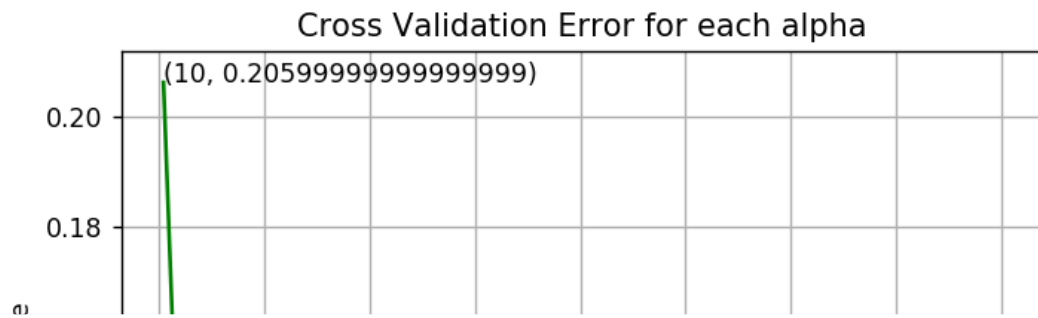
In [0]: 1 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
2
3 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en
4 # -----
5 # default paramters
6 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100,
7 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
8 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_lambda=1,
9 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None)
10
11 # some of methods of RandomForestRegressor()
12 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=0)
13 # get_params([deep]) Get parameters for this estimator.
14 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: If ntree_limit is not None,
15 # get_score(importance_type='weight') -> get the feature importance
16 # -----
17 # video link1: https://www.appliedaicourse.com/course/applied-ai-course-online
18 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online
19 # -----
20
21 alpha=[10,50,100,500,1000,2000]
22 cv_log_error_array=[]
23 for i in alpha:
24     x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
25     x_cfl.fit(X_train,y_train)
26     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
27     sig_clf.fit(X_train, y_train)
28     predict_y = sig_clf.predict_proba(X_cv)
29     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_))
30
31 for i in range(len(cv_log_error_array)):
32     print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
33
34
35 best_alpha = np.argmin(cv_log_error_array)
36
37 fig, ax = plt.subplots()
38 ax.plot(alpha, cv_log_error_array,c='g')
39 for i, txt in enumerate(np.round(cv_log_error_array,3)):
40     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
41 plt.grid()
42 plt.title("Cross Validation Error for each alpha")
43 plt.xlabel("Alpha i's")
44 plt.ylabel("Error measure")
45 plt.show()
46
47 x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
48 x_cfl.fit(X_train,y_train)
49 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
50 sig_clf.fit(X_train, y_train)
51
52 predict_y = sig_clf.predict_proba(X_train)
53 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is", log_loss(y_train, predict_y, labels=x_cfl.classes_))
54 predict_y = sig_clf.predict_proba(X_cv)
55 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is", log_loss(y_cv, predict_y, labels=x_cfl.classes_))
56 predict_y = sig_clf.predict_proba(X_test)

```

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

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

<IPython.core.display.Javascript object>



4.1.5. XgBoost Classification with best hyper parameters using RandomSearch

```
In [0]: 1 # https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-in-machine-learning/
2 x_cfl=XGBClassifier()
3
4 prams={
5     'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
6     'n_estimators':[100,200,500,1000,2000],
7     'max_depth':[3,5,10],
8     'colsample_bytree':[0.1,0.3,0.5,1],
9     'subsample':[0.1,0.3,0.5,1]
10 }
11 random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jobs=10)
12 random_cfl1.fit(X_train,y_train)
```

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

```
[Parallel(n_jobs=-1)]: Done    2 tasks      | elapsed:    26.5s
[Parallel(n_jobs=-1)]: Done    9 tasks      | elapsed:    5.8min
[Parallel(n_jobs=-1)]: Done   19 out of   30 | elapsed:    9.3min remaining:    5.4min
[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[75]: 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]: 1 print (random_cfl1.best_params_)

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

```

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

```

train loss 0.022540976086

cv loss 0.0928710624158

test loss 0.0782688587098

4.2 Modeling with .asm files

There are 10868 files of asm

All the files make up about 150 GB

The asm files contains :

1. Address
2. Segments
3. Opcodes
4. Registers
5. function calls
6. APIs

With the help of parallel processing we extracted all the features. In parallel we can use all the cores that are present in our computer.

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

We read the top solutions and handpicked the features from those papers/v
ideos/blogs.

Refer:<https://www.kaggle.com/c/malware-classification/discussion>

4.2.1 Feature extraction from asm files

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

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

In [0]:

```

1  #http://flint.cs.yale.edu/cs421/papers/x86-asm/asm.html
2  #https://en.wikipedia.org/wiki/X86_instruction_listings
3  #https://en.wikipedia.org/wiki/Data_segment
4
5  """The x86 instruction set refers to the set of instructions that x86-compati
6  executed on the processor. The x86 instruction set has been extended several
7
8  """In computing, a data segment (often denoted .data) is a portion of an obje
9  global variables and static local variables. The size of this segment is dete
10
11
12 def firstprocess():
13     #The prefixes tells about the segments that are present in the asm files
14     #There are 450 segments(approx) present in all asm files.
15     #this prefixes are best segments that gives us best values.
16     #https://en.wikipedia.org/wiki/Data_segment
17
18     prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata
19     #this are opcodes that are used to get best results
20     #https://en.wikipedia.org/wiki/X86_instruction_listings
21
22     opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
23     #best keywords that are taken from different blogs
24     keywords = ['.dll', 'std:', ':dword']
25     #Below taken registers are general purpose registers and special register
26     #All the registers which are taken are best
27     registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
28     file1=open("output/asmsmallfile.txt", "w+")
29     files = os.listdir('first')
30     for f in files:
31         #filling the values with zeros into the arrays
32         prefixescount=np.zeros(len(prefixes),dtype=int)
33         opcodescount=np.zeros(len(opcodes),dtype=int)
34         keywordcount=np.zeros(len(keywords),dtype=int)
35         registerscount=np.zeros(len(registers),dtype=int)
36         features=[]
37         f2=f.split('.')[0] #Contains the file names
38         file1.write(f2+",")
39
40         # https://docs.python.org/3/library/codecs.html#codecs.ignore_errors
41         # https://docs.python.org/3/library/codecs.html#codecs.Codec.encode
42         with codecs.open('first/'+f,encoding='cp1252',errors='replace') as f:
43             for lines in fli:
44                 # https://www.tutorialspoint.com/python3/string_rstrip.htm
45                 line=lines.rstrip().split()
46                 l=line[0]
47                 #counting the prefixes in each and every line
48                 for i in range(len(prefixes)):
49                     if prefixes[i] in line[0]:
50                         prefixescount[i]+=1
51                 line=line[1:]
52                 #counting the opcodes in each and every line
53                 for i in range(len(opcodes)):
54                     if any(opcodes[i]==li for li in line):
55                         features.append(opcodes[i])
56                         opcodescount[i]+=1

```



```

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

```

```

171         file1.write(str(register)+"")
172     for key in keywordcount:
173         file1.write(str(key)+"")
174     file1.write("\n")
175 file1.close()
176
177
178 def fourthprocess():
179     prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:']
180     opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'su
181     keywords = ['.dll', 'std::', ':dword']
182     registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
183     file1=open("output/hugeasmfile.txt", "w+")
184     files = os.listdir('fourth/')
185     for f in files:
186         prefixescount=np.zeros(len(prefixes),dtype=int)
187         opcodescount=np.zeros(len(opcodes),dtype=int)
188         keywordcount=np.zeros(len(keywords),dtype=int)
189         registerscount=np.zeros(len(registers),dtype=int)
190         features=[]
191         f2=f.split('.')[0]
192         file1.write(f2+"")
193
194         with codecs.open('fourth/'+f,encoding='cp1252',errors='replace') as f:
195             for lines in fli:
196                 line=lines.rstrip().split()
197                 l=line[0]
198                 for i in range(len(prefixes)):
199                     if prefixes[i] in line[0]:
200                         prefixescount[i]+=1
201                 line=line[1:]
202                 for i in range(len(opcodes)):
203                     if any(opcodes[i]==li for li in line):
204                         features.append(opcodes[i])
205                         opcodescount[i]+=1
206                 for i in range(len(registers)):
207                     for li in line:
208                         if registers[i] in li and ('text' in l or 'CODE' in l):
209                             registerscount[i]+=1
210                 for i in range(len(keywords)):
211                     for li in line:
212                         if keywords[i] in li:
213                             keywordcount[i]+=1
214         for prefix in prefixescount:
215             file1.write(str(prefix)+"")
216         for opcode in opcodescount:
217             file1.write(str(opcode)+"")
218         for register in registerscount:
219             file1.write(str(register)+"")
220         for key in keywordcount:
221             file1.write(str(key)+"")
222         file1.write("\n")
223     file1.close()
224
225
226 def fifthprocess():
227     prefixes = ['HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:']

```

```

228 opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sc
229 keywords = ['.dll', 'std::', ':dword']
230 registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
231 file1=open("output/trainasmfile.txt", "w+")
232 files = os.listdir('fifth/')
233 for f in files:
234     prefixescount=np.zeros(len(prefixes), dtype=int)
235     opcodescount=np.zeros(len(opcodes), dtype=int)
236     keywordcount=np.zeros(len(keywords), dtype=int)
237     registerscount=np.zeros(len(registers), dtype=int)
238     features=[]
239     f2=f.split('.')[0]
240     file1.write(f2+",")
241
242     with codecs.open('fifth/'+f, encoding='cp1252', errors='replace') as f:
243         for lines in fli:
244             line=lines.rstrip().split()
245             l=line[0]
246             for i in range(len(prefixes)):
247                 if prefixes[i] in line[0]:
248                     prefixescount[i]+=1
249             line=line[1:]
250             for i in range(len(opcodes)):
251                 if any(opcodes[i]==li for li in line):
252                     features.append(opcodes[i])
253                     opcodescount[i]+=1
254             for i in range(len(registers)):
255                 for li in line:
256                     if registers[i] in li and ('text' in l or 'CODE' in l):
257                         registerscount[i]+=1
258             for i in range(len(keywords)):
259                 for li in line:
260                     if keywords[i] in li:
261                         keywordcount[i]+=1
262         for prefix in prefixescount:
263             file1.write(str(prefix)+",")
264         for opcode in opcodescount:
265             file1.write(str(opcode)+",")
266         for register in registerscount:
267             file1.write(str(register)+",")
268         for key in keywordcount:
269             file1.write(str(key)+",")
270         file1.write("\n")
271     file1.close()
272
273
274 def main():
275     #the below code is used for multiprocessing
276     #the number of process depends upon the number of cores present System
277     #process is used to call multiprocessing
278     manager=multiprocessing.Manager()
279     p1=Process(target=firstprocess)
280     p2=Process(target=secondprocess)
281     p3=Process(target=thirdprocess)
282     p4=Process(target=fourthprocess)
283     p5=Process(target=fifthprocess)
284     #p1.start() is used to start the thread execution

```

```

285     p1.start()
286     p2.start()
287     p3.start()
288     p4.start()
289     p5.start()
290     #After completion all the threads are joined
291     p1.join()
292     p2.join()
293     p3.join()
294     p4.join()
295     p5.join()
296
297     if __name__=="__main__":
298         main()

```

In [66]:

```

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

```

Out[66]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...
0	01kcPWA9K2BOxQeS5Rju	19	744	0	127	57	0	323	0	3	...
1	1E93CpP60RHFNiT5Qfvn	17	838	0	103	49	0	0	0	3	...
2	3ekVow2ajZHbTnBcsDfX	17	427	0	50	43	0	145	0	3	...
3	3X2nY7iQaPBIWDrAZqJe	17	227	0	43	19	0	0	0	3	...
4	46OZZdsSKDCFV8h7XWxf	17	402	0	59	170	0	0	0	3	...

5 rows × 53 columns



4.2.1.1 Files sizes of each .asm file

```

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

```

100%|██████████| 10868/10868 [00:48<00:00, 223.06it/s]

```

Out[23]:

```

	ID	size	Class
0	01azqd4InC7m9JpocGv5	56.229886	9
1	01lsoiSMh5gxyDYTI4CB	13.999378	2
2	01jsnpXSAlgW6aPeDxrU	8.507785	9
3	01kcPWA9K2BOxQeS5Rju	0.078190	1
4	01SuzwMJEIXsK7A8dQbl	0.996723	8

```

In [27]: 1 #Save ASM file size dataframe
2 if not os.path.isdir("features"):
3     os.makedirs("features")
4
5 asm_size_byte.to_csv("features/asm_size_df.csv", index=False)

```

In [34]:

```

1 #Adding the file size feature to previous extracted features
2 asm_features_with_size=pd.merge(asm_size_byte,result_asm.drop(columns=["Class
3 asm_features_with_size.head()
```

Out[34]:

		ID	size	Class	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rd
0	01azqd4lnC7m9JpocGv5	56.229886	9	18	22430	0	1158	1366754	0	1	
1	01IsoiSMh5gxyDYTI4CB	13.999378	2	0	109939	0	616	24568	0	26	
2	01jsnpXSAlgW6aPeDxrU	8.507785	9	18	68883	0	304	662	0	1	
3	01kcPWA9K2BOxQeS5Rju	0.078190	1	19	744	0	127	57	0		
4	01SuzwMJEIXsK7A8dQbl	0.996723	8	18	10368	0	206	4595	92		

5 rows × 54 columns



In [35]:

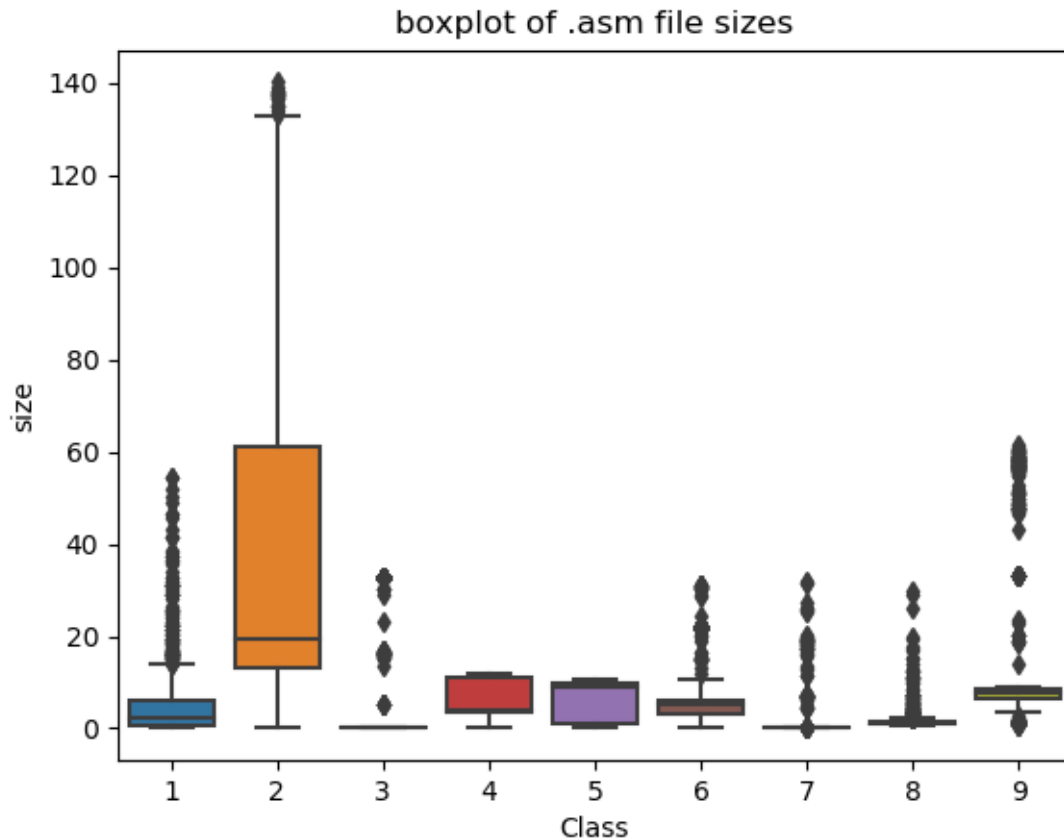
```

1 #Save the ASM Dataframe
2 asm_features_with_size.to_csv("features/asm_features_with_size.csv",index=Non
```

4.2.1.2 Distribution of .asm file sizes

```
In [36]: 1 #boxplot of asm files
2 ax = sns.boxplot(x="Class", y="size", data=asm_size_byte)
3 plt.title("boxplot of .asm file sizes")
4 plt.show()
```

<IPython.core.display.Javascript object>



By looking at the distribution of the ASM file sizes, we can see that class 2 type of malwares can be easily separated from the rest the classes. Class 2 type malwares has the highest spread of file sizes. Among the rest of the classes 7,8 have the least file sizes. But they cannot be strictly separated from the file sizes alone.


```
In [69]: 1 # add the file size feature to previous extracted features
2 result_asm=asm_features_with_size
3 print(result_asm.shape)
4 print(asm_size_byte.shape)
5 result_asm.head()
```

```
(10868, 54)
(10868, 3)
```

```
Out[69]:
```

	ID	size	Class	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rd
0	01azqd4InC7m9JpocGv5	56.229886	9	18	22430	0	1158	1366754	0	1
1	01lsoiSMh5gxyDyTI4CB	13.999378	2	0	109939	0	616	24568	0	26
2	01jsnpXSAlgW6aPeDxrU	8.507785	9	18	68883	0	304	662	0	1
3	01kcPWA9K2BOxQeS5Rju	0.078190	1	19	744	0	127	57	0	
4	01SuzwMJEIXsK7A8dQbl	0.996723	8	18	10368	0	206	4595	92	

5 rows × 54 columns

```
In [78]: 1 byte_features_with_size.head()
```

```
Out[78]:
```

	ID	0	1	2	3	4	5	6	7	8	...	f9
0	01azqd4InC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...	3101
1	01lsoiSMh5gxyDyTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...	439
2	01jsnpXSAlgW6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	9107	...	2242
3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	1078	...	485
4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	422	...	350

5 rows × 260 columns

```
In [70]: 1 asm_features_with_size.head(3)
```

```
Out[70]:
```

	ID	size	Class	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rd
0	01azqd4InC7m9JpocGv5	56.229886	9	18	22430	0	1158	1366754	0	17
1	01lsoiSMh5gxyDyTI4CB	13.999378	2	0	109939	0	616	24568	0	264
2	01jsnpXSAlgW6aPeDxrU	8.507785	9	18	68883	0	304	662	0	10

3 rows × 54 columns

```
In [85]: 1 #Combined the byte and unigrams features along with their sizes and save it i
2 labels=byte_features_with_size["Class"] #Get class labels and ID from byte un
3 with open('features/class_labels.pkl', 'wb') as file:
4     pickle.dump(labels, file)
5
6 combined_unigrams = pd.merge(byte_features_with_size.drop(columns=["Class"]),
7 combined_unigrams=combined_unigrams.drop(columns=["ID"])
8 combined_unigrams.to_csv("features/combined_asm_byte_unigram.csv",index=None)
9 combined_unigrams.head()
```

```
Out[85]:
```

	0	1	2	3	4	5	6	7	8	9	...	:dword	edx	esi	eax
0	601905	3905	2816	3832	3345	3242	3650	3201	2965	3205	...	4371	808	2290	1281
1	39755	8337	7249	7186	8663	6844	8420	7589	9291	358	...	1446	260	1090	391
2	93506	9542	2568	2438	8925	9330	9007	2342	9107	2457	...	903	5	547	5
3	21091	1213	726	817	1257	625	550	523	1078	473	...	137	18	66	15
4	19764	710	302	433	559	410	262	249	422	223	...	1220	18	1228	24

5 rows × 307 columns

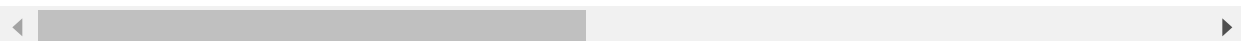


```
In [91]: 1 #We normalize the data each column
2 result_asm = normalize(result_asm)
3 result_asm.head()
```

```
Out[91]:
```

	ID	size	Class	HEADER:	.text:	.Pav:	.idata:	.data:	.b
0	01azqd4lnC7m9JpocGv5	0.400910	9	0.101695	0.032927	0.0	0.006937	0.542847	0.0000
1	01lsoiSMh5gxyDYTI4CB	0.099719	2	0.000000	0.161391	0.0	0.003690	0.009758	0.0000
2	01jsnpXSAIgw6aPeDxrU	0.060553	9	0.101695	0.101121	0.0	0.001821	0.000263	0.0000
3	01kcPWA9K2BOxQeS5Rju	0.000432	1	0.107345	0.001092	0.0	0.000761	0.000023	0.0000
4	01SuzwMJEIXsK7A8dQbl	0.006983	8	0.101695	0.015220	0.0	0.001234	0.001825	0.0120

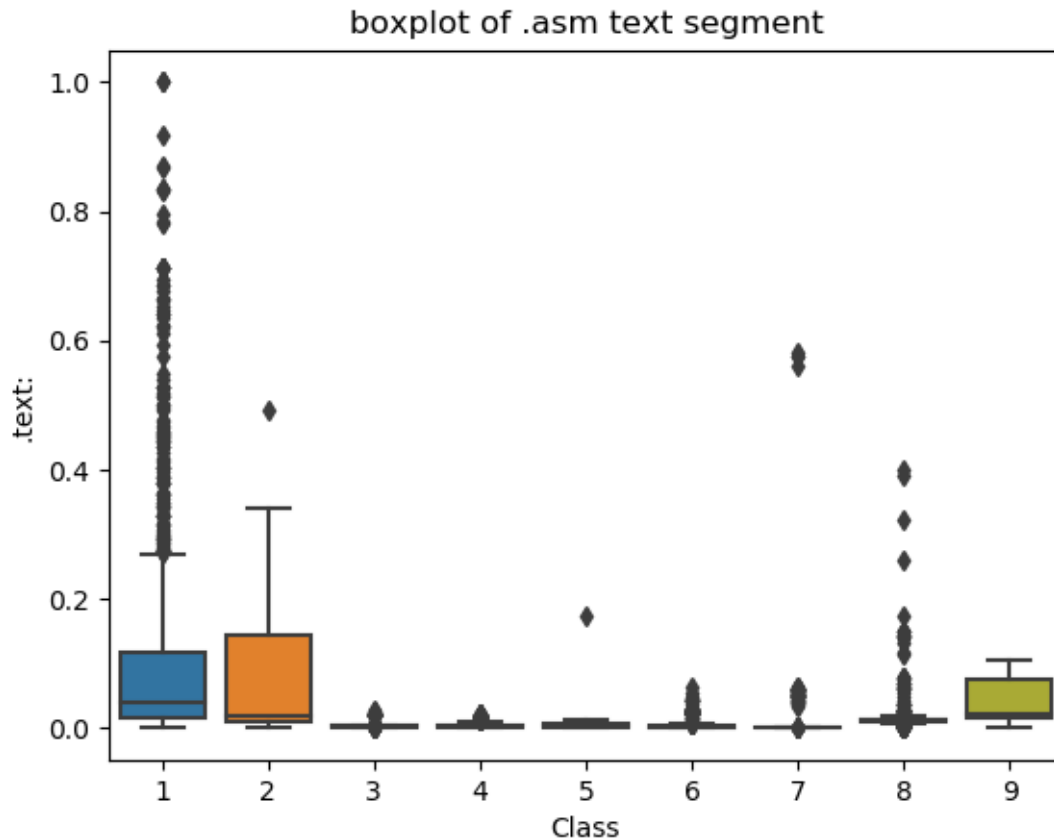
5 rows × 54 columns



4.2.2 Univariate analysis on asm file features

```
In [92]: 1 #Distribution of text keyword segment
2 ax = sns.boxplot(x="Class", y=".text:", data=result_asm)
3 plt.title("boxplot of .asm text segment")
4 plt.show()
```

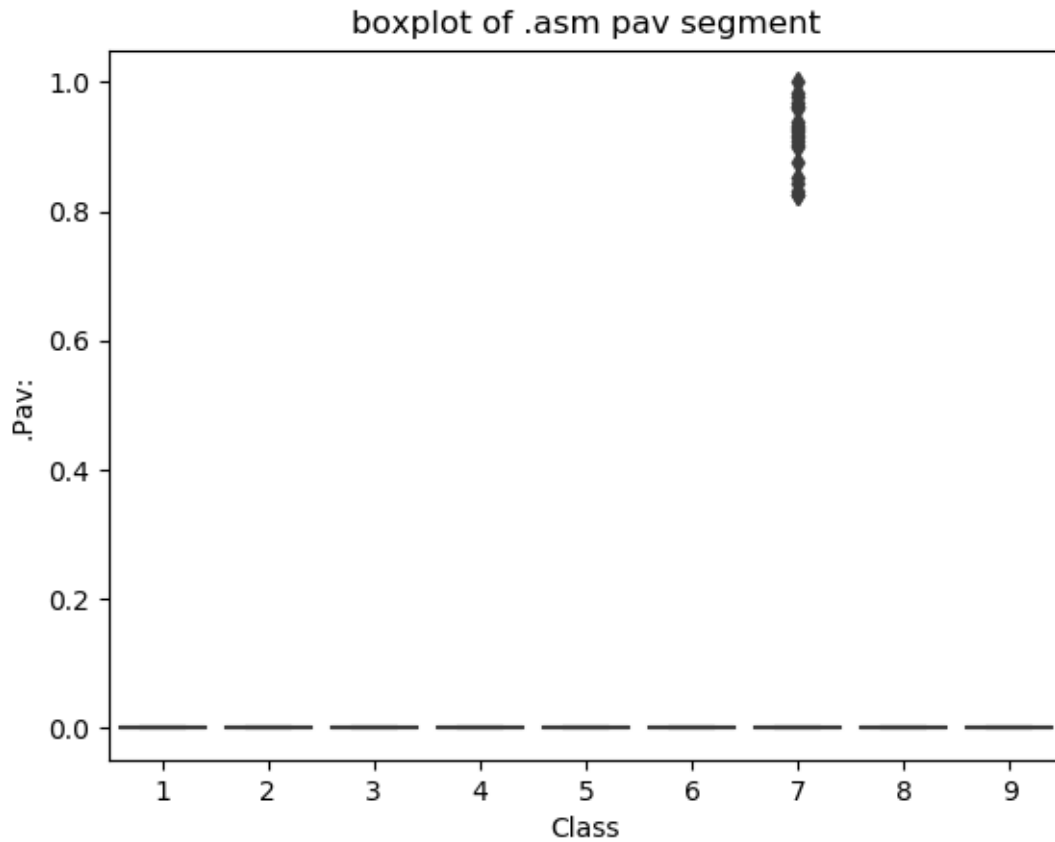
<IPython.core.display.Javascript object>



This plot is between the text segment features vs the class labels. Here we can see that classes 1,2 and 9 can be easily separated using this feature. The rest of the classes aren't easily separated.

```
In [93]: 1 #Distribution of .Pav segments
2 ax = sns.boxplot(x="Class", y=".Pav:", data=result_asm)
3 plt.title("boxplot of .asm pav segment")
4 plt.show()
```

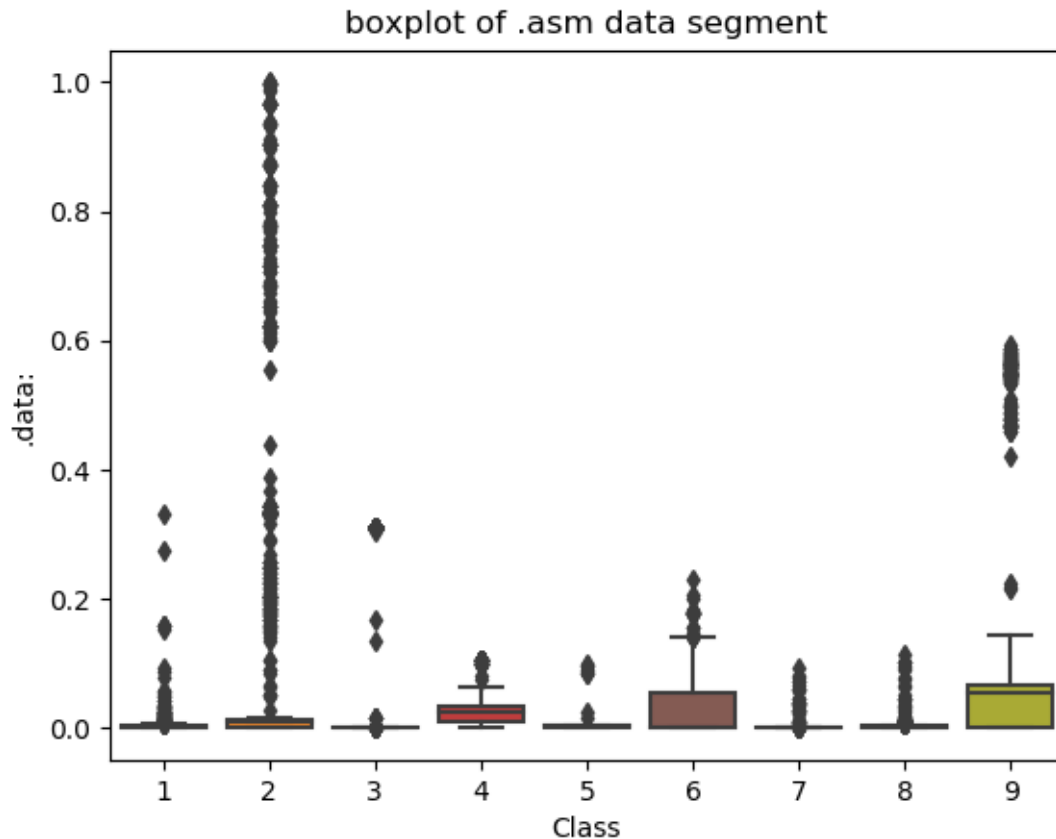
<IPython.core.display.Javascript object>



Here none of the classes can be well separated. However, an important point to note is class 7 type of malware are having the highest spread of the .pav segments.

```
In [94]: 1 #Distribution of .data segments
2 ax = sns.boxplot(x="Class", y=".data:", data=result_asm)
3 plt.title("boxplot of .asm data segment")
4 plt.show()
```

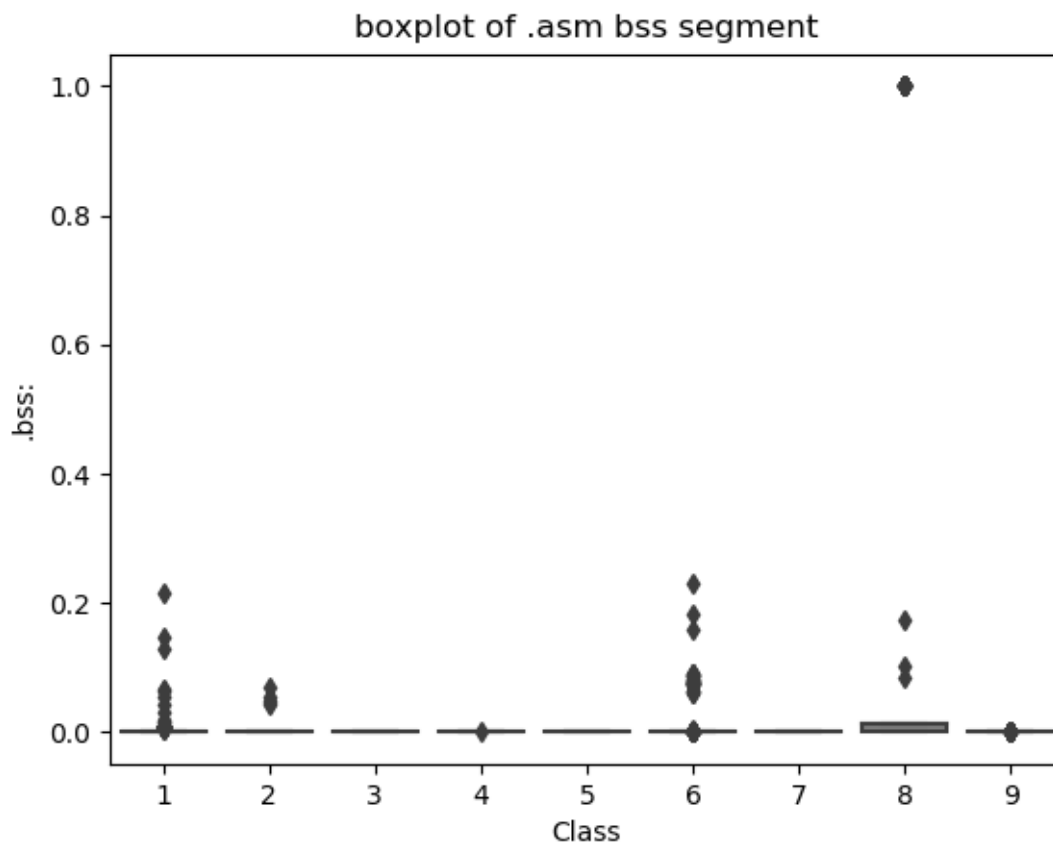
<IPython.core.display.Javascript object>



In this plot, class 6 and 9 can be well separated from the rest of the classes. For the rest of the classes the spread is extremely small to come to a conclusion solely based on this feature alone.

```
In [95]: 1 #Distribution of .bss segments
2 ax = sns.boxplot(x="Class", y=".bss:", data=result_asm)
3 plt.title("boxplot of .asm bss segment")
4 plt.show()
```

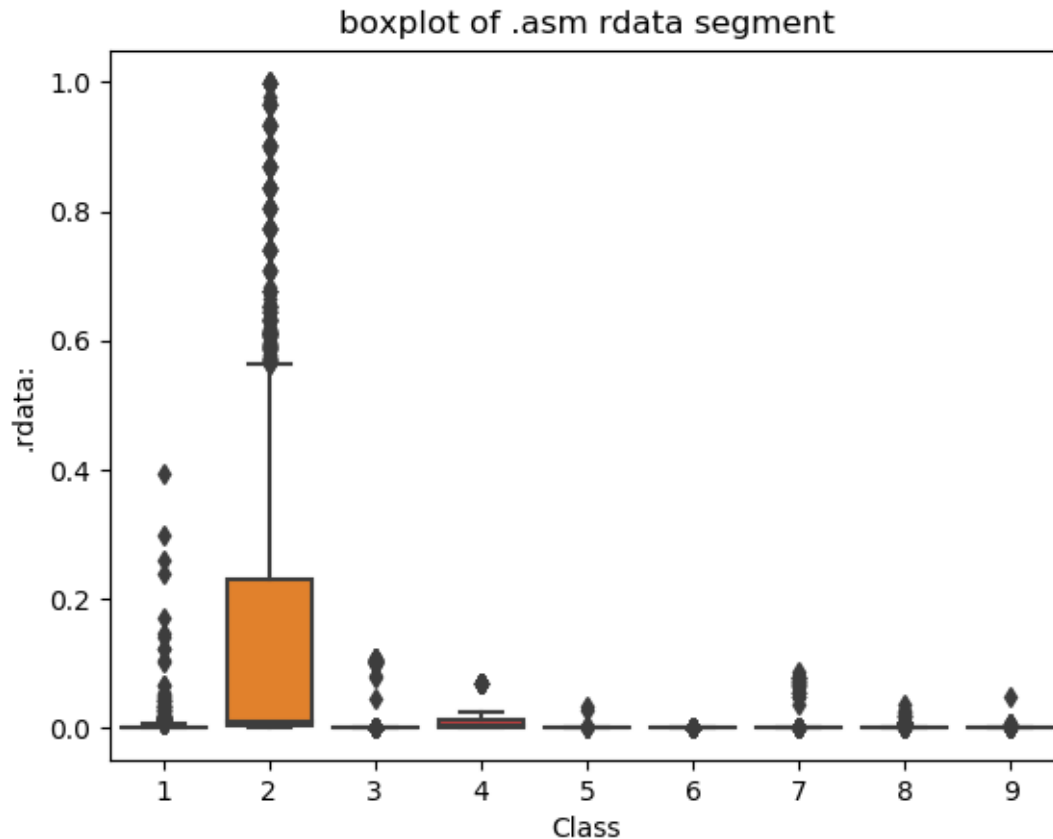
<IPython.core.display.Javascript object>



Here we can see that very less number of files has .bss segments. The classes cannot be separated well from this.

```
In [96]: 1 #Distribution of .rdata segments
2 ax = sns.boxplot(x="Class", y=".rdata:", data=result_asm)
3 plt.title("boxplot of .asm rdata segment")
4 plt.show()
```

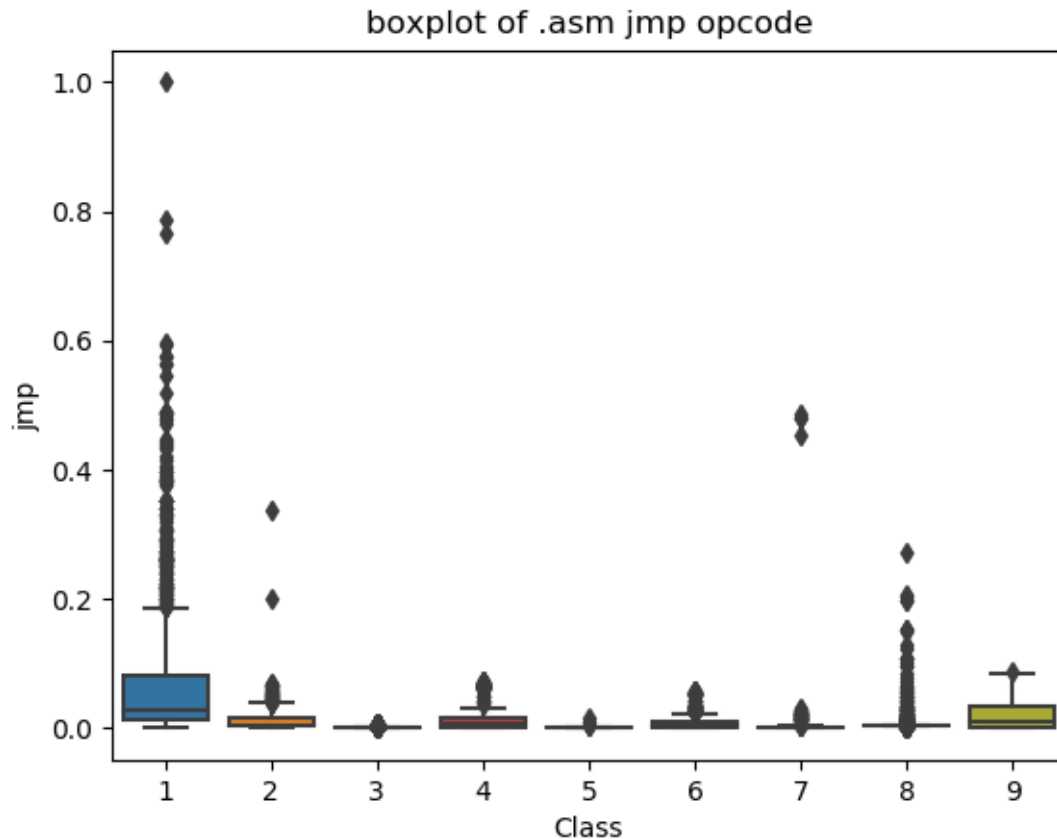
<IPython.core.display.Javascript object>



Here we can see that using .rdata segments, class 2 can be very well separated from the rest of the malware classes. However for rest of the malware classes the .rdata segment distribution looks pretty off.

```
In [97]: 1 #Distribution of .jmp segments
2 ax = sns.boxplot(x="Class", y="jmp", data=result_asm)
3 plt.title("boxplot of .asm jmp opcode")
4 plt.show()
```

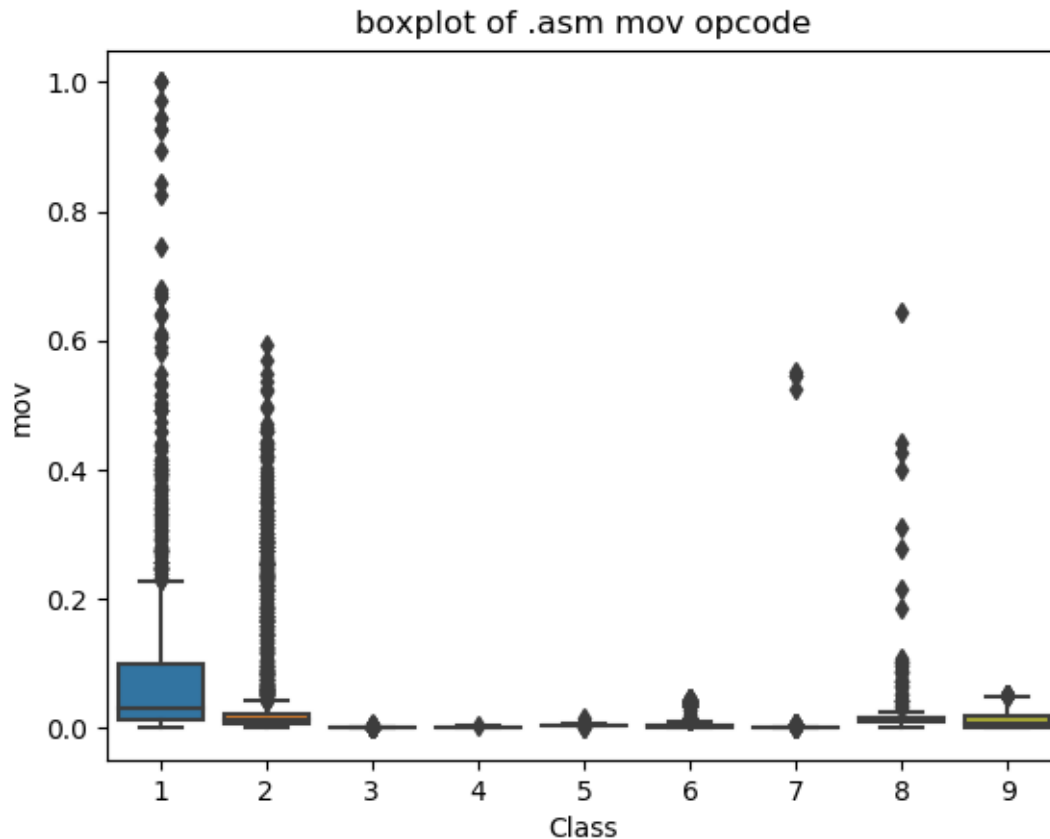
<IPython.core.display.Javascript object>



Here class 1 can be well separated from the rest of the files. However, this feature cannot separate other classes too well. We see almost 75% of points in class 2 have approximately 2000 jmp segments.


```
In [98]: 1 #Distribution of mov segments
2 ax = sns.boxplot(x="Class", y="mov", data=result_asm)
3 plt.title("boxplot of .asm mov opcode")
4 plt.show()
```

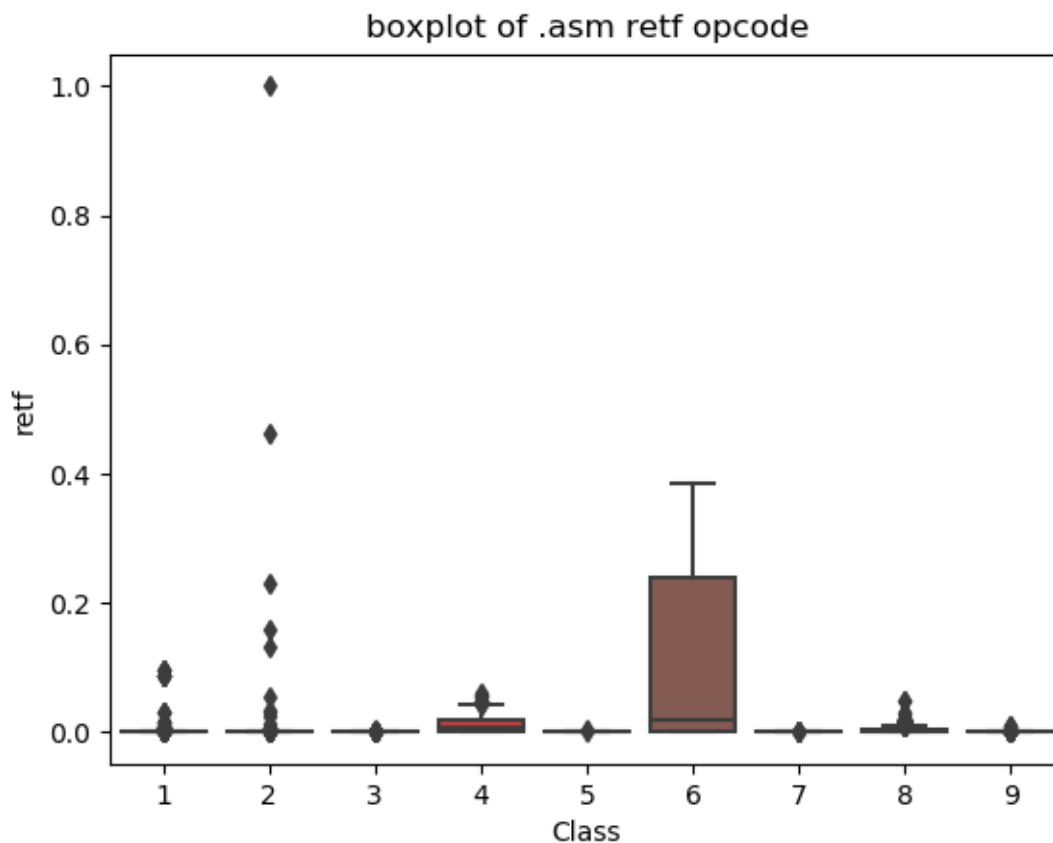
<IPython.core.display.Javascript object>



Here class 1 can be well separated from the rest of the files. However, this feature cannot separate other classes too well. We see almost 75% of points in class 2 have approximately 2000 jmp segments.

```
In [99]: 1 #Distribution of retf vs class labels
2 ax = sns.boxplot(x="Class", y="retf", data=result_asm)
3 plt.title("boxplot of .asm retf opcode")
4 plt.show()
```

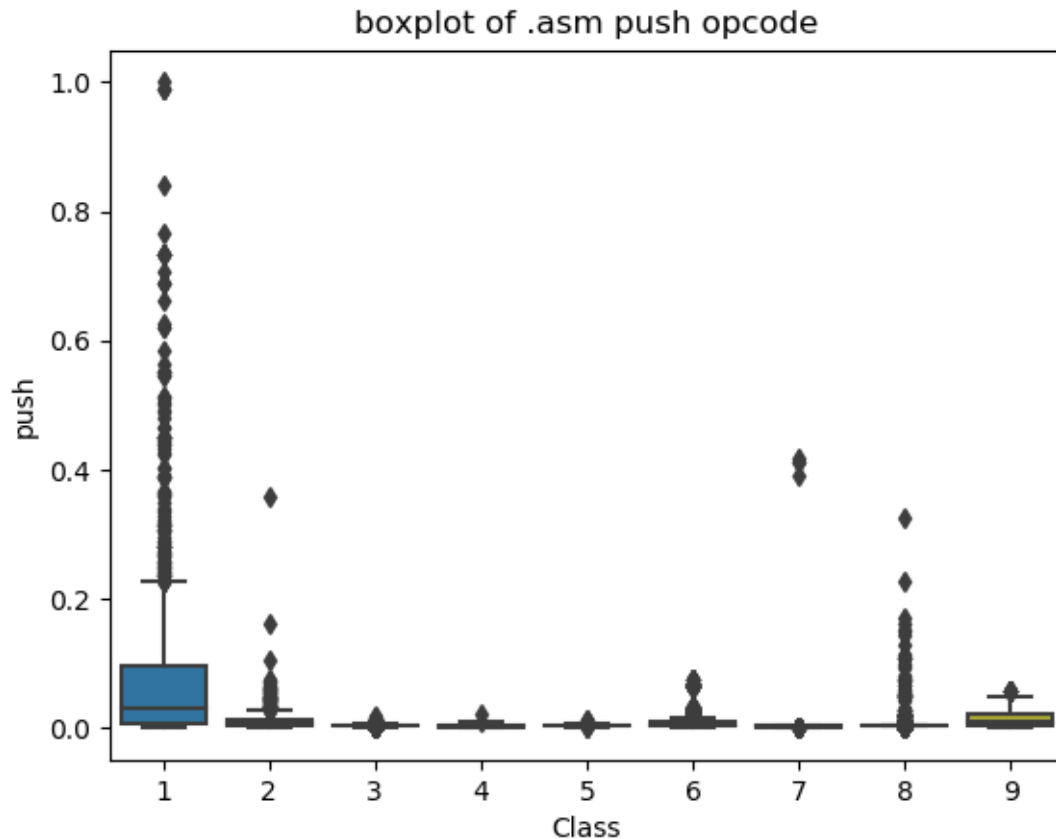
<IPython.core.display.Javascript object>



Here we see that using retf features, class 6 type of malware can be very easily separated from the rest of the malware classes.

```
In [100]: 1 #Distribution of push vs class labels
2 ax = sns.boxplot(x="Class", y="push", data=result_asm)
3 plt.title("boxplot of .asm push opcode")
4 plt.show()
```

<IPython.core.display.Javascript object>



Here class 1 can be very well separated from the rest of classes using push features. The variance is very low for all other classes for this feature.

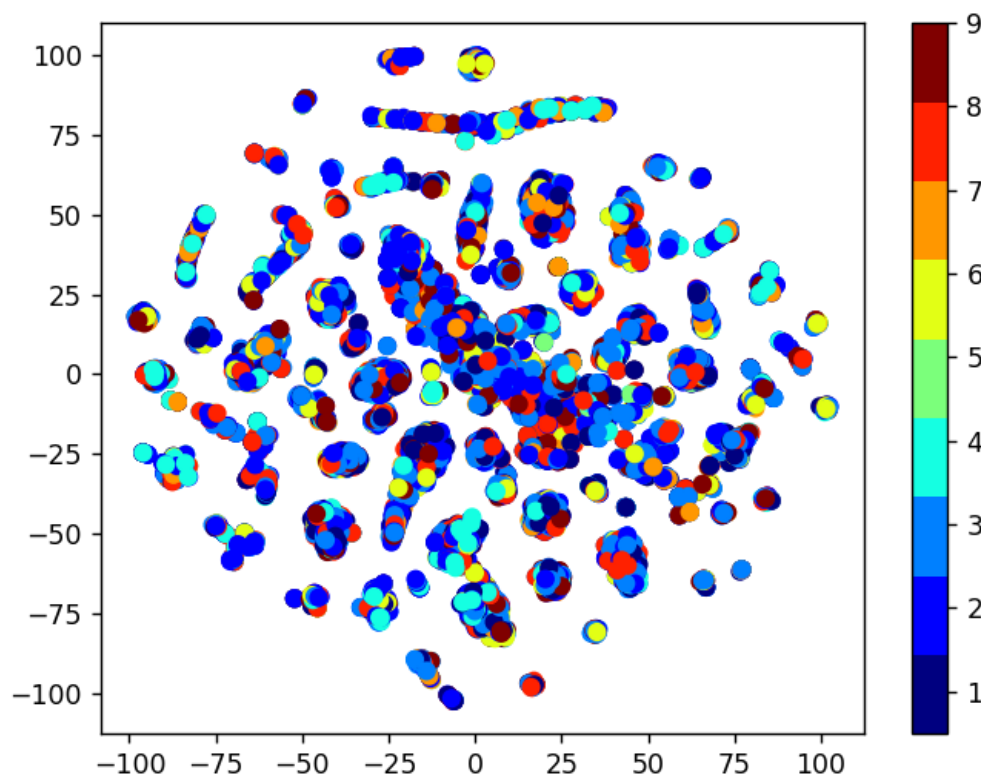
4.2.2 Multivariate Analysis on .asm file features

```

In [0]: 1 # by univariate analysis on the .asm file features we are getting very neglig
2 # 'rtn', '.BSS:', '.CODE' features, so here we are trying multivariate analysis
3 # the plot looks very messy
4
5 #Multivariate analysis on asm files features extracted using unigrams. We use
6
7 #Multivariate analysis on byte files features extracted using unigrams.
8 def draw_tsne_asm(p):
9     xtsne=TSNE(perplexity=p)
10    results=xtsne.fit_transform(result.drop(['ID','Class'], axis=1))
11    vis_x = results[:, 0]
12    vis_y = results[:, 1]
13    plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
14    plt.colorbar(ticks=range(10))
15    plt.clim(0.5, 9)
16    plt.show()
17
18 draw_tsne_asm(30)

```

<IPython.core.display.Javascript object>



From the above TSNE plot, it's clear that the features we have extracted are certainly useful in determining the classes. There is a partial separability amongst all the features.

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.

Key Take-aways from the EDA section.

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 [51]: 1 #We will drop .BSS, .rtn, .CODE features from the dataset because we have seen
2 asm_y = result_asm['Class']
3 asm_x = result_asm.drop(['ID', 'Class', '.BSS:', 'rtn', '.CODE'], axis=1)
```

```
In [52]: 1 #Train: 64%, Cross Validation 16%, Test 20%.
2 X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x, asm_y, test_size=0.2)
3 X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_asm, test_size=0.2)
```

```
In [53]: 1 print( X_cv_asm.isnull().all())
```

```
size      False
HEADER:    False
.text:     False
.Pav:      False
.idata:    False
.data:     False
.bss:      False
.rdata:    False
.edata:    False
.rsrc:     False
.tls:      False
.reloc:    False
jmp        False
mov        False
retf       False
push       False
pop        False
xor        False
retn       False
...
```

4.4. Machine Learning models on features of .asm files

4.4.1 K-Nearest Neighbors

```

In [0]: 1 # find more about KNeighborsClassifier() here http://scikit-learn.org/stable/
2 # -----
3 # default parameter
4 # KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', le
5 # metric='minkowski', metric_params=None, n_jobs=1, **kwargs)
6
7 # methods of
8 # fit(X, y) : Fit the model using X as training data and y as target values
9 # predict(X):Predict the class labels for the provided data
10 # predict_proba(X):Return probability estimates for the test data X.
11 #-----
12 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online
13 #-----
14
15
16 # find more about CalibratedClassifierCV here at http://scikit-learn.org/stab
17 # -----
18 # default paramters
19 # sklearn.calibration.CalibratedClassifierCV(base_estimator=None, method='sig
20 #
21 # some of the methods of CalibratedClassifierCV()
22 # fit(X, y[, sample_weight]) Fit the calibrated model
23 # get_params([deep]) Get parameters for this estimator.
24 # predict(X) Predict the target of new samples.
25 # predict_proba(X) Posterior probabilities of classification
26 #-----
27 # video link:
28 #-----
29
30 alpha = [x for x in range(1, 21,2)]
31 cv_log_error_array=[]
32 for i in alpha:
33     k_cfl=KNeighborsClassifier(n_neighbors=i)
34     k_cfl.fit(X_train_asm,y_train_asm)
35     sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
36     sig_clf.fit(X_train_asm, y_train_asm)
37     predict_y = sig_clf.predict_proba(X_cv_asm)
38     cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=k_cfl.clas
39
40 for i in range(len(cv_log_error_array)):
41     print ('log_loss for k = ',alpha[i],'is',cv_log_error_array[i])
42
43 best_alpha = np.argmin(cv_log_error_array)
44
45 fig, ax = plt.subplots()
46 ax.plot(alpha, cv_log_error_array,c='g')
47 for i, txt in enumerate(np.round(cv_log_error_array,3)):
48     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
49 plt.grid()
50 plt.title("Cross Validation Error for each alpha")
51 plt.xlabel("Alpha i's")
52 plt.ylabel("Error measure")
53 plt.show()
54
55 k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
56 k_cfl.fit(X_train_asm,y_train_asm)

```

```

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

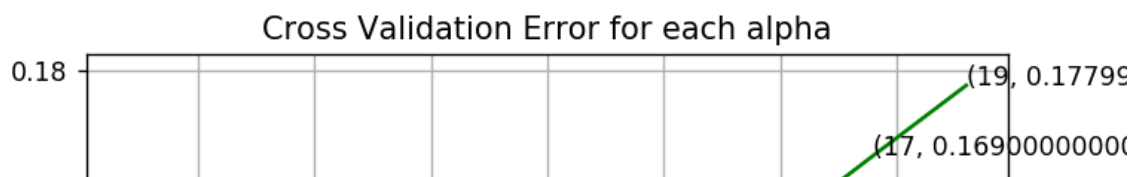
```

```

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

```

<IPython.core.display.Javascript object>



4.4.2 Logistic Regression

```

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

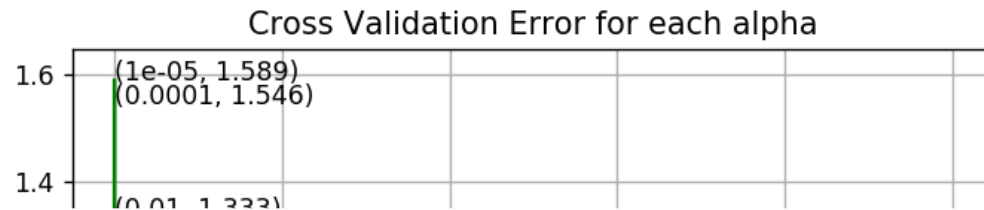
```

log_loss for c = 1e-05 is 1.58867274165

log_loss for c = 0.0001 is 1.54560797884


```
log_loss for c = 0.001 is 1.30137786807
log_loss for c = 0.01 is 1.33317456931
log_loss for c = 0.1 is 1.16705751378
log_loss for c = 1 is 0.757667807779
log_loss for c = 10 is 0.546533939819
log_loss for c = 100 is 0.438414998062
log_loss for c = 1000 is 0.424423536526
```

<IPython.core.display.Javascript object>



4.4.3 Random Forest Classifier

In [0]:

```

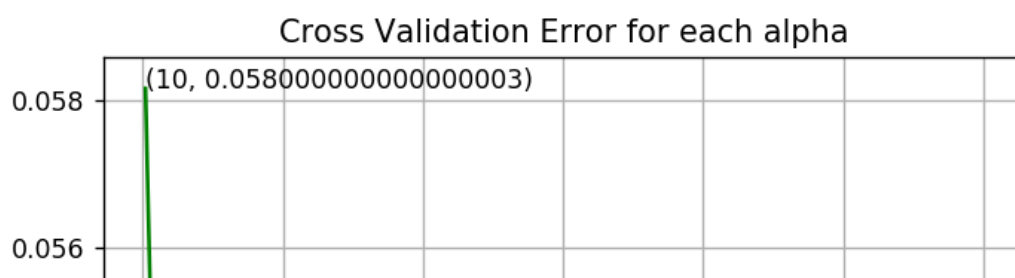
1  # -----
2  # default parameters
3  # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini',
4  # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_
5  # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_
6  # class_weight=None)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

```

```
57 plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```

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

<IPython.core.display.Javascript object>



4.4.4 XgBoost Classifier

```

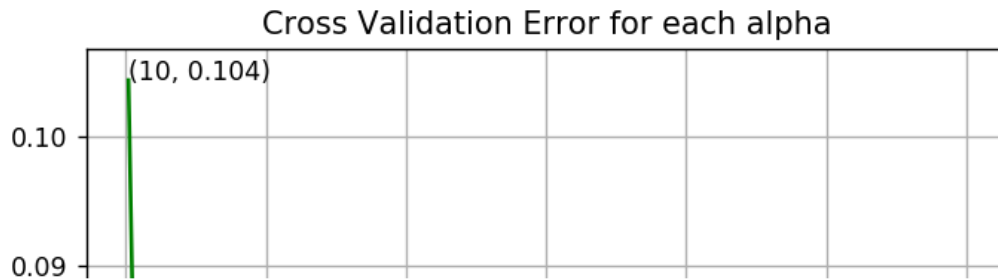
In [0]: 1 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
2
3 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en
4 # -----
5 # default paramters
6 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100,
7 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
8 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_lambda=1,
9 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None)
10
11 # some of methods of RandomForestRegressor()
12 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None)
13 # get_params([deep]) Get parameters for this estimator.
14 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: If ntree_limit is not None,
15 # get_score(importance_type='weight') -> get the feature importance
16 # -----
17 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online
18 # -----
19
20 alpha=[10,50,100,500,1000,2000,3000]
21 cv_log_error_array=[]
22 for i in alpha:
23     x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
24     x_cfl.fit(X_train_asm,y_train_asm)
25     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
26     sig_clf.fit(X_train_asm, y_train_asm)
27     predict_y = sig_clf.predict_proba(X_cv_asm)
28     cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=x_cfl.classes_))
29
30 for i in range(len(cv_log_error_array)):
31     print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
32
33
34 best_alpha = np.argmin(cv_log_error_array)
35
36 fig, ax = plt.subplots()
37 ax.plot(alpha, cv_log_error_array,c='g')
38 for i, txt in enumerate(np.round(cv_log_error_array,3)):
39     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
40 plt.grid()
41 plt.title("Cross Validation Error for each alpha")
42 plt.xlabel("Alpha i's")
43 plt.ylabel("Error measure")
44 plt.show()
45
46 x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
47 x_cfl.fit(X_train_asm,y_train_asm)
48 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
49 sig_clf.fit(X_train_asm, y_train_asm)
50
51 predict_y = sig_clf.predict_proba(X_train_asm)
52
53 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is", log_loss(y_train_asm, predict_y))
54 predict_y = sig_clf.predict_proba(X_cv_asm)
55 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is", log_loss(y_cv_asm, predict_y))
56 predict_y = sig_clf.predict_proba(X_test_asm)

```

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

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

<IPython.core.display.Javascript object>



4.4.5 Xgboost Classifier with best hyperparameters

```
In [0]: 1 x_cfl=XGBClassifier()
2
3 prams={
4     'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
5     'n_estimators':[100,200,500,1000,2000],
6     'max_depth':[3,5,10],
7     'colsample_bytree':[0.1,0.3,0.5,1],
8     'subsample':[0.1,0.3,0.5,1]
9 }
10 random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jc
11 random_cfl.fit(X_train_asm,y_train_asm)
```

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

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

```
Out[163]: RandomizedSearchCV(cv=None, error_score='raise',
        estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, 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]: 1 print (random_cfl.best_params_)

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

```

In [0]: 1 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
2
3 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en
4 # -----
5 # default paramters
6 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100
7 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamm
8 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg
9 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None
10
11 # some of methods of RandomForestRegressor()
12 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppi
13 # get_params([deep]) Get parameters for this estimator.
14 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE
15 # get_score(importance_type='weight') -> get the feature importance
16 # -----
17 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-onlin
18 # -----
19
20 x_cfl=XGBClassifier(n_estimators=200,subsample=0.5,learning_rate=0.15,colsamp
21 x_cfl.fit(X_train_asm,y_train_asm)
22 c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
23 c_cfl.fit(X_train_asm,y_train_asm)
24
25 predict_y = c_cfl.predict_proba(X_train_asm)
26 print ('train loss',log_loss(y_train_asm, predict_y))
27 predict_y = c_cfl.predict_proba(X_cv_asm)
28 print ('cv loss',log_loss(y_cv_asm, predict_y))
29 predict_y = c_cfl.predict_proba(X_test_asm)
30 print ('test loss',log_loss(y_test_asm, predict_y))

```

train loss 0.0102661325822

cv loss 0.0501201796687

test loss 0.0483908764397

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

4.5.1. Merging both asm and byte file features

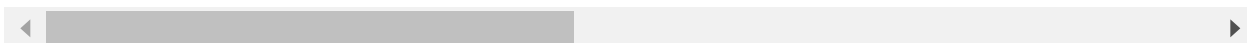
In [60]:

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

Out[60]:

		ID	0	1	2	3	4	5	6
0	01azqd4InC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	
2	01jsnpXSAlgW6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	

5 rows × 260 columns



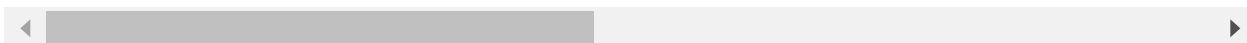
In [61]:

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

Out[61]:

		ID	size	Class	HEADER:	.text:	.Pav:	.idata:	.data:	.b
0	01azqd4InC7m9JpocGv5	0.400910		9	0.101695	0.032927	0.0	0.006937	0.542847	0.0001
1	01lsoiSMh5gxyDYTI4CB	0.099719		2	0.000000	0.161391	0.0	0.003690	0.009758	0.0001
2	01jsnpXSAlgW6aPeDxrU	0.060553		9	0.101695	0.101121	0.0	0.001821	0.000263	0.0001
3	01kcPWA9K2BOxQeS5Rju	0.000432		1	0.107345	0.001092	0.0	0.000761	0.000023	0.0001
4	01SuzwMJEIXsK7A8dQbl	0.006983		8	0.101695	0.015220	0.0	0.001234	0.001825	0.0121

5 rows × 54 columns



In [62]:

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

(10868, 260)

(10868, 54)

In [63]:

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

Out[63]:

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

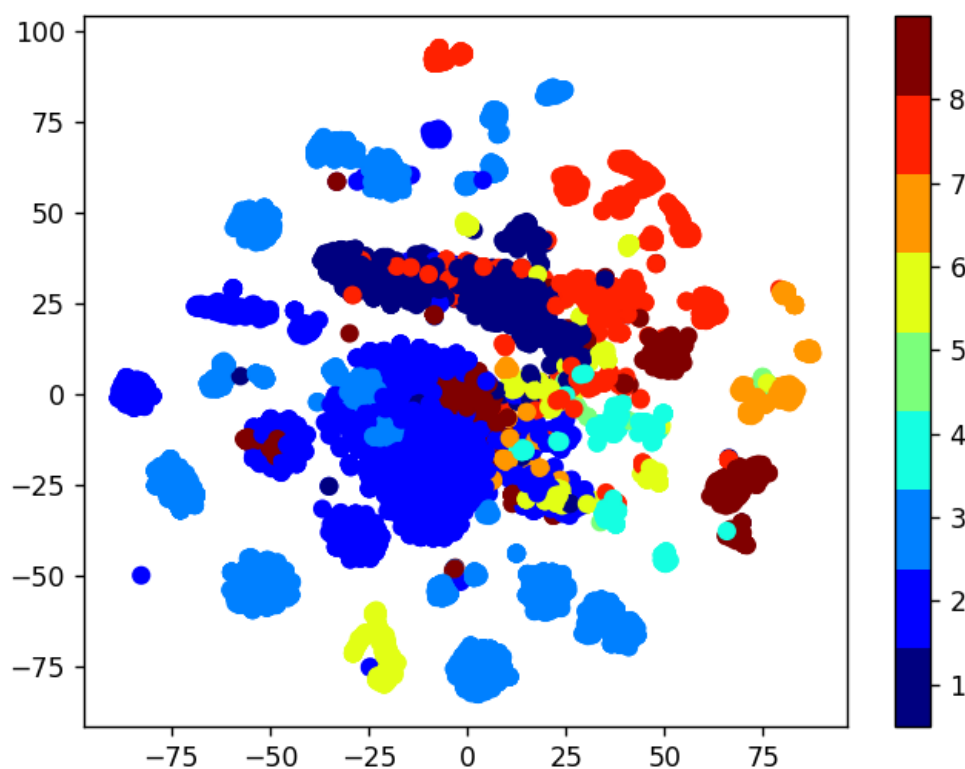
5 rows × 307 columns



4.5.2. Multivariate Analysis on final features

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

<IPython.core.display.Javascript object>



4.5.3. Train and Test split

```
In [0]: 1 X_train, X_test_merge, y_train, y_test_merge = train_test_split(result_x, res
2 X_train_merge, X_cv_merge, y_train_merge, y_cv_merge = train_test_split(X_tra
```

4.5.4. Random Forest Classifier on final features

In [0]:

```

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

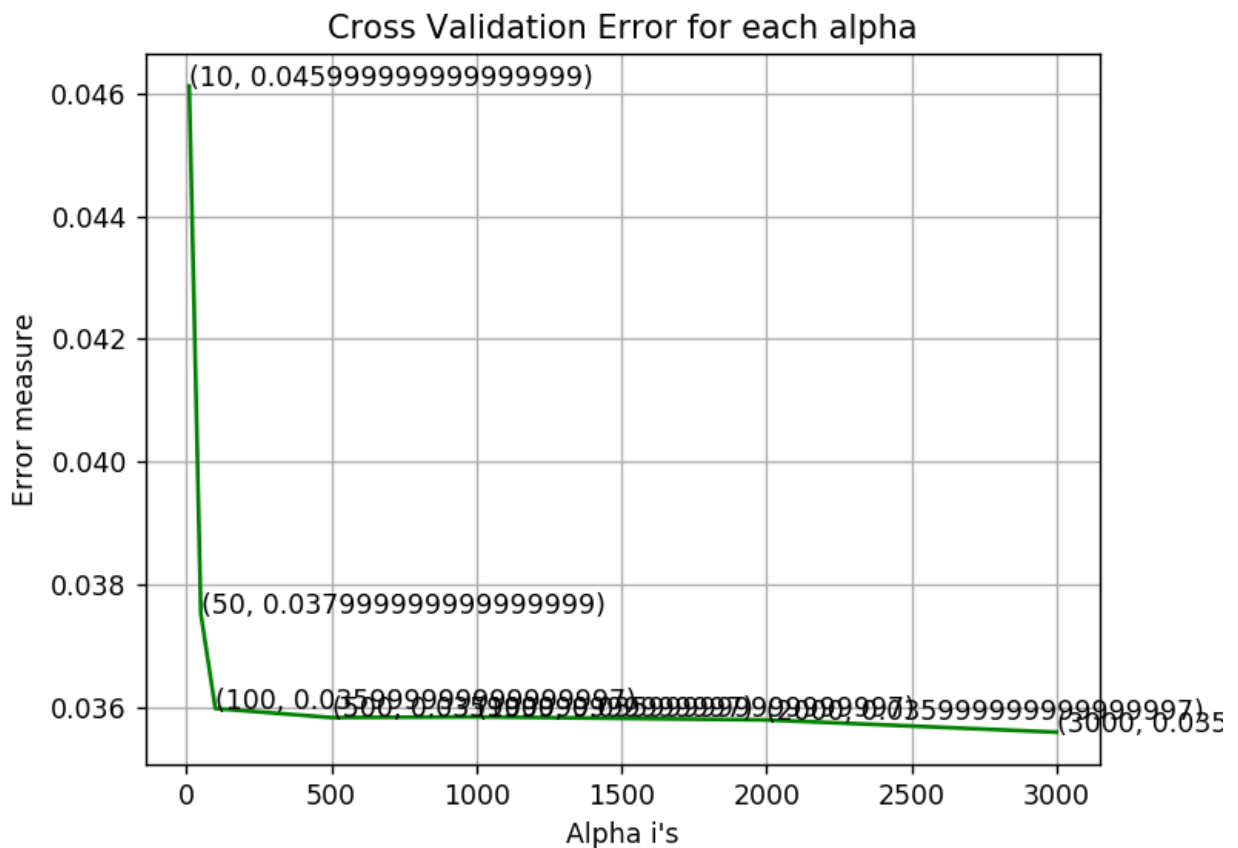
```

```

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

```

<IPython.core.display.Javascript object>



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]: 1 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
2
3 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en
4 # -----
5 # default paramters
6 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100,
7 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0,
8 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_lambda=1,
9 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None)
10
11 # some of methods of RandomForestRegressor()
12 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None)
13 # get_params([deep]) Get parameters for this estimator.
14 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: If ntree_limit is not None,
15 # get_score(importance_type='weight') -> get the feature importance
16 # -----
17 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online
18 # -----
19
20 alpha=[10,50,100,500,1000,2000,3000]
21 cv_log_error_array=[]
22 for i in alpha:
23     x_cfl=XGBClassifier(n_estimators=i)
24     x_cfl.fit(X_train_merge,y_train_merge)
25     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
26     sig_clf.fit(X_train_merge, y_train_merge)
27     predict_y = sig_clf.predict_proba(X_cv_merge)
28     cv_log_error_array.append(log_loss(y_cv_merge, predict_y, labels=x_cfl.classes_))
29
30 for i in range(len(cv_log_error_array)):
31     print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
32
33
34 best_alpha = np.argmin(cv_log_error_array)
35
36 fig, ax = plt.subplots()
37 ax.plot(alpha, cv_log_error_array,c='g')
38 for i, txt in enumerate(np.round(cv_log_error_array,3)):
39     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
40 plt.grid()
41 plt.title("Cross Validation Error for each alpha")
42 plt.xlabel("Alpha i's")
43 plt.ylabel("Error measure")
44 plt.show()
45
46 x_cfl=XGBClassifier(n_estimators=3000,nthread=-1)
47 x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
48 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
49 sig_clf.fit(X_train_merge, y_train_merge)
50
51 predict_y = sig_clf.predict_proba(X_train_merge)
52 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is", cv_log_error_array[best_alpha])
53 predict_y = sig_clf.predict_proba(X_cv_merge)
54 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is", cv_log_error_array[best_alpha])
55 predict_y = sig_clf.predict_proba(X_test_merge)
56 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is", log_loss(y_test_merge, predict_y, labels=x_cfl.classes_))

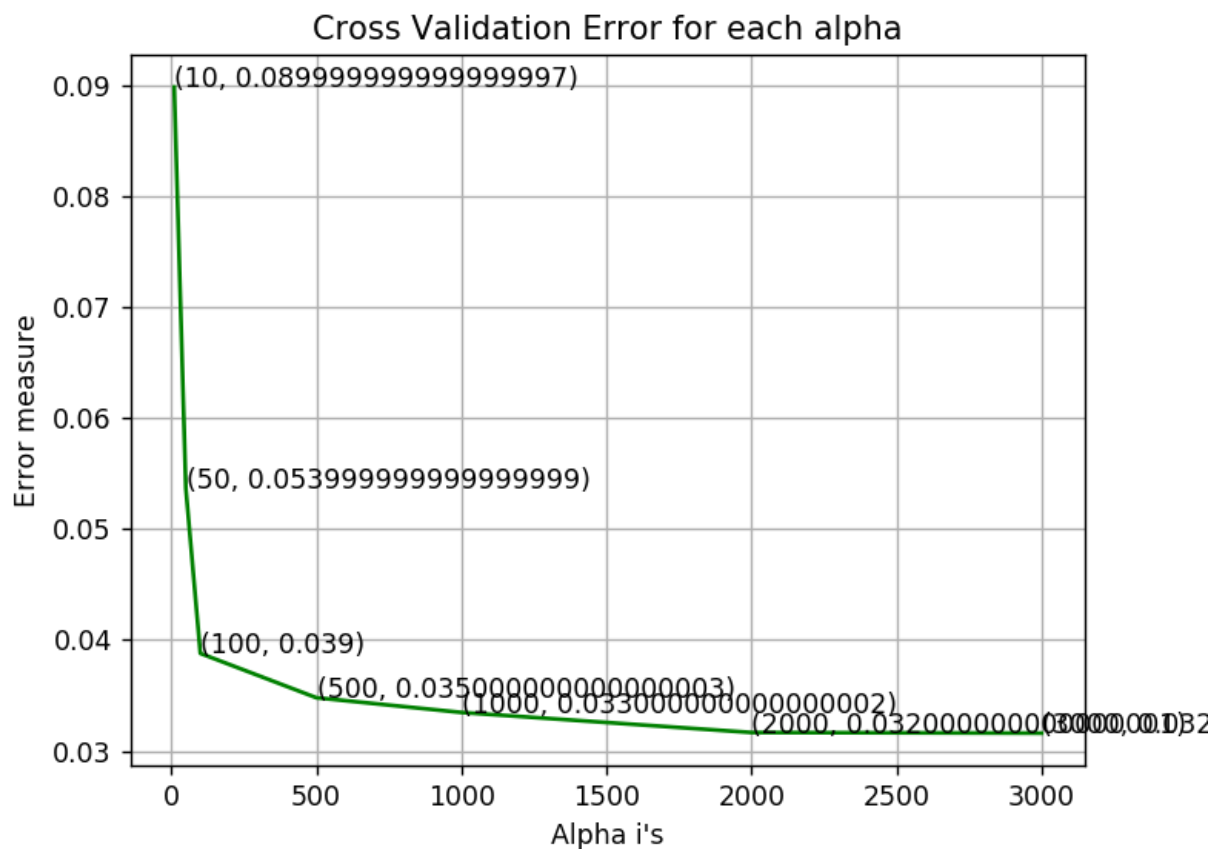
```

```

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

```

<IPython.core.display.Javascript object>



```

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: 0.0315972694
477
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]: 1 x_cfl=XGBClassifier()
2
3 prams={
4     'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
5     'n_estimators':[100,200,500,1000,2000],
6     'max_depth':[3,5,10],
7     'colsample_bytree':[0.1,0.3,0.5,1],
8     'subsample':[0.1,0.3,0.5,1]
9 }
10 random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jc
11 random_cfl.fit(X_train_merge, y_train_merge)
```

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

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

```
Out[187]: RandomizedSearchCV(cv=None, error_score='raise',
        estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, 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], 'co
        lsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
        pre_dispatch='2*n_jobs', random_state=None, refit=True,
        return_train_score=True, scoring=None, verbose=10)
```

```
In [0]: 1 print (random_cfl.best_params_)

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

In [0]:

```

1
2 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en
3 # -----
4 # default paramters
5 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100
6 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma
7 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg
8 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None)
9
10 # some of methods of RandomForestRegressor()
11 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stoppin
12 # get_params([deep]) Get parameters for this estimator.
13 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE
14 # get_score(importance_type='weight') -> get the feature importance
15 # -----
16 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online
17 # -----
18
19 x_cfl=XGBClassifier(n_estimators=1000,max_depth=10,learning_rate=0.15,colsample
20 x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
21 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
22 sig_clf.fit(X_train_merge, y_train_merge)
23
24 predict_y = sig_clf.predict_proba(X_train_merge)
25 print('For values of best alpha = ', alpha[best_alpha], "The train log loss
26 predict_y = sig_clf.predict_proba(X_cv_merge)
27 print('For values of best alpha = ', alpha[best_alpha], "The cross validation
28 predict_y = sig_clf.predict_proba(X_test_merge)
29 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
30 plot_confusion_matrix(y_test_as, 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

```

In [2]: 1 import warnings
2 warnings.filterwarnings("ignore")
3 import shutil
4 import os
5 import pandas as pd
6 import matplotlib
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 import numpy as np
10 import pickle
11 from sklearn.manifold import TSNE
12 from sklearn import preprocessing
13 import pandas as pd
14 from multiprocessing import Process# this is used for multithreading
15 import multiprocessing
16 import codecs# this is used for file operations
17 import random as r
18 from xgboost import XGBClassifier
19 from sklearn.model_selection import RandomizedSearchCV
20 from sklearn.tree import DecisionTreeClassifier
21 from sklearn.calibration import CalibratedClassifierCV
22 from sklearn.neighbors import KNeighborsClassifier
23 from sklearn.metrics import log_loss
24 from sklearn.metrics import confusion_matrix
25 from sklearn.model_selection import train_test_split
26 from sklearn.linear_model import LogisticRegression
27 from sklearn.ensemble import RandomForestClassifier
28 # matplotlib.use(u'nbAgg')
29 # %matplotlib inline

```

5.1 Bigrams Feature extraction from byte files

Get the bigrams corpus

```

In [ ]: 1 byte_vocab = "00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,

```

```

In [2]: 1 #This function will return all the possible 257*257 combinations of bigrams t
2 def get_bigram_tokens(tokens):
3     sent=""
4     byte_bigram_vocab=[]
5     for i in range(len(tokens)):
6         for j in range(len(tokens)):
7             bigram=tokens[i]+" "+tokens[j]
8             sent=sent+bigram+", "
9             byte_bigram_vocab.append(bigram)
10    return byte_bigram_vocab
11
12 byte_bigram_vocab=get_bigram_tokens(tokens) #This will contain all the possib

```

Extract the byte bigram features and save the dataframe containing bi-

gram features

```
In [5]: 1 #Extract the byte bigram features and save the dataframe containing bi-gram f
2 vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False,ng
3 file_lists=os.listdir('byteFiles')
4 features=["ID"]+vectorizer.get_feature_names()
5 byte_bigram_df=pd.DataFrame(columns=features)
6
7 with open("features/byte_bigram_df.csv", mode='w') as byte_bigram_df:
8     byte_bigram_df.write(','.join(map(str, features)))
9     byte_bigram_df.write('\n')
10    for _, file in tqdm(enumerate(file_lists)):
11        file_id=file.split(".")[0] #ID of each file
12        file = open('byteFiles/' + file)
13        corpus=[file.read().replace('\n', ' ').lower()] #This will contain al
14        bigrams=vectorizer.transform(corpus) #This will return a sparse vecto
15        row=scipy.sparse.csr_matrix(bigrams).toarray() #Update each row of th
16        byte_bigram_df.write(','.join(map(str, [file_id]+list(row[0])))) #Wri
17        byte_bigram_df.write('\n')
18        file.close()
```

10868it [3:29:50, 3.07it/s]

Read the byte bigrams count dataset

```
In [ ]: 1 #Use dask dataframe to avoid memory problems
2 byte_bigram_df=dd.read_csv("features/byte_bigram_df.csv",sample=256000000)
```

5.2 Registers Bigrams Feature extraction from ASM files

Get the registers sequence from the ASM files

```
In [6]: 1 registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
2 registers_dict = dict(zip(registers, [1 for i in range(len(registers))]))
```

```

In [7]: 1 if not os.path.isdir("asmFilesRegisters"):
2         os.mkdir('asmFilesRegisters')
3
4         #Get registers sequences for each of the files and save it as a text file. Ea
5         #at regular intervals with some words between them. So we have to extract the
6         #a bigram matrix of vectors can be calculated which will give us the 2 grams
7     def get_registers_seq():
8         filenames=os.listdir('asmFiles')
9         for asmfile in tqdm(filenames):
10            registers_file = open("asmFilesRegisters/{}_registers_bigrams.txt".fo
11            registers_sequence = ""
12            with codecs.open('asmFiles/' + asmfile, encoding='cp1252', errors ='r
13                for lines in file:
14                    line = lines.rstrip().split()
15                    for word in line:
16                        if registers_dict.get(word)==1:
17                            registers_sequence += word + ' '
18            registers_file.write(registers_sequence + "\n")
19            registers_file.close()
20
21    get_registers_seq()
22
23    registers_bigram_vocab=get_bigram_tokens(registers) #This will contain all th

```

100%|██████████| 10868/10868 [9:05:09<00:00, 1.03s/it]

Extract the registers bigram features and save the dataframe containing bi-gram features

```

In [8]: 1 #Extract the bigram features and save the dataframe containing bi-gram featur
2     vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False,ng
3     file_lists=os.listdir('asmFilesRegisters')
4     features=["ID"]+vectorizer.get_feature_names()
5     registers_bigram_df=pd.DataFrame(columns=features)
6
7     with open("features/registers_bigram_df.csv", mode='w') as registers_bigram_d
8         registers_bigram_df.write(','.join(map(str, features)))
9         registers_bigram_df.write('\n')
10    for _, file in tqdm(enumerate(file_lists)):
11        file_id=file.split("_")[0] #ID of each file
12        file = open('asmFilesRegisters/' + file)
13        corpus=[file.read().replace('\n', ' ').lower()] #This will contain al
14        bigrams=vectorizer.transform(corpus) #This will return a sparse vecto
15        row=scipy.sparse.csr_matrix(bigrams).toarray() #Update each row of th
16        registers_bigram_df.write(','.join(map(str, [file_id]+list(row[0]))))
17        registers_bigram_df.write('\n')
18        file.close()

```

10868it [03:18, 54.69it/s]

Display the registers bigrams count dataset

```
In [129]: 1 registers_bigram_df=pd.read_csv("features/registers_bigram_df.csv")
          2 registers_bigram_df.head()
```

Out[129]:

	ID	edx edx	edx esi	edx eax	edx ebx	edx ecx	edx edi	edx ebp	edx esp	edx eip	...	esp eip	eip edx	eip esi	eip eax
0	01azqd4lnC7m9JpocGv5	0	0	0	0	0	0	0	0	0	...	0	0	0	0
1	01lsoiSMh5gxyDYTI4CB	31	12	72	15	68	23	5	8	0	...	0	0	0	0
2	01jsnpXSAIgw6aPeDxrU	66	0	72	1	50	0	5	0	0	...	0	0	0	0
3	01kcPWA9K2BOxQeS5Rju	1	1	4	4	0	0	0	0	0	...	0	0	0	0
4	01SuzwMJEIXsK7A8dQbl	8	1	374	0	12	1	8	3	0	...	0	0	0	0

5 rows × 82 columns



```
In [ ]: 1 del(registers_bigram_df)
```

5.3 Registers Trigrams Feature extraction from ASM files

Get the trigram corpus

```
In [9]: 1 #TRIGRAMS
          2 registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
          3
          4 #This function will return all the possible n*n combinations of trigrams that
          5 def get_trigram_tokens(tokens):
          6     sent=""
          7     registers_trigram_vocab=[]
          8     for i in range(len(tokens)):
          9         for j in range(len(tokens)):
         10             for k in range(len(tokens)):
         11                 trigram=tokens[i]+" "+tokens[j]+" "+tokens[k]
         12                 registers_trigram_vocab.append(trigram)
         13     return registers_trigram_vocab
         14
         15 registers_trigram_vocab=get_trigram_tokens(registers) #This will contain all
```

Extract the registers trigram features and save the dataframe containing tri-gram features

```
In [10]: 1 #Extract the trigram features and save the dataframe containing tri-gram feat
2 vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False,ng
3 file_lists=os.listdir('asmFilesRegisters')
4 features=["ID"]+vectorizer.get_feature_names()
5 registers_trigram_df=pd.DataFrame(columns=features)
6
7 with open("features/registers_trigram_df.csv", mode='w') as registers_trigram
8     registers_trigram_df.write(','.join(map(str, features)))
9     registers_trigram_df.write('\n')
10    for _, file in tqdm(enumerate(file_lists)):
11        file_id=file.split("_")[0] #ID of each file
12        file = open('asmFilesRegisters/' + file)
13        corpus=[file.read().replace('\n', ' ').lower()] #This will contain al
14        trigrams=vectorizer.transform(corpus) #This will return a sparse vect
15        row=scipy.sparse.csr_matrix(trigrams).toarray() #Update each row of t
16        registers_trigram_df.write(','.join(map(str, [file_id]+list(row[0]))))
17        registers_trigram_df.write('\n')
18        file.close()
```

10868it [01:12, 148.92it/s]

Display the registers trigrams count dataset

```
In [138]: 1 registers_trigram_df=pd.read_csv("features/registers_trigram_df.csv")
2 registers_trigram_df.head()
```

Out[138]:

	ID	edx edx	edx esi	edx eax	edx ebx	edx ecx	edx edi	edx ebp	edx esp	edx eip	...	eip esp	eip eip	eip esi	eip eax
0	01azqd4lnC7m9JpocGv5	0	0	0	0	0	0	0	0	0	...	0	0	0	0
1	01lsoiSMh5gxyDYTI4CB	6	1	10	0	12	2	0	0	0	...	0	0	0	0
2	01jsnpXSAlgW6aPeDxrU	23	0	18	0	23	0	2	0	0	...	0	0	0	0
3	01kcPWA9K2BOxQeS5Rju	0	0	0	1	0	0	0	0	0	...	0	0	0	0
4	01SuzwMJEIXsK7A8dQbl	3	0	1	0	2	0	2	0	0	...	0	0	0	0

5 rows × 730 columns

```
In [ ]: 1 del(registers_trigram_df)
```

5.4 Opcodes Bigrams Feature extraction from ASM files

```
In [11]: 1 #Putting the list in a dictionary to decrease time complexity
2 opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub',
3 opcode_dict = dict(zip(opcodes, [1 for i in range(len(opcodes))]))
```

Get the opcodes sequence from the ASM files

```
In [3]: 1 opcode_dict = dict(zip(opcodes, [1 for i in range(len(opcodes))]))
2 if not os.path.isdir("asmFilesOpcodes"):
3     os.mkdir('asmFilesOpcodes')
4
5 #Get opcode sequences for each of the files and save it as a text file. Each
6 #at regular intervals with some words between them. So we have to extract the
7 #a bigram matrix of vectors can be calculated which will give us the 2 grams
8 def get_opcode_seq():
9     filenames=os.listdir('asmFiles')
10    for asmfile in tqdm(filenames):
11        opcode_file = open("asmFilesOpcodes/{}_opcode_bigrams.txt".format(asm
12        opcode_sequence = ""
13        with codecs.open('asmFiles/' + asmfile, encoding='cp1252', errors='r
14            for lines in file:
15                line = lines.rstrip().split()
16                for word in line:
17                    if opcode_dict.get(word)==1:
18                        opcode_sequence += word + ' '
19        opcode_file.write(opcode_sequence + "\n")
20        opcode_file.close()
21
22 get_opcode_seq()
23
24 opcodes_bigram_vocab=get_bigram_tokens(opcodes)
```

100%|██████████| 10868/10868 [10:12:45<00:00, 1.09s/it]

Extract the opcodes bigram features and save the dataframe containing bi-gram features

```
In [14]: 1 #Extract the bigram features and save the dataframe containing bi-gram featur
2 vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False,ng
3 file_lists=os.listdir('asmFilesOpcodes')
4 features=["ID"]+vectorizer.get_feature_names()
5 opcodes_bigram_df=pd.DataFrame(columns=features)
6
7 with open("features/opcodes_bigram_df.csv", mode='w') as opcodes_bigram_df:
8     opcodes_bigram_df.write(','.join(map(str, features)))
9     opcodes_bigram_df.write('\n')
10    for _, file in tqdm(enumerate(file_lists)):
11        file_id=file.split("_")[0] #ID of each file
12        file = open('asmFilesOpcodes/' + file)
13        corpus=[file.read().replace('\n', ' ').lower()] #This will contain al
14        bigrams=vectorizer.transform(corpus) #This will return a sparse vecto
15        row=scipy.sparse.csr_matrix(bigrams).toarray() #Update each row of th
16        opcodes_bigram_df.write(','.join(map(str, [file_id]+list(row[0])))) #
17        opcodes_bigram_df.write('\n')
18        file.close()
```

10868it [03:47, 47.81it/s]

Display the opcodes bigrams count dataset

```
In [2]: 1 opcodes_bigram_df=dd.read_csv("features/opcodes_bigram_df.csv")
2 opcodes_bigram_df.head()
```

Out[2]:

	ID	jmp jmp	jmp mov	jmp retf	jmp push	jmp pop	jmp xor	jmp retn	jmp nop	jmp sub	...	movzx cmp	movzx call
0	01azqd4InC7m9JpocGv5	440	192	0	6	0	17	0	0	24	...	0	0
1	01IsoiSMh5gxyDYTI4CB	0	32	0	3	1	3	1	0	0	...	1	0
2	01jsnpXSAIgw6aPeDxrU	0	0	0	0	0	0	0	0	0	...	12	0
3	01kcPWA9K2BOxQeS5Rju	0	5	0	1	0	2	1	0	0	...	0	0
4	01SuzwMJEIXsK7A8dQbl	5	57	1	4	1	1	0	0	0	...	11	0

5 rows × 677 columns

```
In [ ]: 1 del(opcodes_bigram_df)
```

5.5 Opcodes Trigrams Feature extraction from ASM files

Get the trigram corpus

In [16]:

```

1 #Putting the list in a dictionary to decrease time complexity
2 opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub',
3 opcodes_trigram_vocab=get_trigram_tokens(opcodes) #This will contain all the

```

Extract the opcodes trigrams features and save the dataframe containing tri-gram features

In [17]:

```

1 #Extract the trigram features and save the dataframe containing tri-gram feat
2 vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False,ng
3 file_lists=os.listdir('asmFilesOpcodes')
4 features=["ID"]+vectorizer.get_feature_names()
5 opcodes_trigram_df=pd.DataFrame(columns=features)
6
7 with open("features/opcodes_trigram_df.csv", mode='w') as opcodes_trigram_df:
8     opcodes_trigram_df.write(','.join(map(str, features)))
9     opcodes_trigram_df.write('\n')
10    for _, file in tqdm(enumerate(file_lists)):
11        file_id=file.split("_")[0] #ID of each file
12        file = open('asmFilesOpcodes/' + file)
13        corpus=[file.read().replace('\n', ' ').lower()] #This will contain al
14        trigrams=vectorizer.transform(corpus) #This will return a sparse vect
15        row=scipy.sparse.csr_matrix(trigrams).toarray() #Update each row of t
16        opcodes_trigram_df.write(','.join(map(str, [file_id]+list(row[0]))))
17        opcodes_trigram_df.write('\n')
18        file.close()

```

10868it [06:25, 32.60it/s]

Load the opcodes trigrams count dataset

In [151]:

```

1 opcodes_trigram_df=pd.read_csv("features/opcodes_trigram_df.csv")
2 opcodes_trigram_df.head()

```

Out[151]:

		jmp	jmp	jmp	jmp	jmp	jmp	jmp	jmp	jmp	...	movzx	movzx
	ID	jmp	jmp	jmp	jmp	jmp	jmp	jmp	jmp	jmp		movzx	movzx
		jmp	mov	retf	push	pop	xor	retn	nop	sub		cmp	call
0	01azqd4lnC7m9JpocGv5	437	0	0	0	0	1	0	0	0	...	0	0
1	01lsoiSMh5gxyDYTI4CB	0	0	0	0	0	0	0	0	0	...	0	0
2	01jsnpXSAlgW6aPeDxrU	0	0	0	0	0	0	0	0	0	...	1	0
3	01kcPWA9K2BOxQeS5Rju	0	0	0	0	0	0	0	0	0	...	0	0
4	01SuzwMJEIXsK7A8dQbl	2	1	1	1	0	0	0	0	0	...	1	0

5 rows × 17577 columns

```
In [ ]: 1 del(opcodes_trigram_df)
```

5.6 Extract Image features from the ASM files

Refer: <https://www.kaggle.com/c/malware-classification/discussion/13897#latest-105551>
(<https://www.kaggle.com/c/malware-classification/discussion/13897#latest-105551>)

I have used the code snippets by:

1. Xiaozhou Wang, xiaozhou@ualberta.ca (<mailto:xiaozhou@ualberta.ca>)
2. Jiwei Liu, University of Pittsburgh, aixueer4ever@gmail.com (<mailto:aixueer4ever@gmail.com>)
3. Xueer Chen, University of Pittsburgh, xuer.chen.human@gmail.com
(<mailto:xuer.chen.human@gmail.com>)

I have used the code from the pdf provided at: <https://www.kaggle.com/c/malware-classification/discussion/13897#latest-105551> (<https://www.kaggle.com/c/malware-classification/discussion/13897#latest-105551>)

I have also used this as a reference: <https://github.com/dchad/malware-detection> (<https://github.com/dchad/malware-detection>)

Convert the ASM files and BYTE files to images. (Original Image Dimensions are retained here)

In [19]:

```

1 import numpy as np
2 import os
3 import codecs
4 import imageio
5 import array
6 from datetime import datetime as dt
7
8 if not os.path.isdir("asmFileImages"):
9     os.mkdir("asmFileImages")
10 if not os.path.isdir("byteFileImages"):
11     os.mkdir("byteFileImages")
12
13 asmfile_list=os.listdir("asmFiles/")
14 bytefile_list=os.listdir("byteFiles/")
15
16 #This function will generate images from ASM files and Byte files.
17 def get_images_from_text(filename_list, save_folder):
18     for filename in tqdm(filename_list):
19
20         if(filename.endswith("asm")):
21             file = codecs.open("asmFiles/"+filename, 'rb')
22             file_size = os.path.getsize("asmFiles/"+filename)
23         else:
24             file = open("byteFiles/"+filename, 'rb')
25             file_size = os.path.getsize("byteFiles/"+filename)
26
27         file_width = int(file_size**0.5)
28         rem = file_size%file_width
29         img_array = array.array('B')
30         img_array.fromfile(file,file_size-rem)
31         file.close()
32         img_arr_final = np.reshape(img_array[:file_width * file_width], (file_width, file_width))
33         img_arr_final = np.uint8(img_arr_final)
34         imageio.imsave(save_folder+'/' + filename.split(".")[0] + '.png',img_arr_final)

```

In []:

```

1 #Convert ASM files to Images
2 get_images_from_text(asmfile_list,'asmFileImages')

```

In []:

```

1 #Convert Byte files to Images
2 get_images_from_text(bytefile_list,'byteFileImages')

```

Resize the ASM and Byte image files to 256x256 dimensions. Keep the original images untouched.

Own Reference: <https://stackoverflow.com/questions/44650888/resize-an-image-without-distortion-opencv/49208362#49208362> (<https://stackoverflow.com/questions/44650888/resize-an-image-without-distortion-opencv/49208362#49208362>)

```
In [4]: 1 import cv2
        2 import os
        3 import imageio
        4 from tqdm import tqdm
        5 import random
        6
        7 #Take any image as input, resize it to 256x256, return the resized image
        8 def resize_image(image):
        9     resized_image = cv2.resize(image, (256,256), interpolation = cv2.INTER_AR
       10     return resized_image
```

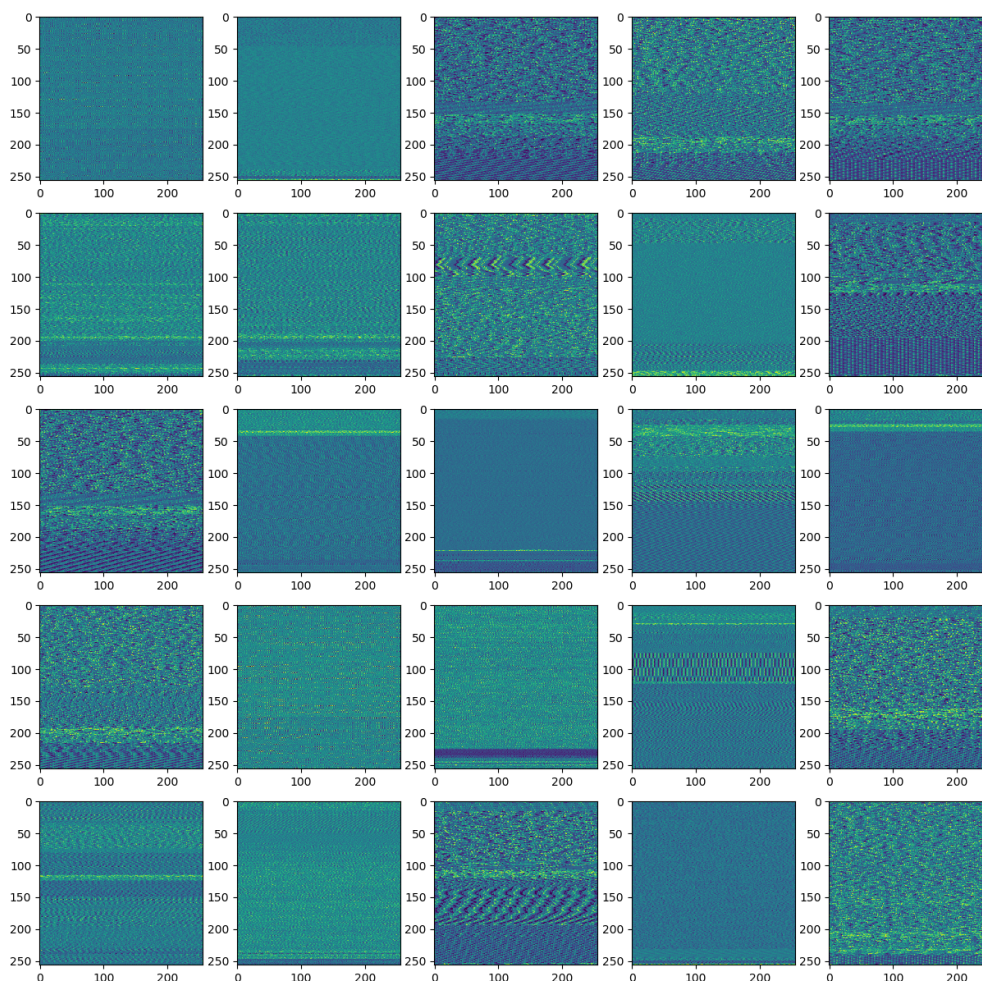
Display some sample images that we have obtained from ASM files.

```

In [6]: 1 image_dir = "asmFileImages/"
        2 filenames=random.sample(os.listdir(image_dir),26)
        3
        4 #Display 25 images from ASM files
        5 plt.figure(figsize=(15,15))
        6 for i in range(1,len(filenames)):
        7     row = i
        8     image = imageio.imread(image_dir+filenames[i]) #Image(filename=image_dir+
        9     image = resize_image(image)
       10     plt.subplot(5,5,row)
       11     plt.imshow(image)
       12 plt.show()

```

<IPython.core.display.Javascript object>



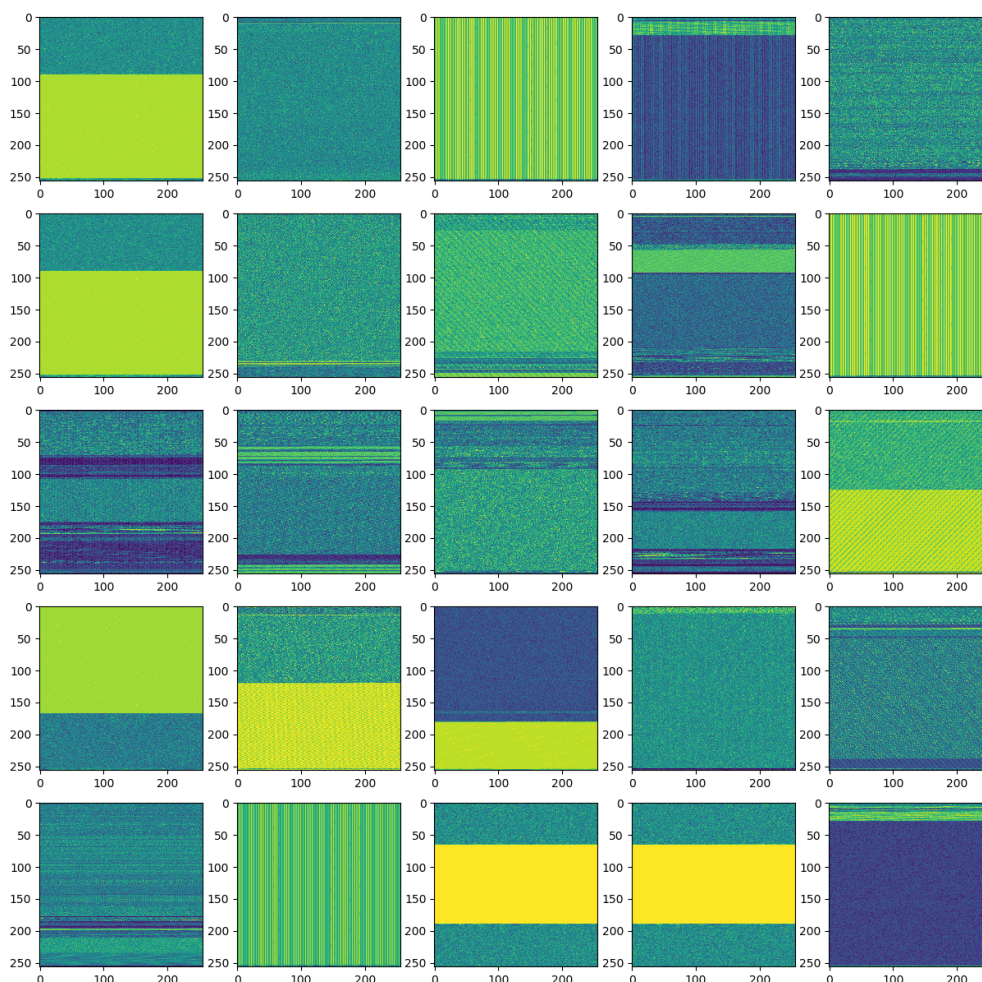
Display some sample images that we have obtained from Byte files.

```

In [5]: 1 image_dir = "byteFileImages/"
        2 filenames=random.sample(os.listdir(image_dir),26)
        3
        4 #Display 25 images from BYTE files
        5 plt.figure(figsize=(15,15))
        6 for i in range(1,len(filenames)):
        7     row = i
        8     image = imageio.imread(image_dir+filenames[i]) #Image(filename=image_dir+
        9     image = resize_image(image)
       10     plt.subplot(5,5,row)
       11     plt.imshow(image)
       12 plt.show()

```

<IPython.core.display.Javascript object>



Load the image, convert it into a numpy array and take the first 800 pixels from each image.

```
In [10]: 1 import os
2 from scipy.misc import imread
3 import imageio
4 from tqdm import tqdm
5
6 import warnings
7 warnings.filterwarnings("ignore")
```

Get first 800 pixel information from ASM File Images

```
In [50]: 1 file_lists=os.listdir('asmFileImages/')
2
3 with open("features/asm_image_df.csv", mode='w') as asm_image_df: #file_lists
4     asm_image_df.write(','.join(map(str, ["ID"]+["Pixel{}".format(i) for i in
5     asm_image_df.write('\n')
6
7     for image in tqdm(file_lists):
8         file_id=image.split(".")[0]
9         image_array=imageio.imread("asmFileImages/"+image) #This will contain
10        image_array=image_array.flatten()[:800] #Taking the first 1000 pixels
11        asm_image_df.write(','.join(map(str, [file_id]+list(image_array))))
12        asm_image_df.write('\n') #Write to the dataframe
```

100%|██████████| 10868/10868 [20:59<00:00, 20.28it/s]

Display the ASM Image dataframe

```
In [51]: 1 asm_image_df=pd.read_csv("features/asm_image_df.csv")
2 asm_image_df.head()
```

```
Out[51]:
```

	ID	Pixel0	Pixel1	Pixel2	Pixel3	Pixel4	Pixel5	Pixel6	Pixel7	Pixel8	...
0	01azqd4InC7m9JpocGv5	72	69	65	68	69	82	58	48	48	...
1	01IsoiSMh5gxyDYTI4CB	46	116	101	120	116	58	48	48	52	...
2	01jsnpXSAIgw6aPeDxrU	72	69	65	68	69	82	58	48	48	...
3	01kcPWA9K2BOxQeS5Rju	72	69	65	68	69	82	58	49	48	...
4	01SuzwMJEIXsK7A8dQbl	72	69	65	68	69	82	58	48	48	...

5 rows × 801 columns



```
In [54]: 1 del(asm_image_df)
```

Get first 800 pixel information from Byte File Images


```
In [52]: 1 file_lists=os.listdir('byteFileImages/')
2
3 with open("features/byte_image_df.csv", mode='w') as byte_image_df: #file_list
4     byte_image_df.write(','.join(map(str, ["ID"]+["Pixel{}".format(i) for i in range(19)])))
5     byte_image_df.write('\n')
6
7     for image in tqdm(file_lists):
8         file_id=image.split('.')[0]
9         image_array=imageio.imread("byteFileImages/"+image) #This will contain the image
10        image_array=image_array.flatten()[:800] #Taking the first 1000 pixels
11        byte_image_df.write(','.join(map(str, [file_id]+list(image_array))))
12        byte_image_df.write('\n') #Write to the dataframe
```

100%|██████████| 10868/10868 [10:51<00:00, 16.67it/s]

Display the BYTE image dataframe

```
In [53]: 1 byte_image_df=pd.read_csv("features/byte_image_df.csv")
2         byte_image_df.head()
```

Out[53]:

	ID	Pixel0	Pixel1	Pixel2	Pixel3	Pixel4	Pixel5	Pixel6	Pixel7	Pixel8	...
0	01azqd4InC7m9JpocGv5	69	56	32	48	66	32	48	48	32	...
1	01lsoiSMh5gxyDYTI4CB	67	55	32	48	49	32	50	52	32	...
2	01jsnpXSAIgw6aPeDxrU	67	66	32	67	66	32	67	66	32	...
3	01kcPWA9K2BOxQeS5Rju	54	65	32	70	70	32	54	56	32	...
4	01SuzwMJElXsK7A8dQbl	65	52	32	65	67	32	52	65	32	...

5 rows × 801 columns



```
In [55]: 1 del(byte_image_df)
```

5.7 Extract CNN Codes from the image files using Pretrained VGG-16 on ImageNet data

```
In [4]: 1 from datetime import datetime as dt
2 from keras.preprocessing.image import ImageDataGenerator
3 from keras import applications
4 import numpy as np
```

Using TensorFlow backend.

Get CNN codes for ASM Files

```

In [4]: 1 global_start=dt.now()
        2
        3 #We will take 256 x 256 image dimensions, because the least dimension that wa
        4 img_width, img_height = 256, 256
        5 batch_size=44
        6
        7 #Declaration of parameters needed for training and validation
        8 data_dir = 'finalASMimages'
        9 global bottleneck_features
       10
       11 #Get the bottleneck features by Weights.T * Xi
       12 def save_bottleneck_features_asm():
       13     datagen = ImageDataGenerator(rescale=1./255)
       14
       15     #Load the pre trained VGG16 model from Keras, we will initialize only the
       16     model = applications.VGG16(include_top=False, weights='imagenet')
       17
       18     generator_tr = datagen.flow_from_directory(data_dir,
       19                                              target_size=(img_width, img_height),
       20                                              batch_size=batch_size,
       21                                              class_mode=None, #class_mode=None
       22                                              shuffle=False) #We won't shuffle
       23     nb_train_samples = len(generator_tr.filesnames)
       24     bottleneck_features = model.predict_generator(generator_tr, nb_train_samples)
       25
       26     np.save('features/final_features/asm_files_bottleneck_features.npy', bottleneck_features)
       27     print("Got the bottleneck features in time: ", dt.now()-global_start)
       28
       29 save_bottleneck_features_asm()

```

Found 10868 images belonging to 1 classes.

Got the bottleneck features in time: 1:41:29.963485

Get CNN codes for Byte Files

```

In [6]: 1 global_start=dt.now()
        2
        3 #We will take 256 x 256 image dimensions, because the least dimension that wa
        4 img_width, img_height = 256, 256
        5 batch_size=44
        6
        7 #Declaration of parameters needed for training and validation
        8 data_dir = "finalBYTEimages"
        9
        10 #Get the bottleneck features by Weights.T * Xi
        11 def save_bottleneck_features_byte():
        12     datagen = ImageDataGenerator(rescale=1./255)
        13
        14     #Load the pre trained VGG16 model from Keras, we will initialize only the
        15     model = applications.VGG16(include_top=False, weights='imagenet')
        16
        17     generator_tr = datagen.flow_from_directory(data_dir,
        18                                              target_size=(img_width, img_height),
        19                                              batch_size=batch_size,
        20                                              class_mode=None, #class_mode=None
        21                                              shuffle=False) #We won't shuffle
        22     nb_train_samples = len(generator_tr.filesnames)
        23     bottleneck_features = model.predict_generator(generator_tr, nb_train_samples)
        24
        25     np.save('features/final_features/byte_files_bottleneck_features.npy', bottleneck_features)
        26     print("Got the bottleneck features in time: ", dt.now()-global_start)
        27
        28 save_bottleneck_features_byte()

```

Found 10868 images belonging to 1 classes.

Got the bottleneck features in time: 1:52:11.927833

Prediction using only CNN Codes

```

In [2]: 1 import numpy as np
        2 asm_bottleneck = np.load("features/final_features/asm_files_bottleneck_features.npy")
        3 asm_bottleneck.shape

```

Out[2]: (10868, 8, 8, 512)

```

In [1]: 1 import numpy as np
        2 asm_bottleneck = np.load("features/final_features/asm_files_bottleneck_features.npy")
        3 byte_bottleneck = np.load("features/final_features/byte_files_bottleneck_features.npy")
        4 bot_feats = np.hstack((asm_bottleneck, byte_bottleneck))
        5 np.save('features/final_features/bot_feats.npy', bot_feats)
        6
        7 import gc
        8 gc.collect()

```

Out[1]: 16

TODO: Trigrams + Fourgrams + Fivegrams + 6Grams Feature extraction

from byte + ASM files

TODO: Get count of malicious words

```
In [114]: 1 #TRIGRAMS
2 """tokens="00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,
3 tokens=tokens.split(",")
4
5 #This function will return all the possible 257*257*257 combinations of trigr
6 def get_trigram_tokens(tokens):
7     sent=""
8     byte_trigram_vocab=[]
9     for i in tqdm(range(len(tokens))):
10         for j in range(len(tokens)):
11             for k in range(len(tokens)):
12                 trigram=tokens[i]+" "+tokens[j]+" "+tokens[k]
13                 sent=sent+trigram+", "
14                 byte_trigram_vocab.append(trigram)
15     return byte_trigram_vocab
16
17 byte_trigram_vocab=get_trigram_tokens(tokens) #This will contain all the poss
18
19 #Extract the trigram features and save the dataframe containing bi-gram featu
20 """vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False
21 file_lists=os.listdir('byteFiles')
22 bytetrigram_vector = scipy.sparse.csr_matrix((1, len(byte_trigram_vocab))) #f
23
24 byte_trigram_df=pd.DataFrame(columns=vectorizer.get_feature_names())
25
26 with open("features/byte_trigram_df.csv", mode='w') as byte_trigram_df:
27     byte_trigram_df.write(','.join(map(str, vectorizer.get_feature_names()))
28     byte_trigram_df.write('\n')
29     for _, file in tqdm(enumerate(file_lists)):
30         f = open('byteFiles/' + file)
31         corpus=[f.read().replace('\n', ' ').lower()] #This will contain all t
32         trigrams=vectorizer.fit_transform(corpus) #This will return a sparse
33         row=scipy.sparse.csr_matrix(trigrams).toarray() #Update each row of t
34         byte_trigram_df.write(','.join(map(str, row[0])))
35         byte_trigram_df.write('\n')
36         f.close() """
37
38 print("TODO: Trigrams Feature extraction from byte files ")
```

TODO: Trigrams Feature extraction from byte files

6.0 Getting the most important features using using SelectKBest with Chi-Square Test

```
In [ ]: 1 from sklearn.feature_selection import SelectKBest, chi2
2
3 if not os.path.isdir("features/final_features"):
4     os.mkdir("features/final_features")
5
6 if not os.path.isdir("features/feature_score"):
7     os.mkdir("features/feature_score")
8
9 if not os.path.isdir("features/final_features"):
10     os.mkdir("features/final_features")
11
12 #Load the class labels for training with random forest feature selector
13 with open('features/class_labels.pkl', 'rb') as file:
14     labels=pkl.load(file)
```

6.1 Getting the 50 most important features for Registers bigrams

```
In [7]: 1 #Load the non normalized bigrams dataset
2 X=pd.read_csv('features/registers_bigram_df.csv')
3 y=labels
4 X.head()
```

Out[7]:

		ID	edx edx	edx esi	edx eax	edx ebx	edx ecx	edx edi	edx ebp	edx esp	edx eip	...	esp eip	eip edx	eip esi	eip eax
0	01azqd4lnC7m9JpocGv5		55	97	211	33	130	40	15	0	0	...	0	0	0	0
1	01lsoiSMh5gxyDYTI4CB		31	12	72	15	68	23	5	8	0	...	0	0	0	0
2	01jsnpXSAIgw6aPeDxrU		66	0	72	1	50	0	5	0	0	...	0	0	0	0
3	01kcPWA9K2BOxQeS5Rju		1	1	4	4	0	0	0	0	0	...	0	0	0	0
4	01SuzwMJEIXsK7A8dQbl		8	1	374	0	12	1	8	3	0	...	0	0	0	0

5 rows × 82 columns



```

In [9]: 1 #Get the best 50 features using SelectKBest. Save the feature scores along wi
2 kbest_object=SelectKBest(score_func=chi2, k=50)
3 best_features=kbest_object.fit(X.drop("ID", axis=1),y)
4 df_scores=pd.DataFrame(best_features.scores_)
5 df_columns=pd.DataFrame(X.columns)
6 feature_score_df=pd.concat([df_columns,df_scores],axis=1)
7 feature_score_df.columns=["Feature_Name", "Feature_Score"]
8
9 #Let's Look at the top 50 features along with their scores + Save the feature
10 feature_score_df=feature_score_df.nlargest(50,"Feature_Score")
11 feature_score_df.to_csv("features/feature_score/registers_bigram_df.csv", ind
12 feature_score_df.head(5)

```

```

Out[9]:

```

	Feature_Name	Feature_Score
47	edi esi	27901.216442
20	eax esi	25442.443807
23	eax ecx	24936.321815
21	eax eax	20482.605329
11	esi esi	19799.204224

```

In [10]: 1 #Get the first 50 feature names in a list
2 top_50_feats=list(feature_score_df["Feature_Name"])
3
4 #Get the dataframe containing the top 50 features
5 reduced_features=pd.concat([X["ID"],X[top_50_feats]], axis=1)
6
7 #Save the dataframe containing the top 50 features
8 reduced_features.to_csv("features/final_features/top50_registers_bigram_df.cs

```

```

In [ ]: 1 del(X, feature_score_df, reduced_features)

```

6.2 Getting the 300 most important features from Registers trigrams

```
In [14]: 1 #Load the non normalized bigrams dataset
2 X=pd.read_csv('features/registers_trigram_df.csv')
3 y=labels
4 X.head()
```

Out[14]:

		edx edx	edx esi	edx eax	edx ebx	edx ecx	edx edi	edx ebp	edx esp	edx eip	...	eip esp	eip eip	eip esi	eip eax
0	01azqd4lnC7m9JpocGv5	7	2	27	6	9	3	1	0	0	...	0	0	0	0
1	01lsoiSMh5gxyDYTI4CB	6	1	10	0	12	2	0	0	0	...	0	0	0	0
2	01jsnpXSAIgw6aPeDxrU	23	0	18	0	23	0	2	0	0	...	0	0	0	0
3	01kcPWA9K2BOxQeS5Rju	0	0	0	1	0	0	0	0	0	...	0	0	0	0
4	01SuzwMJEIXsK7A8dQbl	3	0	1	0	2	0	2	0	0	...	0	0	0	0

5 rows × 730 columns

```
In [15]: 1 #Get the best 300 features using SelectKBest. Save the feature scores along w
2 kbest_object=SelectKBest(score_func=chi2, k=300)
3 best_features=kbest_object.fit(X.drop("ID", axis=1),y)
4 df_scores=pd.DataFrame(best_features.scores_)
5 df_columns=pd.DataFrame(X.columns)
6 feature_score_df=pd.concat([df_columns,df_scores],axis=1)
7 feature_score_df.columns=["Feature_Name", "Feature_Score"]
8
9 #Let's Look at the top 300 features along with their scores + Save the featur
10 feature_score_df=feature_score_df.nlargest(300,"Feature_Score")
11 feature_score_df.to_csv("features/feature_score/registers_trigram_df.csv", in
12 feature_score_df.head(5))
```

Out[15]:

	Feature_Name	Feature_Score
637	esp esp ebp	19633.605891
50	edx edi ecx	11379.497387
189	eax eax eip	11371.244132
182	eax eax esi	11124.759127
410	edi edx ecx	11075.744257

```
In [16]: 1 #Get the first 300 feature names in a list
2 top_300_feats=list(feature_score_df["Feature_Name"])
3
4 #Get the dataframe containing the top 300 features
5 reduced_features=pd.concat([X["ID"],X[top_300_feats]], axis=1)
6
7 #Save the dataframe containing the top 300 features
8 reduced_features.to_csv("features/final_features/top300_registers_trigram_df.
```

```
In [ ]: 1 del(X, feature_score_df, reduced_features)
```

6.3 Getting the 500 most important features from Opcodes bigrams

```
In [17]: 1 #Load the non normalized bigrams dataset
2 X=pd.read_csv('features/opcodes_bigram_df.csv')
3 y=labels
4 X.head()
```

Out[17]:

	ID	jmp jmp	jmp mov	jmp retf	jmp push	jmp pop	jmp xor	jmp retn	jmp nop	jmp sub	...	movzx cmp	movzx call
0	01azqd4InC7m9JpocGv5	440	192	0	6	0	17	0	0	24	...	0	0
1	01lsoiSMh5gxyDYTI4CB	0	32	0	3	1	3	1	0	0	...	1	0
2	01jsnpXSAIgw6aPeDxrU	0	0	0	0	0	0	0	0	0	...	12	0
3	01kcPWA9K2BOxQeS5Rju	0	5	0	1	0	2	1	0	0	...	0	0
4	01SuzwMJEIXsK7A8dQbl	5	57	1	4	1	1	0	0	0	...	11	0

5 rows × 677 columns



```

In [18]: 1 #Get the best 500 features using SelectKBest. Save the feature scores along w
2 kbest_object=SelectKBest(score_func=chi2, k=500)
3 best_features=kbest_object.fit(X.drop("ID", axis=1),y)
4 df_scores=pd.DataFrame(best_features.scores_)
5 df_columns=pd.DataFrame(X.columns)
6 feature_score_df=pd.concat([df_columns,df_scores],axis=1)
7 feature_score_df.columns=["Feature_Name", "Feature_Score"]
8
9 #Let's Look at the top 500 features along with their scores + Save the featur
10 feature_score_df=feature_score_df.nlargest(500,"Feature_Score")
11 feature_score_df.to_csv("features/feature_score/opcodes_bigram_df.csv", index
12 feature_score_df.head(5)

```

```

Out[18]:

```

	Feature_Name	Feature_Score
200	nop call	284660.375541
27	mov jmp	276173.423650
475	shl retn	225563.920296
38	mov add	111179.349251
313	imul jmp	106299.863432

```

In [19]: 1 #Get the first 500 feature names in a List
2 top_500_feats=list(feature_score_df["Feature_Name"])
3
4 #Get the dataframe containing the top 500 features
5 reduced_features=pd.concat([X["ID"],X[top_500_feats]], axis=1)
6
7 #Save the dataframe containing the top 500 features
8 reduced_features.to_csv("features/final_features/top500_opcodes_bigram_df.csv")

```

```

In [ ]: 1 del(X, feature_score_df, reduced_features)

```

6.4 Getting the 500 most important features from Opcodes trigrams

```
In [20]: 1 #Load the non normalized bigrams dataset
2 X=pd.read_csv('features/opcodes_trigram_df.csv')
3 y=labels
4 X.head()
```

```
Out[20]:
```

	ID	jmp jmp	jmp mov	jmp retf	jmp push	jmp pop	jmp xor	jmp retn	jmp nop	jmp sub	...	movzx movzx cmp	movzx movzx call
0	01azqd4lnC7m9JpocGv5	437	0	0	0	0	1	0	0	0	...	0	0
1	01lsoiSMh5gxyDYTI4CB	0	0	0	0	0	0	0	0	0	...	0	0
2	01jsnpXSAIgw6aPeDxrU	0	0	0	0	0	0	0	0	0	...	1	0
3	01kcPWA9K2BOxQeS5Rju	0	0	0	0	0	0	0	0	0	...	0	0
4	01SuzwMJEIXsK7A8dQbl	2	1	1	1	0	0	0	0	0	...	1	0

5 rows × 17577 columns

```
In [21]: 1 #Get the best 500 features using SelectKBest. Save the feature scores along w
2 kbest_object=SelectKBest(score_func=chi2, k=500)
3 best_features=kbest_object.fit(X.drop("ID", axis=1),y)
4 df_scores=pd.DataFrame(best_features.scores_)
5 df_columns=pd.DataFrame(X.columns)
6 feature_score_df=pd.concat([df_columns,df_scores],axis=1)
7 feature_score_df.columns=["Feature_Name", "Feature_Score"]
8
9 #Let's Look at the top 500 features along with their scores + Save the featur
10 feature_score_df=feature_score_df.nlargest(500,"Feature_Score")
11 feature_score_df.to_csv("features/feature_score/opcodes_trigram_df.csv", inde
12 feature_score_df.head(5)
```

```
Out[21]:
```

	Feature_Name	Feature_Score
12368	shl nop call	361070.460585
703	mov mov jmp	239711.417051
8312	imul nop call	201312.047619
5207	nop shl retn	201052.627929
10496	shr xchg call	200793.864581

```
In [22]: 1 #Get the first 500 feature names in a list
2 top_500_feats=list(feature_score_df["Feature_Name"])
3
4 #Get the dataframe containing the top 500 features
5 reduced_features=pd.concat([X["ID"],X[top_500_feats]], axis=1)
6
7 #Save the dataframe containing the top 500 features
8 reduced_features.to_csv("features/final_features/top500_opcodes_trigram_df.csv")
```

```
In [ ]: 1 del(X, feature_score_df, reduced_features)
```

6.5 Getting the 1000 most important features from Byte bigrams

```
In [5]: 1 #Load the non normalized bigrams dataset
2 X=dd.read_csv('features/byte_bigram_df.csv', sample=25600000)
3 y=labels
4 X.head()
```

Out[5]:

	ID	00 00	00 01	00 02	00 03	00 04	00 05	00 06	00 07	00 08	...	?? f7	?? f8
0	01azqd4InC7m9JpocGv5	274425	1269	1029	1469	1227	1144	1437	1263	1174	...	0	0
1	01lsoiSMh5gxyDYTI4CB	21075	752	73	48	175	12	10	11	42	...	0	0
2	01jsnpXSAlg6aPeDxrU	16798	596	159	144	513	595	557	146	528	...	0	0
3	01kcPWA9K2BOxQeS5Rju	10417	225	61	69	114	40	25	22	63	...	0	0
4	01SuzwMJElXsK7A8dQbl	16271	62	22	126	9	11	3	5	11	...	0	0

5 rows × 66050 columns




```

In [6]: 1 #Get the best 1000 features using SelectKBest. Save the feature scores along
2 kbest_object=SelectKBest(score_func=chi2, k=1000)
3 best_features=kbest_object.fit(X.drop("ID", axis=1),y)
4 df_scores=pd.DataFrame(best_features.scores_)
5 df_columns=pd.DataFrame(X.columns)
6 feature_score_df=pd.concat([df_columns,df_scores],axis=1)
7 feature_score_df.columns=["Feature_Name","Feature_Score"]
8
9 #Let's look at the top 1000 features along with their scores + Save the featu
10 feature_score_df=feature_score_df.nlargest(1000,"Feature_Score")
11 feature_score_df.to_csv("features/feature_score/byte_bigram_df.csv", index=No
12 feature_score_df.head(5)

```

```

Out[6]:

```

	Feature_Name	Feature_Score
66048	?? ff	1.269512e+07
22258	56 9b	1.150512e+07
22251	56 94	1.143470e+07
22319	56 d8	1.136046e+07
38379	95 55	1.133120e+07

```

In [ ]: 1 #Get the first 1000 feature names in a List
2 top_1000_feats=list(feature_score_df["Feature_Name"])
3
4 #Get the dataframe containing the top 1000 features
5 reduced_features=dd.concat([X["ID"],X[top_1000_feats]], axis=1)
6
7 #Save the dataframes containing the top 1000 features
8 reduced_features.to_csv("features/final_features/temp/top1000_byte_bigram_df*

```

```

In [ ]: 1 del(X, feature_score_df, reduced_features)

```

```

In [ ]: 1 #Concatenate all the dataframes and save it as "top1000_byte_bigram_df.csv"
2 file_lists=os.listdir('features/final_features/temp/')
3 final_byte_bigram_df=pd.read_csv('features/final_features/temp/'+file_lists[0])
4
5 for i in range(1,len(file_lists)):
6     df = pd.read_csv('features/final_features/temp/'+file_lists[i])
7     final_byte_bigram_df=pd.concat([final_byte_bigram_df,df], axis=0)
8
9 final_byte_bigram_df.to_csv("features/final_features/top1000_byte_bigram_df.c

```

```

In [ ]: 1 del(final_byte_bigram_df, file_lists)

```

Code for bigram creation using sparse matrix.

```
In [ ]: 1 """
2
3 #Extract the byte bigram features and save the dataframe containing bi-gram f
4 vectorizer = CountVectorizer(tokenizer=lambda x: x.split(),lowercase=False,ng
5 file_lists=os.listdir('byteFiles/'))
6 byte_bigram_data=sp.csr_matrix((1, 66049))
7
8 file = open('byteFiles/' + file_lists[0])
9 corpus = [file.read().replace('\n', ' ').lower()] #Read the file content. Rep
10 bigrams=vectorizer.fit_transform(corpus)
11 byte_bigram_data = sp.csr_matrix(bigrams)
12
13 for i in tqdm(range(1,10868)): ##The open() function returns a file object, w
14     file = open('byteFiles/' + file_lists[i])
15     corpus = [file.read().replace('\n', ' ').lower()] #Read the file content.
16     bigrams=vectorizer.fit_transform(corpus)
17     byte_bigram_data = sp.vstack((byte_bigram_data,bigrams))
18     file.close()
19 """
```

Use this code section in future for extremely large datasets

```

In [91]: 1 #This code snippet iteratively builds a dataframe containing the most feature
2 """
3
4 final_score_df=pd.DataFrame()
5
6 start=0
7 end=100
8 count=1
9 while(end<66051):
10     #Get the col names we want for this iteration.
11     col=columns[start:end]
12     X=df[col]
13
14     best_features=SelectKBest(score_func=chi2, k=X.shape[1])
15     fit=best_features.fit(X,y)
16
17     df_scores=pd.DataFrame(fit.scores_)
18     df_columns=pd.DataFrame(X.columns)
19
20     feature_score=pd.concat([df_columns,df_scores],axis=1)
21     feature_score.columns=["Feature Name","Feature Score"]
22
23
24     start=start+100
25     end=end+100
26
27     if(end>66049):
28         end=66050
29
30     final_score_df=pd.concat([final_score_df,feature_score],ignore_index=True)
31
32
33     del(X)
34
35     if(start>end):
36         break
37
38
39     print("Iteration {} completed".format(count))
40     count+=1
41     gc.collect()
42     final_score_df.nlargest(50,"Score")
43
44     #Load the class labels for training with random forest feature selector
45     with open('features/class_labels.pkl', 'rb') as file:
46         labels=pkl.load(file)
47
48     #Load the sparse data as a sparse matrix
49     byte_bigram_stage1=normalize(scipy.sparse.load_npz('features/sparse/byte_bigram_stage1.npz'))
50
51
52 #TODO: Correct the issues and parallalize the dask dataframes
53
54 """import pandas as pd
55
56 from dask.distributed import Client

```

```

57 from sklearn.tree import DecisionTreeClassifier
58 from sklearn.ensemble import RandomForestClassifier
59 import dask.dataframe as dd
60 from sklearn.externals.joblib import parallel_backend
61
62 client = Client() # start a local Dask client
63 byte_bigram_df=dd.read_csv("features/byte_bigram_df.csv", sample=2560000)
64 labels=pd.read_csv('features/byte_unigrams_df.csv')['Class']
65
66 with parallel_backend('dask'):
67     #Scikit learn code here
68     rf_clf=DecisionTreeClassifier(random_state=0)
69     rf_clf.fit(byte_bigram_df,labels)
70     feature_index=np.argsort(rf_clf.feature_importances_)[::-1]
71     most_imp_feat_idx=feature_index[:5000]
72
73 most_imp_feat_idx"""
74
75 print("Execute this code for iteratively getting feature scores..")

```

Execute this code for iteratively getting feature scores..

6.6 Use the CNN Codes features

Ignore the bottleneck features for now. The final dimensions of the the bottleneck features is 20 less than the actual number of images. TODO: Debug this.

```

In [2]: 1 X_cnn = np.load('final_features_cnn')
        2 Y_cnn = pd.read_csv('final_features_cnn')['Class']
        3 X_cnn.shape

```

Out[2]: (10868, 65536)

```

In [3]: 1 #Split the data into test and train by maintaining same distribution of output
        2 X_train, X_test, y_train, y_test = train_test_split(X_cnn,Y_cnn, stratify=y,
        3 #Split the train data into train and cross validation by maintaining same distribution
        4 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y,

```

```

In [5]: 1 from sklearn.model_selection import RandomizedSearchCV
        2 from sklearn.ensemble import RandomForestClassifier
        3 from sklearn.calibration import CalibratedClassifierCV
        4 from xgboost import XGBClassifier

```

```

In [ ]: 1 estimators=[10,50,100,500,1000,2000,3000,4000]
2 cv_log_error_array=[]
3 for i in estimators:
4     x_cfl=XGBClassifier(n_estimators=i, n_jobs=-1)
5     x_cfl.fit(X_train , y_train )
6     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
7     sig_clf.fit(X_train , y_train )
8     predict_y = sig_clf.predict_proba(X_cv )
9     cv_log_error_array.append(log_loss(y_cv , predict_y, labels=x_cfl.classes_))
10
11 for i in range(len(cv_log_error_array)):
12     print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_array[i])
13
14
15 best_estimators = np.argmin(cv_log_error_array)
16
17 fig, ax = plt.subplots()
18 ax.plot(estimators, cv_log_error_array,c='g')
19 for i, txt in enumerate(np.round(cv_log_error_array,3)):
20     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_array[i]))
21 plt.grid()
22 plt.title("Cross Validation Error for each estimators")
23 plt.xlabel("Estimators i's")
24 plt.ylabel("Error measure")
25 plt.show()
26
27 x_cfl=XGBClassifier(n_estimators=estimators[best_estimators],nthread=-1,n_jobs=-1)
28 x_cfl.fit(X_train , y_train ,verbose=True)
29 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
30 sig_clf.fit(X_train , y_train )
31
32 predict_y = sig_clf.predict_proba(X_train )
33 print ('For values of best estimators = ', estimators[best_estimators], "The error is", cv_log_error_array[best_estimators])
34 predict_y = sig_clf.predict_proba(X_cv )
35 print('For values of best estimators = ', estimators[best_estimators], "The error is", cv_log_error_array[best_estimators])
36 predict_y = sig_clf.predict_proba(X_test )
37 print('For values of best estimators = ', estimators[best_estimators], "The error is", cv_log_error_array[best_estimators])

```

7.0 Stack all the features together and save it as one single large feature dataframe

```
In [3]: 1 #Load the byte unigrams dataframe
2 byte_features_with_size = pd.read_csv("features/byte_features_with_size.csv")
3
4 #Load the ASM unigrams dataframe
5 asm_features_with_size = pd.read_csv("features/asm_features_with_size.csv").d
6
7 #Load the byte bigram dataframe with reduced features
8 top1000_byte_bigram_df = pd.read_csv("features/final_features/top1000_byte_bi
9
10 #Load the ASM Opcodes bigram dataframe with reduced features
11 top500_opcodes_bigram_df = pd.read_csv("features/final_features/top500_opcode
12
13 #Load the ASM Registers bigram dataframe with reduced features
14 top50_registers_bigram_df = pd.read_csv("features/final_features/top50_regist
15
16 #Load the ASM Opcodes trigram dataframe with reduced features
17 top500_opcodes_trigram_df = pd.read_csv("features/final_features/top500_opcod
18
19 #Load the ASM Registers trigram dataframe with reduced features
20 top300_registers_trigram_df = pd.read_csv("features/final_features/top300_reg
21
22 #Load the ASM Image dataframe
23 asm_image_df = pd.read_csv("features/asm_image_df.csv")
24
25 #Load the Byte Image dataframe
26 byte_image_df = pd.read_csv("features/byte_image_df.csv")
27
28 #Join all the dataframes together on "ID" and save it as one single dataframe
29 final_features_df=byte_features_with_size
30 dfs = [asm_features_with_size,
31         top1000_byte_bigram_df,
32         top500_opcodes_bigram_df,
33         top50_registers_bigram_df,
34         top500_opcodes_trigram_df,
35         top300_registers_trigram_df,
36         asm_image_df,
37         byte_image_df]
38
39 for df in dfs:
40     final_features_df=pd.merge(final_features_df,df,on="ID",how="left")
41
42 final_features_df.to_csv("features/final_features/final_features_df.csv", ind
```

```

In [7]: 1 #Load the byte unigrams dataframe
2 byte_features_with_size = pd.read_csv("features/byte_features_with_size.csv")
3
4 #Load the ASM unigrams dataframe
5 asm_features_with_size = pd.read_csv("features/asm_features_with_size.csv").d
6
7 #Load the byte bigram dataframe with reduced features
8 top1000_byte_bigram_df = pd.read_csv("features/final_features/top1000_byte_bi
9
10 #Load the ASM Opcodes bigram dataframe with reduced features
11 top500_opcodes_bigram_df = pd.read_csv("features/final_features/top500_opcode
12
13 #Load the ASM Registers bigram dataframe with reduced features
14 top50_registers_bigram_df = pd.read_csv("features/final_features/top50_regist
15
16 #Load the ASM Opcodes trigram dataframe with reduced features
17 top500_opcodes_trigram_df = pd.read_csv("features/final_features/top500_opcod
18
19 #Load the ASM Registers trigram dataframe with reduced features
20 top300_registers_trigram_df = pd.read_csv("features/final_features/top300_reg
21
22 #Join all the dataframes together on "ID" and save it as one single dataframe
23 final_features_df=byte_features_with_size
24 dfs = [asm_features_with_size,
25         top1000_byte_bigram_df,
26         top500_opcodes_bigram_df,
27         top50_registers_bigram_df,
28         top500_opcodes_trigram_df,
29         top300_registers_trigram_df]
30
31 for df in dfs:
32     final_features_df=pd.merge(final_features_df,df,on="ID",how="left")
33
34 final_features_df.to_csv("features/final_features/final_features_df.csv", ind

```

```

In [8]: 1 import gc
2 del(final_features_df,byte_features_with_size,asm_features_with_size,top1000_
3 gc.collect()

```

Out[8]: 219

8.1 Train Test Split. 64% Train, 16% Cross Validation, 20% Test

```

In [ ]: 1 X = pd.read_csv('final_features_df.csv')
2 #normalizing it
3
4 X = normalize(X)

```

Load the normalized DF along with the class labels

```
In [4]: 1 X=pd.read_csv("features/final_features/final_features_df_normalized.csv").fil  
2 y=pd.read_csv("features/final_features/final_features_df.csv")["Class"]
```

```
In [5]: 1 #Split the data into test and train by maintaining same distribution of outpu  
2 X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, test_si  
3  
4 #Split the train data into train and cross validation by maintaining same dis  
5 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_
```

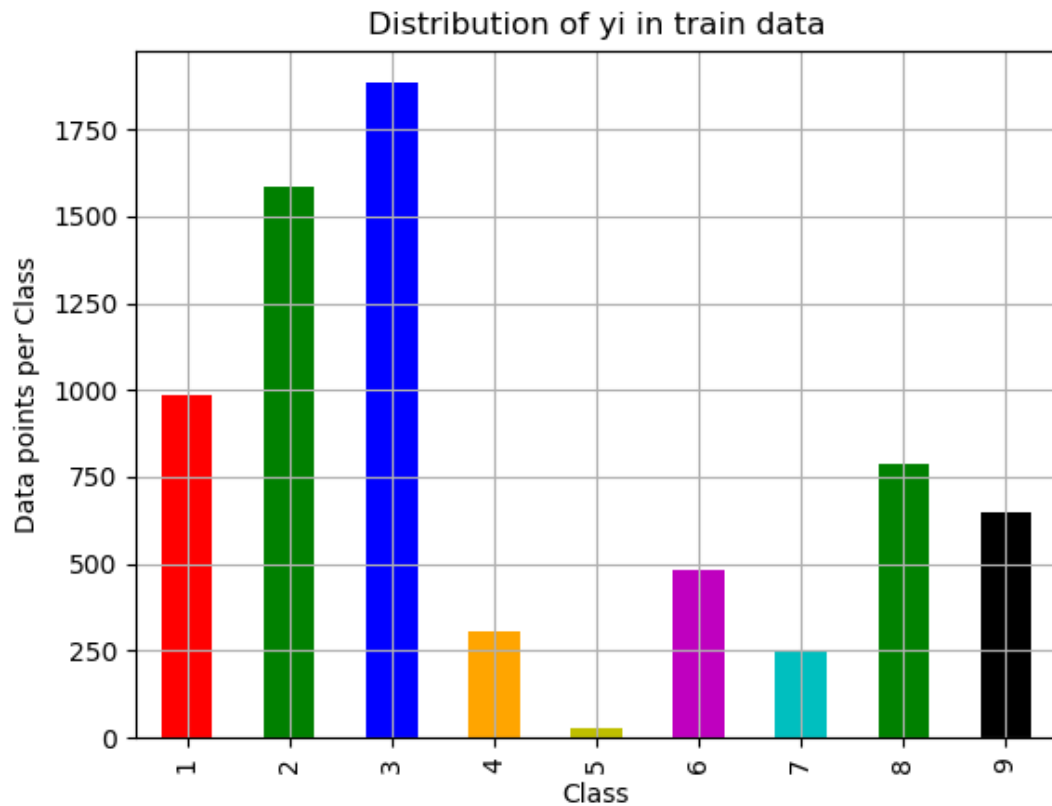
8.2 Check the distribution of Labels in Train, Test and Cross Validation Dataset


```

In [7]: 1 #It returns a dict, keys as class labels and values as the number of data poi
2 train_class_distribution = y_train.value_counts().sortlevel()
3 test_class_distribution = y_test.value_counts().sortlevel()
4 cv_class_distribution = y_cv.value_counts().sortlevel()
5
6 my_colors = ["r","g","b","orange","y","m","c","g","black"]
7 train_class_distribution.plot(kind='bar', color=my_colors)
8 plt.xlabel('Class')
9 plt.ylabel('Data points per Class')
10 plt.title('Distribution of yi in train data')
11 plt.grid()
12 plt.show()
13
14 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.arg
15 # -(train_class_distribution.values): the minus sign will give us in decreasi
16 sorted_yi = np.argsort(-train_class_distribution.values)
17 for i in sorted_yi:
18     print('Number of data points in class', i+1, ':',train_class_distribution
19
20
21 print('-'*80)
22 test_class_distribution.plot(kind='bar', color=my_colors)
23 plt.xlabel('Class')
24 plt.ylabel('Data points per Class')
25 plt.title('Distribution of yi in test data')
26 plt.grid()
27 plt.show()
28
29 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.arg
30 # -(train_class_distribution.values): the minus sign will give us in decreasi
31 sorted_yi = np.argsort(-test_class_distribution.values)
32 for i in sorted_yi:
33     print('Number of data points in class', i+1, ':',test_class_distribution.
34
35 print('-'*80)
36 cv_class_distribution.plot(kind='bar', color=my_colors)
37 plt.xlabel('Class')
38 plt.ylabel('Data points per Class')
39 plt.title('Distribution of yi in cross validation data')
40 plt.grid()
41 plt.show()
42
43 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.arg
44 # -(train_class_distribution.values): the minus sign will give us in decreasi
45 sorted_yi = np.argsort(-train_class_distribution.values)
46 for i in sorted_yi:
47     print('Number of data points in class', i+1, ':',cv_class_distribution.va

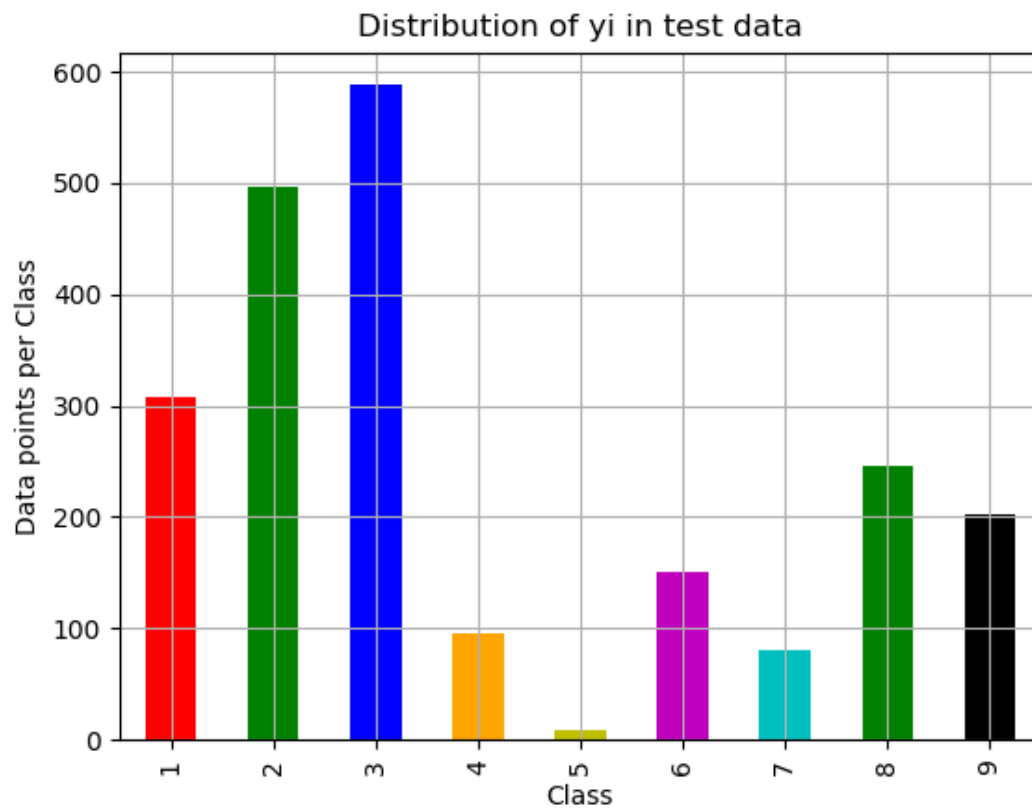
```

<IPython.core.display.Javascript object>



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

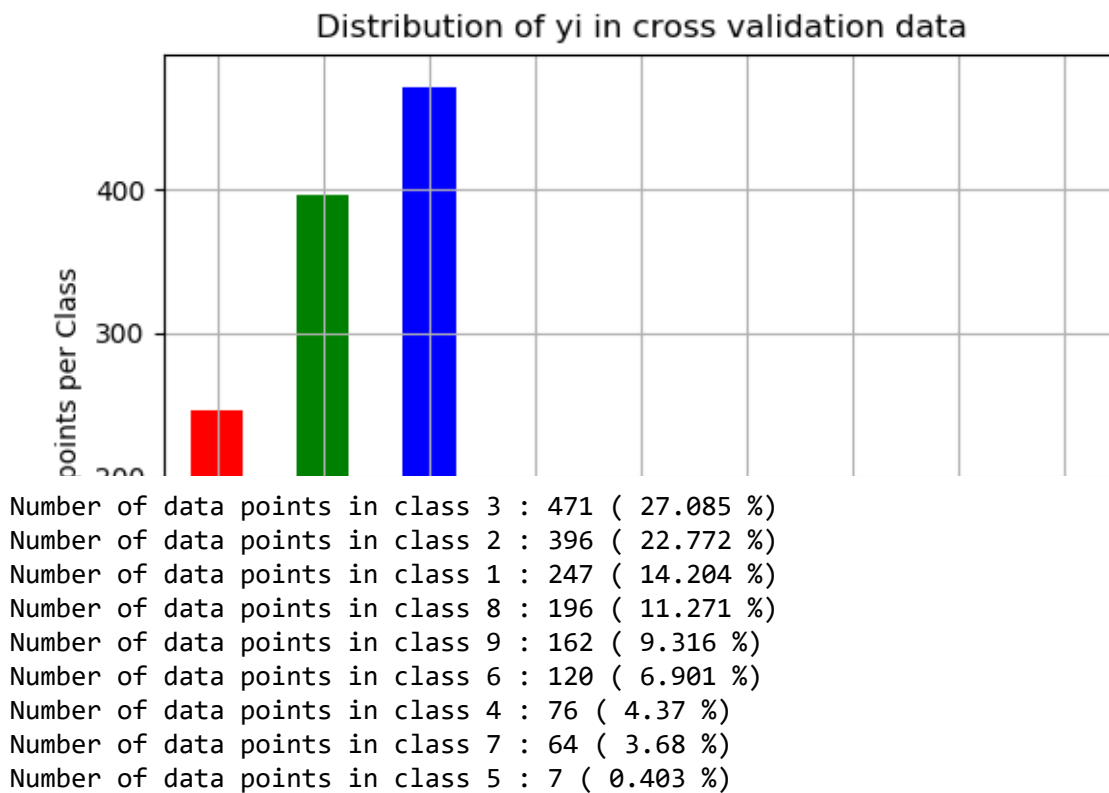
<IPython.core.display.Javascript object>



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

-

<IPython.core.display.Javascript object>



8.3 Code for Confusion, Precision and Recall matrix

```

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

```

Machine learning models

10.1 Training a Random Forest Classifier on the final sets of features

```
In [3]: 1 # import warnings filter
        2 from warnings import simplefilter
        3 # ignore all future warnings
        4 simplefilter(action='ignore', category=FutureWarning)
```

```
In [4]: 1 from sklearn.model_selection import RandomizedSearchCV
        2 from sklearn.ensemble import RandomForestClassifier
        3 from sklearn.calibration import CalibratedClassifierCV
        4 from xgboost import XGBClassifier
        5 import pandas as pd
```

Tune the hyperparameter `n_estimators`.

```

In [6]: 1 estimator=[10,50,100,500,1000,2000,3000]
2 cv_log_error_array=[]
3 from sklearn.ensemble import RandomForestClassifier
4 for i in estimator:
5     r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
6     r_cfl.fit(X_train , y_train )
7     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
8     sig_clf.fit(X_train , y_train )
9     predict_y = sig_clf.predict_proba(X_cv )
10    cv_log_error_array.append(log_loss(y_cv , predict_y, labels=r_cfl.classes_)
11
12    for i in range(len(cv_log_error_array)):
13        print ('log_loss for c = ',estimator[i],'is',cv_log_error_array[i])
14
15
16    best_estimator = np.argmin(cv_log_error_array)
17
18    r_cfl=RandomForestClassifier(n_estimators=estimator[best_estimator],random_state=42,n_jobs=-1)
19    r_cfl.fit(X_train , y_train )
20    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
21    sig_clf.fit(X_train , y_train )
22    predict_y = sig_clf.predict_proba(X_train )
23    print ('For values of best estimator = ', estimator[best_estimator], "The train log loss is: ", log_loss(y_train , predict_y, labels=r_cfl.classes_))
24    predict_y = sig_clf.predict_proba(X_cv )
25    print('For values of best estimator = ', estimator[best_estimator], "The cross validation log loss is: ", log_loss(y_cv , predict_y, labels=r_cfl.classes_))
26    predict_y = sig_clf.predict_proba(X_test )
27    print('For values of best estimator = ', estimator[best_estimator], "The test log loss is: ", log_loss(y_test , predict_y, labels=r_cfl.classes_))
28

```

```

log_loss for c = 10 is 0.05860248748022682
log_loss for c = 50 is 0.046068787750972805
log_loss for c = 100 is 0.04794121453575634
log_loss for c = 500 is 0.04627045596291302
log_loss for c = 1000 is 0.046136037788736546
log_loss for c = 2000 is 0.046206151064847226
log_loss for c = 3000 is 0.04622147586641312

```

```

For values of best alpha = 50 The train log loss is: 0.01773713171122862
For values of best alpha = 50 The cross validation log loss is: 0.046068787750972805
For values of best alpha = 50 The test log loss is: 0.051558208326363784

```

10.2 Training an XGBoost Classifier on the final sets of features

```

In [12]: 1 estimators=[10,50,100,500,1000,2000,3000]
2 cv_log_error_array=[]
3 for i in estimators:
4     x_cfl=XGBClassifier(n_estimators=i, n_jobs=-1)
5     x_cfl.fit(X_train , y_train )
6     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
7     sig_clf.fit(X_train , y_train )
8     predict_y = sig_clf.predict_proba(X_cv )
9     cv_log_error_array.append(log_loss(y_cv , predict_y, labels=x_cfl.classes_))
10
11 for i in range(len(cv_log_error_array)):
12     print ('log_loss for c = ',estimators[i],'is',cv_log_error_array[i])
13
14
15 best_estimators = np.argmin(cv_log_error_array)
16
17 fig, ax = plt.subplots()
18 ax.plot(estimators, cv_log_error_array,c='g')
19 for i, txt in enumerate(np.round(cv_log_error_array,3)):
20     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_
21 plt.grid()
22 plt.title("Cross Validation Error for each estimators")
23 plt.xlabel("Estimators i's")
24 plt.ylabel("Error measure")
25 plt.show()
26
27 x_cfl=XGBClassifier(n_estimators=estimators[best_estimators],nthread=-1,n_job
28 x_cfl.fit(X_train , y_train ,verbose=True)
29 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
30 sig_clf.fit(X_train , y_train )
31
32 predict_y = sig_clf.predict_proba(X_train )
33 print ('For values of best estimators = ', estimators[best_estimators], "The
34 predict_y = sig_clf.predict_proba(X_cv )
35 print('For values of best estimators = ', estimators[best_estimators], "The c
36 predict_y = sig_clf.predict_proba(X_test )
37 print('For values of best estimators = ', estimators[best_estimators], "The t

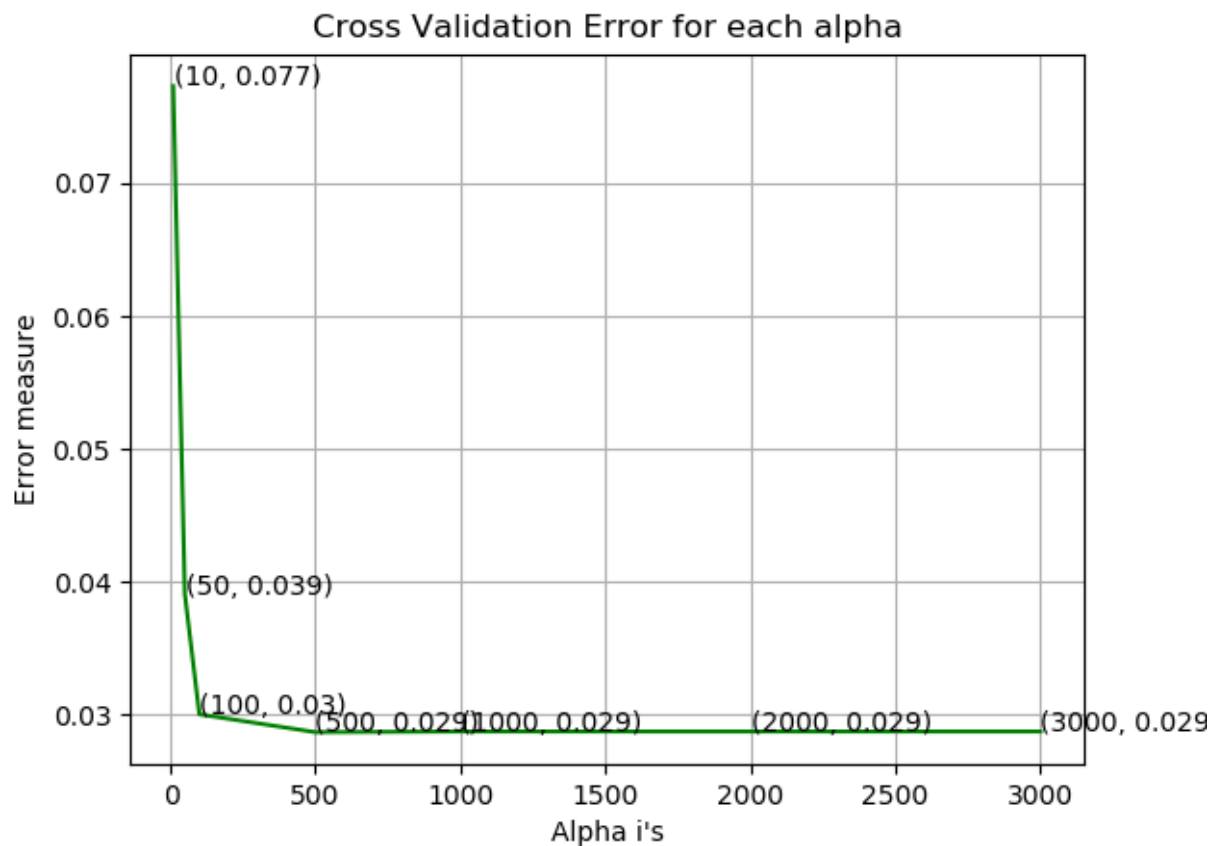
```

```

log_loss for c = 10 is 0.07733012736364236
log_loss for c = 50 is 0.03905875498959229
log_loss for c = 100 is 0.02999076721501786
log_loss for c = 500 is 0.02864663745787661
log_loss for c = 1000 is 0.028697083167181243
log_loss for c = 2000 is 0.02869700239044487
log_loss for c = 3000 is 0.02869688851297395

```

<IPython.core.display.Javascript object>



For values of best alpha = 500 The train log loss is: 0.013021228130231977
 For values of best alpha = 500 The cross validation log loss is: 0.02864663745787661
 For values of best alpha = 500 The test log loss is: 0.030512075737233454

10.3 Selecting top features using Random Forest + Apply XGBoost on top of the 1500 selected fetaures

In [9]: `1 from xgboost import XGBClassifier`

10.3.1 Code for Selecting best features with random forest.

```

In [10]: 1 from sklearn.ensemble import RandomForestClassifier
2 import numpy as np
3
4 def get_feature_importance(X, y, nb_imp_feats): #For selecting top 1500 features
5     n_estimators=[10,50,80,100,125,250,500,1000,2000,3000]
6     cv_log_error_array=[]
7     for i in tqdm(n_estimators):
8         rf_clf=RandomForestClassifier(n_estimators=i,n_jobs=-1, random_state=
9         rf_clf.fit(X,y)
10        sig_clf = CalibratedClassifierCV(rf_clf, method="sigmoid")
11        sig_clf.fit(X_train , y_train )
12        predict_y = sig_clf.predict_proba(X_cv )
13        cv_log_error_array.append(log_loss(y_cv , predict_y, labels=rf_clf.cl
14
15
16        best_estimator = np.argmin(cv_log_error_array)
17        rf_clf=RandomForestClassifier(n_estimators=n_estimators[best_estimator],r
18        rf_clf.fit(X_train ,y_train)
19
20        feature_index=np.argsort(rf_clf.feature_importances_)[::-1]
21        most_imp_feat_idx=feature_index[:nb_imp_feats]
22        return most_imp_feat_idx

```

10.3.2. Using the best features to construct train, test and cross-validation dataset

```

In [11]: 1 #Get all the feature names in a List
2 features = list(X_train.columns)
3
4 #Get the top 1500 features indexes
5 top_1500_features_index = get_feature_importance(X_train, y_train, 1500)
6
7 #Get the top 1500 features
8 top_1500_features = np.take(features,top_1500_features_index)
9
10 #Create a train, test and cv dataset with top 1500 features obtained using ra
11 X_train=X_train[list(top_1500_features)]
12 X_test=X_test[list(top_1500_features)]
13 X_cv=X_cv[list(top_1500_features)]

```

100%|██████████| 10/10 [05:17<00:00, 66.72s/it]

10.3.3. Training an XGBoost classifier on the final set of 1500 features. (Tuning the number of estimators)

```

In [13]: 1 estimators=[10,50,100,250,500,1000,2000,3000]
2 cv_log_error_array=[]
3 for i in tqdm(estimators):
4     x_cfl=XGBClassifier(n_estimators=i, n_jobs=-1)
5     x_cfl.fit(X_train , y_train )
6     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
7     sig_clf.fit(X_train , y_train )
8     predict_y = sig_clf.predict_proba(X_cv )
9     cv_log_error_array.append(log_loss(y_cv , predict_y, labels=x_cfl.classes)
10
11 for i in range(len(cv_log_error_array)):
12     print ('log_loss for c = ',estimators[i],'is',cv_log_error_array[i])
13
14
15 best_estimators = np.argmin(cv_log_error_array)
16
17 fig, ax = plt.subplots()
18 ax.plot(estimators, cv_log_error_array,c='g')
19 for i, txt in enumerate(np.round(cv_log_error_array,3)):
20     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_
21 plt.grid()
22 plt.title("Cross Validation Error for each estimators")
23 plt.xlabel("Estimators i's")
24 plt.ylabel("Error measure")
25 plt.show()
26
27 x_cfl=XGBClassifier(n_estimators=estimators[best_estimators],nthread=-1,n_job
28 x_cfl.fit(X_train , y_train ,verbose=True)
29 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
30 sig_clf.fit(X_train , y_train )
31
32 predict_y = sig_clf.predict_proba(X_train )
33 print ('For values of best estimators = ', estimators[best_estimators], "The
34 predict_y = sig_clf.predict_proba(X_cv )
35 print('For values of best estimators = ', estimators[best_estimators], "The c
36 predict_y = sig_clf.predict_proba(X_test )
37 print('For values of best estimators = ', estimators[best_estimators], "The t

```

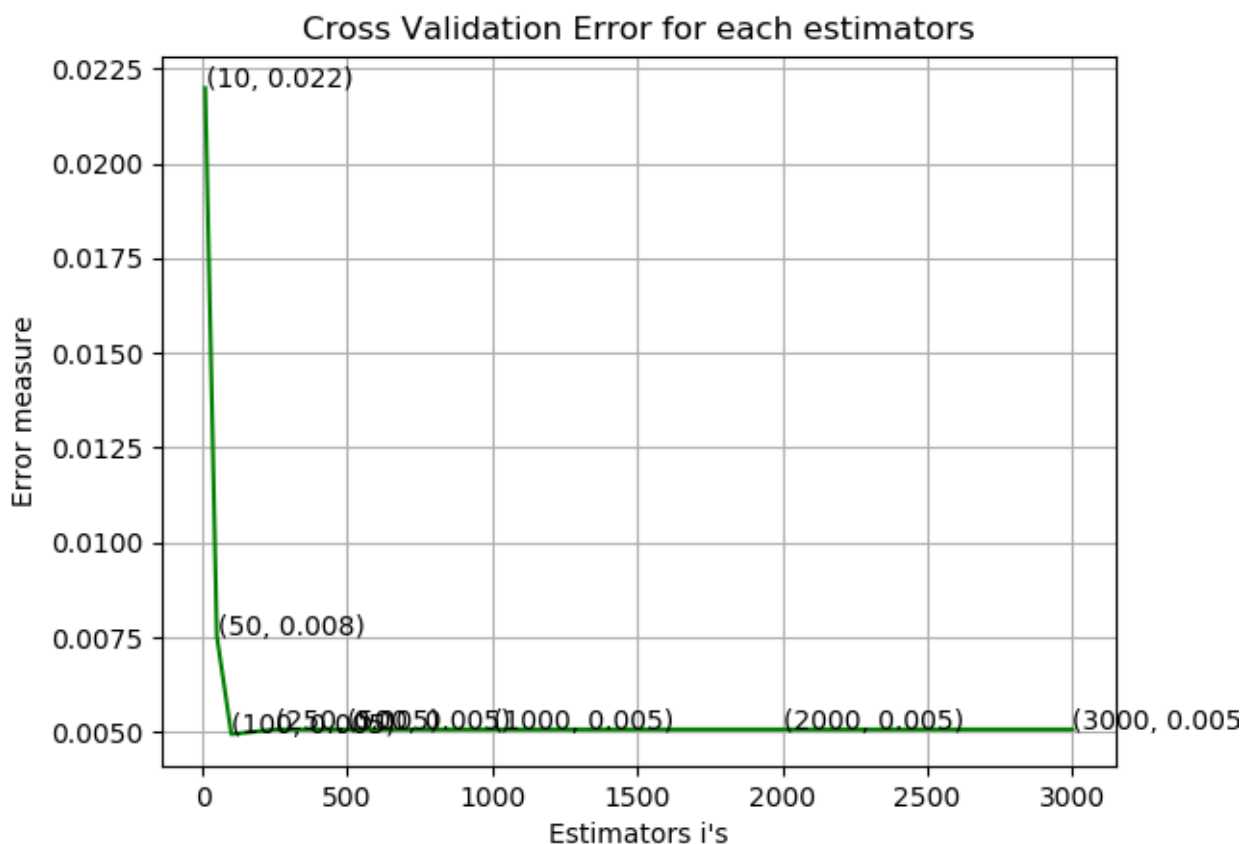
100%|██████████| 8/8 [1:33:17<00:00, 1238.39s/it]

```

log_loss for c = 10 is 0.021982833032912527
log_loss for c = 50 is 0.0075346142355446205
log_loss for c = 100 is 0.004945830182961919
log_loss for c = 250 is 0.005071604155408008
log_loss for c = 500 is 0.00507270568428335
log_loss for c = 1000 is 0.005072926974678231
log_loss for c = 2000 is 0.005072718941999207
log_loss for c = 3000 is 0.005072960627925119

```

<IPython.core.display.Javascript object>



For values of best estimators = 100 The train log loss is: 0.004581871383625352

For values of best estimators = 100 The cross validation log loss is: 0.004945830182961919

For values of best estimators = 100 The test log loss is: 0.004481605018364259

11. Experimenting with best 500 selected features

11.1 Using the best features to construct train, test and cross-validation dataset

```
In [11]: 1 #Get all the feature names in a list
2 features = list(X_train.columns)
3
4 #Get the top 500 features indexes
5 top_500_features_index = get_feature_importance(X_train, y_train, 500)
6
7 #Get the top 500 features
8 top_500_features = np.take(features, top_500_features_index)
9
10 #Create a train, test and cv dataset with top 500 features obtained using ran
11 X_train=X_train[list(top_500_features)]
12 X_test=X_test[list(top_500_features)]
13 X_cv=X_cv[list(top_500_features)]
14
15 #Save the list of features
16 with open('features/top_500_features.pkl', 'wb') as file:
17     pickle.dump(top_500_features, file)
```

100%|██████████| 10/10 [05:13<00:00, 66.20s/it]

11.2.1 XGBoost: Hyper-parameter tuning + Predicting on Test Data (Tuning the number of estimators)

In [13]:

```

1  from tqdm import tqdm
2  from xgboost import XGBClassifier
3
4  estimators=[70,80,90,105,110,115,130,150,180,200,250,300]
5  cv_log_error_array=[]
6  for i in tqdm(estimators):
7      x_cfl=XGBClassifier(n_estimators=i, n_jobs=-1)
8      x_cfl.fit(X_train , y_train )
9      sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
10     sig_clf.fit(X_train , y_train )
11     predict_y = sig_clf.predict_proba(X_cv )
12     cv_log_error_array.append(log_loss(y_cv , predict_y, labels=x_cfl.classes_))
13
14  for i in range(len(cv_log_error_array)):
15     print ('log_loss for n_estimators = ',estimators[i], 'is', cv_log_error_array[i])
16
17
18  best_estimators = np.argmin(cv_log_error_array)
19
20  fig, ax = plt.subplots()
21  ax.plot(estimators, cv_log_error_array, c='g')
22  for i, txt in enumerate(np.round(cv_log_error_array,3)):
23     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_array[i]))
24  plt.grid()
25  plt.title("Cross Validation Error for each estimators")
26  plt.xlabel("Estimators i's")
27  plt.ylabel("Error measure")
28  plt.show()
29
30  x_cfl=XGBClassifier(n_estimators=estimators[best_estimators], nthread=-1, n_jobs=-1)
31  x_cfl.fit(X_train , y_train , verbose=True)
32  sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
33  sig_clf.fit(X_train , y_train )
34
35  predict_y = sig_clf.predict_proba(X_train )
36  print ('For values of best estimators = ', estimators[best_estimators], "The cross validation error is", cv_log_error_array[best_estimators])
37  predict_y = sig_clf.predict_proba(X_cv )
38  print('For values of best estimators = ', estimators[best_estimators], "The cross validation error is", cv_log_error_array[best_estimators])
39  predict_y = sig_clf.predict_proba(X_test )
40  print('For values of best estimators = ', estimators[best_estimators], "The test error is", log_loss(y_test, predict_y))
41
42  plot_confusion_matrix(y_test, x_cfl.predict(X_test))

```

100%|██████████| 12/12 [15:08<00:00, 97.20s/it]

```
log_loss for n_estimators = 70 is 0.005502279196150079
log_loss for n_estimators = 80 is 0.005085534323274856
log_loss for n_estimators = 90 is 0.004872238584347242
log_loss for n_estimators = 105 is 0.004776526887835807
log_loss for n_estimators = 110 is 0.004787429118714052
log loss for n estimators = 115 is 0.004806338024363514
```

11.5 RandomForestClassifier: Hyper-parameter tuning + Predicting on Test Data

In [19]:

```

1  from sklearn.ensemble import RandomForestClassifier
2
3  estimators=[10,50,100,200,350,500,1000,2000]
4  cv_log_error_array=[]
5  for i in tqdm(estimators):
6      rf_clf=RandomForestClassifier(n_estimators=i, n_jobs=-1, class_weight='balanced')
7      rf_clf.fit(X_train ,y_train )
8      sig_clf = CalibratedClassifierCV(rf_clf, method="sigmoid")
9      sig_clf.fit(X_train , y_train )
10     predict_y = sig_clf.predict_proba(X_cv )
11     cv_log_error_array.append(log_loss(y_cv , predict_y, labels=rf_clf.classes_))
12
13  for i in range(len(cv_log_error_array)):
14      print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_array[i])
15
16
17  best_estimators = np.argmin(cv_log_error_array)
18
19  fig, ax = plt.subplots()
20  ax.plot(estimators, cv_log_error_array,c='g')
21  for i, txt in enumerate(np.round(cv_log_error_array,3)):
22      ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_array[i]))
23  plt.grid()
24  plt.title("Cross Validation Error for each estimators")
25  plt.xlabel("Estimators i's")
26  plt.ylabel("Error measure")
27  plt.show()
28
29  rf_clf=RandomForestClassifier(n_estimators=estimators[best_estimators],n_jobs=-1)
30  rf_clf.fit(X_train ,y_train)
31  sig_clf = CalibratedClassifierCV(rf_clf, method="sigmoid")
32  sig_clf.fit(X_train , y_train )
33
34  predict_y = sig_clf.predict_proba(X_train )
35  print ('For values of best estimators = ', estimators[best_estimators], "The cross validation error is", cv_log_error_array[best_estimators])
36  predict_y = sig_clf.predict_proba(X_cv )
37  print('For values of best estimators = ', estimators[best_estimators], "The cross validation error is", cv_log_error_array[best_estimators])
38  predict_y = sig_clf.predict_proba(X_test )
39  print('For values of best estimators = ', estimators[best_estimators], "The test error is", log_loss(y_test, predict_y, labels=rf_clf.classes_))
40
41  #Plot confusion matrices
42  plot_confusion_matrix(y_test, rf_clf.predict(X_test))

```

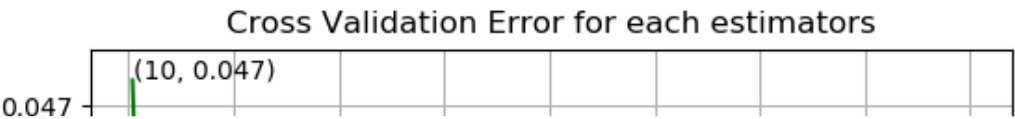
100%|██████████| 8/8 [02:26<00:00, 31.69s/it]

```

log_loss for n_estimators = 10 is 0.04731188287459294
log_loss for n_estimators = 50 is 0.04132885135890258
log_loss for n_estimators = 100 is 0.04423992031938117
log_loss for n_estimators = 200 is 0.03984525650634593
log_loss for n_estimators = 350 is 0.04156142917464634
log_loss for n_estimators = 500 is 0.04162118116948644
log_loss for n_estimators = 1000 is 0.041954036786431295
log_loss for n_estimators = 2000 is 0.04113752448585498

```

<IPython.core.display.Javascript object>



11.6 Stacking Classifiers : XGBClassifier meta classifier

In [29]:

```

1  from mlxtend.classifier import StackingClassifier
2
3  xgb_1 = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
4                        colsample_bytree=1, gamma=0, learning_rate=0.1, max_depth=3, min_child_weight=1, missing=None, n_estimators=100, n_jobs=-1, nthread=-1, objective='multi:softprob', random_state=None, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None, silent=True, subsample=1)
5
6  xgb_2 = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
7                        colsample_bytree=1, gamma=0, learning_rate=0.1, max_depth=3, min_child_weight=1, missing=None, n_estimators=100, n_jobs=-1, nthread=-1, objective='multi:softprob', random_state=None, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None, silent=True, subsample=1)
8
9
10 xgb_2 = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
11                       colsample_bytree=1, gamma=0, learning_rate=0.1, max_depth=3, min_child_weight=1, missing=None, n_estimators=100, n_jobs=-1, nthread=-1, objective='multi:softprob', random_state=None, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None, silent=True, subsample=1)
12
13
14 lgb_1=LGBMClassifier(boosting_type='gbdt', class_weight=None, colsample_bytree=1, importance_type='split', learning_rate=0.1, max_depth=-1, min_child_samples=20, min_child_weight=0.001, min_split_gain=0.01, n_estimators=190, n_jobs=-1, nthread=-1, num_leaves=31, objective=None, random_state=None, reg_alpha=0.0, reg_lambda=1.0, silent=True, subsample=1.0, subsample_for_bin=200000, subsample_freq=0)
15
16
17 lgb_2=LGBMClassifier(boosting_type='gbdt', class_weight=None, colsample_bytree=1, importance_type='split', learning_rate=0.1, max_depth=-1, min_child_samples=20, min_child_weight=0.001, min_split_gain=0.01, n_estimators=190, n_jobs=-1, nthread=-1, num_leaves=31, objective=None, random_state=None, reg_alpha=0.0, reg_lambda=1.0, silent=True, subsample=1.0, subsample_for_bin=200000, subsample_freq=0)
18
19
20 cv_log_error_array=[]
21 estimators = [10,50,100,150,250,500,750,1000,2000,3000,4000]
22 for i in estimators:
23     meta_clf = XGBClassifier(n_estimators=i, n_jobs=-1)
24     sig_clf = StackingClassifier(classifiers=[xgb_1,xgb_2,lgb_1,lgb_2], meta_classifier=meta_clf)
25     sig_clf.fit(X_train,y_train)
26     predict_y = sig_clf.predict_proba(X_cv)
27     cv_log_error_array.append(log_loss(y_cv, predict_y, eps=1e-15))
28
29
30 for i in range(len(cv_log_error_array)):
31     print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_array[i])
32
33
34 best_estimators = np.argmin(cv_log_error_array)
35
36 fig, ax = plt.subplots()
37 ax.plot(estimators, cv_log_error_array,c='g')
38 for i, txt in enumerate(np.round(cv_log_error_array,3)):
39     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_array[i]))
40 plt.grid()
41 plt.title("Cross Validation Error for each estimators")
42 plt.xlabel("Estimators i's")

```

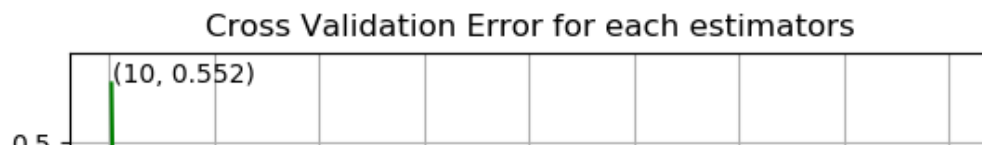
```

57 plt.ylabel("Error measure")
58 plt.show()
59
60 meta_clf = XGBClassifier(n_estimators=estimators[best_estimators], n_jobs=-1)
61 sig_clf=StackingClassifier(classifiers=[xgb_1,xgb_2,lgb_1,lgb_2], meta_classifier=meta_clf)
62 sig_clf.fit(X_train,y_train)
63
64 predict_y = sig_clf.predict_proba(X_train )
65 print ('For values of best_estimators = ', estimators[best_estimators], "The cross validation error is ", log_loss)
66 predict_y = sig_clf.predict_proba(X_cv )
67 print('For values of best_estimators = ', estimators[best_estimators], "The cross validation error is ", log_loss)
68 predict_y = sig_clf.predict_proba(X_test )
69 print('For values of best_estimators = ', estimators[best_estimators], "The test error is ", test_error)
70
71 plot_confusion_matrix(y_test,sig_clf.predict(X_test))

```

log_loss for n_estimators = 10 is 0.5518787449430363
log_loss for n_estimators = 50 is 0.010307589801839195
log_loss for n_estimators = 100 is 0.0008983027434716662
log_loss for n_estimators = 150 is 0.0007903983087084497
log_loss for n_estimators = 250 is 0.0007623319212748128
log_loss for n_estimators = 500 is 0.0007594183390671436
log_loss for n_estimators = 750 is 0.000759415734148798
log_loss for n_estimators = 1000 is 0.000759415734148798
log_loss for n_estimators = 2000 is 0.0007594162214998865
log_loss for n_estimators = 3000 is 0.0007594162214998865
log_loss for n_estimators = 4000 is 0.000759418826418232

<IPython.core.display.Javascript object>



12. Experimenting with Image features + Best 500 features.

12.1 Using the best features to construct train, test and cross-validation dataset

```

In [48]: 1 #Load and merge the dataframes
          2 asm_image_df=pd.read_csv("features/asm_image_df.csv")
          3 asm_image_df_norm = normalize(asm_image_df).fillna(0)
          4
          5 combined_df = pd.concat([X[list(top_500_features)],asm_image_df_norm.drop(["I

```

```
In [54]: 1 X = combined_df.drop(["Class"],axis=1)
2 y = combined_df["Class"]
3
4 #Split the data into test and train by maintaining same distribution of output
5 X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, test_size=0.2)
6
7 #Split the train data into train and cross validation by maintaining same distribution
8 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)
9
10 print('Number of data points in train data:', X_train.shape[0])
11 print('Number of data points in test data:', X_test.shape[0])
12 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 [5]: 1 del(best_features_df,asm_image_df,combined_df, X, y)
```

12.2 XGBoost: Hyper-parameter tuning + Predicting on Test Data

```

In [55]: 1 from tqdm import tqdm
2 estimators=[10,50,100,250,500,1000,2000,3000]
3 cv_log_error_array=[]
4 for i in tqdm(estimators):
5     x_cfl=XGBClassifier(n_estimators=i, n_jobs=-1)
6     x_cfl.fit(X_train , y_train )
7     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
8     sig_clf.fit(X_train , y_train )
9     predict_y = sig_clf.predict_proba(X_cv )
10    cv_log_error_array.append(log_loss(y_cv , predict_y, labels=x_cfl.classes)
11
12 for i in range(len(cv_log_error_array)):
13     print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_arr
14
15
16 best_estimators = np.argmin(cv_log_error_array)
17
18 fig, ax = plt.subplots()
19 ax.plot(estimators, cv_log_error_array,c='g')
20 for i, txt in enumerate(np.round(cv_log_error_array,3)):
21     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_
22 plt.grid()
23 plt.title("Cross Validation Error for each estimators")
24 plt.xlabel("Estimators i's")
25 plt.ylabel("Error measure")
26 plt.show()
27
28 x_cfl=XGBClassifier(n_estimators=estimators[best_estimators],nthread=-1,n_job
29 x_cfl.fit(X_train , y_train ,verbose=True)
30 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
31 sig_clf.fit(X_train , y_train )
32
33 predict_y = sig_clf.predict_proba(X_train )
34 print ('For values of best estimators = ', estimators[best_estimators], "The
35 predict_y = sig_clf.predict_proba(X_cv )
36 print('For values of best estimators = ', estimators[best_estimators], "The c
37 predict_y = sig_clf.predict_proba(X_test )
38 print('For values of best estimators = ', estimators[best_estimators], "The t
39
40 #Plot confusion matrices
41 plot_confusion_matrix(y_test, x_cfl.predict(X_test))

```

100%|██████████| 8/8 [1:18:49<00:00, 978.67s/it]

```

log_loss for n_estimators = 10 is 0.1065423079653175
log_loss for n_estimators = 50 is 0.061644552270603836
log_loss for n_estimators = 100 is 0.05031276599984848
log_loss for n_estimators = 250 is 0.048096991796449724
log_loss for n_estimators = 500 is 0.04772363821774151
log_loss for n_estimators = 1000 is 0.04759632347924493
log_loss for n_estimators = 2000 is 0.047776135523249774
log_loss for n_estimators = 3000 is 0.04777639490629287

```

<IPython.core.display.Javascript object>

Cross Validation Error for each estimators

(10, 0.107)					
-------------	--	--	--	--	--

12.4 ExtraTreesClassifier: Hyper-parameter tuning + Predicting on Test Data

In [57]:

```

1  from sklearn.ensemble import ExtraTreesClassifier
2
3  estimators=[200,210,220,230,240,250,260,270,280,290,300]
4  cv_log_error_array=[]
5  for i in tqdm(estimators):
6      extra_clf=ExtraTreesClassifier(n_estimators=i, n_jobs=-1, class_weight='b
7      extra_clf.fit(X_train ,y_train )
8      sig_clf = CalibratedClassifierCV(extra_clf, method="sigmoid")
9      sig_clf.fit(X_train , y_train )
10     predict_y = sig_clf.predict_proba(X_cv )
11     cv_log_error_array.append(log_loss(y_cv , predict_y, labels=extra_clf.cla
12
13  for i in range(len(cv_log_error_array)):
14      print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_arr
15
16
17  best_estimators = np.argmin(cv_log_error_array)
18
19  fig, ax = plt.subplots()
20  ax.plot(estimators, cv_log_error_array,c='g')
21  for i, txt in enumerate(np.round(cv_log_error_array,3)):
22      ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_
23  plt.grid()
24  plt.title("Cross Validation Error for each estimators")
25  plt.xlabel("Estimators i's")
26  plt.ylabel("Error measure")
27  plt.show()
28
29  extra_clf=ExtraTreesClassifier(n_estimators=estimators[best_estimators],n_job
30  extra_clf.fit(X_train ,y_train)
31  sig_clf = CalibratedClassifierCV(extra_clf, method="sigmoid")
32  sig_clf.fit(X_train , y_train )
33
34  predict_y = sig_clf.predict_proba(X_train )
35  print ('For values of best estimators = ', estimators[best_estimators], "The
36  predict_y = sig_clf.predict_proba(X_cv )
37  print('For values of best estimators = ', estimators[best_estimators], "The c
38  predict_y = sig_clf.predict_proba(X_test )
39  print('For values of best estimators = ', estimators[best_estimators], "The t
40
41  #Plot confusion matrices
42  plot_confusion_matrix(y_test, extra_clf.predict(X_test))

```

100%|██████████| 11/11 [01:33<00:00, 9.36s/it]

```

log_loss for n_estimators = 200 is 0.05821910745413007
log_loss for n_estimators = 210 is 0.05664206175410856
log_loss for n_estimators = 220 is 0.056079021386598156
log_loss for n_estimators = 230 is 0.05564309916566966
log_loss for n_estimators = 240 is 0.05747641782647607
log_loss for n_estimators = 250 is 0.05539290671043567
log_loss for n_estimators = 260 is 0.057051025361702375
log_loss for n_estimators = 270 is 0.0567654098411044
log_loss for n_estimators = 280 is 0.056530599130382136
log_loss for n_estimators = 290 is 0.05629446841052641
log_loss for n_estimators = 300 is 0.056685125738838577

```

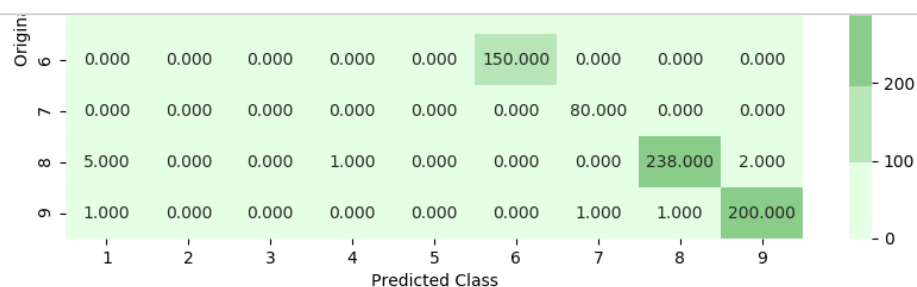
<IPython.core.display.Javascript object>

12.5 RandomForestClassifier: Hyper-parameter tuning + Predicting on Test Data


```

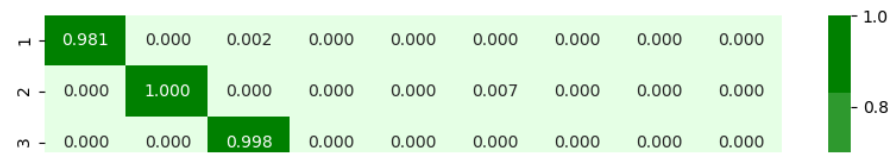
In [58]: 1 from sklearn.ensemble import RandomForestClassifier
2
3 estimators=[150,160,170,180,190,200,210,220,230,240]
4 cv_log_error_array=[]
5 for i in tqdm(estimators):
6     rf_clf=RandomForestClassifier(n_estimators=i, n_jobs=-1, class_weight='balanced')
7     rf_clf.fit(X_train ,y_train )
8     sig_clf = CalibratedClassifierCV(rf_clf, method="sigmoid")
9     sig_clf.fit(X_train , y_train )
10    predict_y = sig_clf.predict_proba(X_cv )
11    cv_log_error_array.append(log_loss(y_cv , predict_y, labels=rf_clf.classes_))
12
13 for i in range(len(cv_log_error_array)):
14     print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_array[i])
15
16
17 best_estimators = np.argmin(cv_log_error_array)
18
19 fig, ax = plt.subplots()
20 ax.plot(estimators, cv_log_error_array,c='g')
21 for i, txt in enumerate(np.round(cv_log_error_array,3)):
22     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_array[i]))
23 plt.grid()
24 plt.title("Cross Validation Error for each estimators")
25 plt.xlabel("Estimators i's")
26 plt.ylabel("Error measure")
27 plt.show()
28
29 rf_clf=RandomForestClassifier(n_estimators=estimators[best_estimators],n_jobs=-1)
30 rf_clf.fit(X_train ,y_train)
31 sig_clf = CalibratedClassifierCV(rf_clf, method="sigmoid")
32 sig_clf.fit(X_train , y_train )
33
34 predict_y = sig_clf.predict_proba(X_train )
35 print ('For values of best estimators = ', estimators[best_estimators], "The cross validation error is", cv_log_error_array[best_estimators])
36 predict_y = sig_clf.predict_proba(X_cv )
37 print('For values of best estimators = ', estimators[best_estimators], "The cross validation error is", cv_log_error_array[best_estimators])
38 predict_y = sig_clf.predict_proba(X_test )
39 print('For values of best estimators = ', estimators[best_estimators], "The test error is", log_loss(y_test, predict_y, labels=rf_clf.classes_))
40
41 #Plot confusion matrices
42 plot_confusion_matrix(y_test, rf_clf.predict(X_test))

```



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

<IPython.core.display.Javascript object>



12.6 Stacking Classifiers : XGBClassifier meta classifier

In [59]:

```

1  from mlxtend.classifier import StackingClassifier
2
3  xgb_1 = x_cfl
4  xgb_2 = x_cfl
5
6
7  lgb_1 = lgbm_clf
8  lgb_2 = lgbm_clf
9
10 cv_log_error_array=[]
11 estimators = [10,50,100,150,250,500,750,1000,2000,3000,4000]
12 for i in estimators:
13     meta_clf = XGBClassifier(n_estimators=i, n_jobs=-1)
14     sig_clf = StackingClassifier(classifiers=[xgb_1,xgb_2,lgb_1,lgb_2], meta_
15     sig_clf.fit(X_train,y_train)
16     predict_y = sig_clf.predict_proba(X_cv)
17     cv_log_error_array.append(log_loss(y_cv, predict_y, eps=1e-15))
18
19
20 for i in range(len(cv_log_error_array)):
21     print ('log_loss for n_estimators = ',estimators[i],'is',cv_log_error_arr
22
23
24 best_estimators = np.argmin(cv_log_error_array)
25
26 fig, ax = plt.subplots()
27 ax.plot(estimators, cv_log_error_array,c='g')
28 for i, txt in enumerate(np.round(cv_log_error_array,3)):
29     ax.annotate((estimators[i],np.round(txt,3)), (estimators[i],cv_log_error_
30 plt.grid()
31 plt.title("Cross Validation Error for each estimators")
32 plt.xlabel("Estimators i's")
33 plt.ylabel("Error measure")
34 plt.show()
35
36 meta_clf = XGBClassifier(n_estimators=estimators[best_estimators], n_jobs=-1)
37 sig_clf=StackingClassifier(classifiers=[xgb_1,xgb_2,lgb_1,lgb_2], meta_classi
38 sig_clf.fit(X_train,y_train)
39
40 predict_y = sig_clf.predict_proba(X_train )
41 print ('For values of best estimators = ', estimators[best_estimators], "The
42 predict_y = sig_clf.predict_proba(X_cv )
43 print('For values of best estimators = ', estimators[best_estimators], "The c
44 predict_y = sig_clf.predict_proba(X_test )
45 print('For values of best estimators = ', estimators[best_estimators], "The t
46
47 plot_confusion_matrix(y_test,sig_clf.predict(X_test))

```

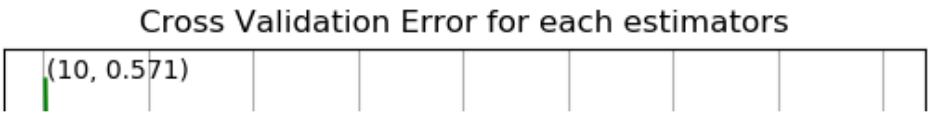
```

log_loss for n_estimators = 10 is 0.5710341889601867
log_loss for n_estimators = 50 is 0.06353912034991523
log_loss for n_estimators = 100 is 0.07225381782059646
log_loss for n_estimators = 150 is 0.07304243437228607
log_loss for n_estimators = 250 is 0.07339369495270898
log_loss for n_estimators = 500 is 0.07343474542690308
log_loss for n_estimators = 750 is 0.07343476643762605
log_loss for n_estimators = 1000 is 0.07343476753443377

```

```
log_loss for n_estimators = 2000 is 0.073434769048971
log_loss for n_estimators = 3000 is 0.07343476795216328
log_loss for n_estimators = 4000 is 0.07343476997440251

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



Display the results in a table.

```

In [7]: 1 from prettytable import PrettyTable
2
3 #Table 1
4 print("Initial Models using simple BYTE features:")
5 table =PrettyTable()
6 table.field_names = ["Model", "Train Log Loss", "Test Log loss", "%tage Misca
7 table.add_row(["Random Model", 2.4561 , 2.4851, 88.5004])
8 table.add_row(["KNN", 0.0782 , 0.2415, 4.5078])
9 table.add_row(["Logistic Regression", 0.4989 , 0.5283, 12.3275])
10 table.add_row(["Random Forest", 0.0266 , 0.0858, 2.0239])
11 table.add_row(["XGBoost", 0.0225 , 0.0792, 1.2419])
12 table.add_row(["XGBoost (Parameter Tuned)", 0.0225 , 0.0782, "----"])
13 print(table)
14
15 #Table 2
16 print("\n\nInitial Models using simple ASM features:")
17 table =PrettyTable()
18 table.field_names = ["Model", "Train Log Loss", "Test Log loss", "%tage Misca
19 table.add_row(["KNN", 0.0476 , 0.0894, 2.0239])
20 table.add_row(["Logistic Regression", 0.3962 , 0.4156, 9.6136])
21 table.add_row(["Random Forest", 0.0166 , 0.0571, 1.1499])
22 table.add_row(["XGBoost", 0.0117 , 0.0491, 0.8739])
23 table.add_row(["XGBoost (Parameter Tuned)", 0.0102 , 0.0483, "----"])
24 print(table)
25
26 #Table 3
27 print("\n\nMerging both ASM and BYTE features:")
28 table =PrettyTable()
29 table.field_names = ["Model", "Train Log Loss", "Test Log loss", "%tage Misca
30 table.add_row(["Random Forest", 0.0166 , 0.0401, "----"])
31 table.add_row(["XGBoost", 0.0111 , 0.0323, "----"])
32 table.add_row(["XGBoost (Parameter Tuned)", 0.0121 , 0.0317, "----"])
33 print(table)
34
35 print("\n\nAdvanced features - top 1500 features:")
36 table =PrettyTable()
37 table.field_names = ["Model", "Train Log Loss", "Test Log loss", "%tage Misca
38 table.add_row(["Random Forest on advanced features", 0.0177 , 0.0515, "----"])
39 table.add_row(["XGBoost on advanced features", 0.0131 , 0.0305, 0.5059])
40 table.add_row(["XGBoost on top 1500 advanced features", 0.0045 , 0.0044, 0.00
41
42 print(table)
43
44 print("\n\nTop 500 features + ASM top 800 pixels:")
45 table =PrettyTable()
46 table.field_names = ["Model", "Train Log Loss", "Test Log loss", "%tage Misca
47 table.add_row(["XGBoost", 0.0117 , 0.0212, 0.2759])
48 table.add_row(["RandomForest", 0.0159 , 0.0393, 0.6899])
49 table.add_row(["EnsembleModel1", 0.0103 , 0.0285, 0.2759])
50 print(table)
51
52 print("\n\nAdvanced features - top 500 features:")
53 table =PrettyTable()
54 table.field_names = ["Model", "Train Log Loss", "Test Log loss", "%tage Misca
55 table.add_row(["XGBoost", 0.0045 , 0.0044, 0.0000])
56 table.add_row(["RandomForest", 0.0125 , 0.0284, 0.5519])

```

```

57 table.add_row(["EnsembleModel1", 0.0007 , 0.0007, 0.0000])
58
59 print(table)

```

Initial Models using simple BYTE features:

Model	Train Log Loss	Test Log loss	%tage Miscalssified points
Random Model	2.4561	2.4851	88.5004
KNN	0.0782	0.2415	4.5078
Logistic Regression	0.4989	0.5283	12.3275
Random Forest	0.0266	0.0858	2.0239
XGBoost	0.0225	0.0792	1.2419
XGBoost (Parameter Tuned)	0.0225	0.0782	----

Initial Models using simple ASM features:

Model	Train Log Loss	Test Log loss	%tage Miscalssified points
KNN	0.0476	0.0894	2.0239
Logistic Regression	0.3962	0.4156	9.6136
Random Forest	0.0166	0.0571	1.1499
XGBoost	0.0117	0.0491	0.8739
XGBoost (Parameter Tuned)	0.0102	0.0483	----

Merging both ASM and BYTE features:

Model	Train Log Loss	Test Log loss	%tage Miscalssified points
Random Forest	0.0166	0.0401	----
XGBoost	0.0111	0.0323	----

XGBoost (Parameter Tuned)	0.0121	0.0317	----
-----+-----+-----+-----			
-----+			

Advanced features - top 1500 features:

Model	Train Log Loss	Test Log loss	%tage Miscalssified points
Random Forest on advanced features	0.0177	0.0515	
XGBoost on advanced features	0.0131	0.0305	0.5059
XGBoost on top 1500 advanced features	0.0045	0.0044	0.0

Top 500 features + ASM top 800 pixels:

Model	Train Log Loss	Test Log loss	%tage Miscalssified points
XGBoost	0.0117	0.0212	0.2759
RandomForest	0.0159	0.0393	0.6899
EnsembleModel1	0.0103	0.0285	0.2759

Advanced features - top 500 features:

Model	Train Log Loss	Test Log loss	%tage Miscalssified points
XGBoost	0.0045	0.0044	0.0
RandomForest	0.0125	0.0284	0.5519
EnsembleModel1	0.0007	0.0007	0.0

What is the objective?

The main purpose of this case study was to classify a given .byte or a .asm file in one of the 9 categories of malware. The 9 categories are Ramnit, Lollipop, Kelihos_ver3, Vundo, Simda, Tracur, Kelihos_ver1, Obfuscator, ACY and Gatak.

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 software 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 malware.

Since it's a multiclass classification problem, we have to choose our KPI (Key Performance Indicator) in such a way so that it minimizes the overall error on all the 9 classes. We want to know not only which file belongs to which class but also what is the probability of that belonging to a particular class. Without probability scores that model would just say that a file belongs to class 9(say). So the ideal KPI for this problem would be multi-class log loss.

Conclusions:

1. First i tried with Byte bigrams and got the best score of 0.0782 for XGBoost
2. Then ASM image 500 features and got 0.483 for XGboost
3. Next for ASM image features(500)+byte bigrams(1000)+byte trigrams(1000) i got 0.037 for boosted trees
4. Next, finally for advanced top 5000 features i got a log loss of 0.0007