# Microsoft Malware detection

# 1.Business/Real-world Problem

## 1.1. What is Malware?

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

## 1.2. Problem Statement

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

## 1.3 Source/Useful Links

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

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

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

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

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

# 2. Machine Learning Problem

## 2.1. Data

### 2.1.1. Data Overview

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

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

### 2.1.2. Example Data Point

**.asm file**

```
.text:00401000                                          assume
es:nothing, ss: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::exc
eption(char const * const &)
.text:0040100D C7 06 08    BB 42 00
mov     dword ptr [esi],    offset off_42BB08
.text:00401013 8B C6                                    mov
eax, esi
.text:00401015 5E                                       pop
esi
.text:00401016 C2 04 00                                        r
etn    4
.text:00401016                                  ;
----------------------------------------------------------------
---------------
.text:00401019 CC CC CC    CC CC CC CC
align 10h
.text:00401020 C7 01 08    BB 42 00
mov     dword ptr [ecx],    offset off_42BB08
.text:00401026 E9 26 1C    00 00
jmp     sub_402C51
.text:00401026                                  ;
----------------------------------------------------------------
---------------
.text:0040102B CC CC CC    CC CC
align 10h
.text:00401030 56                                       push
esi
.text:00401031 8B F1                                    mov
esi, ecx
.text:00401033 C7 06 08    BB 42 00
mov     dword ptr [esi],    offset off_42BB08
.text:00401039 E8 13 1C    00 00
call    sub_402C51
.text:0040103E F6 44 24    08 01
test    byte ptr    [esp+8], 1
.text:00401043 74 09                                    jz
short loc_40104E
.text:00401045 56                                       push
esi
.text:00401046 E8 6C 1E    00 00
```

```
call    ??3@YAXPAX@Z   ; operator delete(void *)
.text:0040104B 83 C4 04                                        a
dd     esp, 4
.text:0040104E
.text:0040104E                                loc_40104E:
; CODE XREF: .text:00401043j
.text:0040104E 8B C6                                          mov
eax, esi
.text:00401050 5E                                             pop
esi
.text:00401051 C2 04 00                                        r
etn    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 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
```

## 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 Learing Objectives and Constraints

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

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

## 2.3. Train and Test Dataset

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

## 2.4. Useful blogs, videos and reference papers

http://blog.kaggle.com/2015/05/26/microsoft-malware-winners-interview-1st-place-no-to-overfitting/
https://arxiv.org/pdf/1511.04317.pdf
First place solution in Kaggle competition: https://www.youtube.com /watch?v=VLQTRILGz5Y
https://github.com/dchad/malware-detection
http://vizsec.org/files/2011/Nataraj.pdf
https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EelnEjvvuQg2nu_pIB6ua?dl=0
" Cross validation is more trustworthy than domain knowledge."

# 3. Exploratory Data Analysis

```
In [0]:  import warnings
         warnings.filterwarnings("ignore")
         import shutil
         import os
         import pandas as pd
         import matplotlib
         matplotlib.use(u'nbAgg')
         import matplotlib.pyplot as plt
         import seaborn as sns
         import numpy as np
         import pickle
         from sklearn.manifold import TSNE
         from sklearn import preprocessing
         import pandas as pd
         from multiprocessing import Process# this is used for multithreadin
         import multiprocessing
         import codecs# this is used for file operations
         import random as r
```

```python
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

In [0]:
```python
#separating byte files and asm files

source = 'train'
destination = 'byteFiles'

# we will check if the folder 'byteFiles' exists if it not there we
if not os.path.isdir(destination):
    os.makedirs(destination)

# if we have folder called 'train' (train folder contains both .asm
# for every file that we have in our 'asmFiles' directory we check
# 'byteFiles' folder

# so by the end of this snippet we will separate all the .byte file
if os.path.isdir(source):
    os.rename(source,'asmFiles')
    source='asmFiles'
    data_files = os.listdir(source)
    for file in asm_files:
        if (file.endswith("bytes")):
            shutil.move(source+file,destination)
```

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

In [0]:
```python
Y=pd.read_csv("trainLabels.csv")
total = len(Y)*1.
ax=sns.countplot(x="Class", data=Y)
for p in ax.patches:
        ax.annotate('{:.1f}%'.format(100*p.get_height()/total), (p.
        
#put 11 ticks (therefore 10 steps), from 0 to the total number of r
ax.yaxis.set_ticks(np.linspace(0, total, 11))

#adjust the ticklabel to the desired format, without changing the p
ax.set_yticklabels(map('{:.1f}%'.format, 100*ax.yaxis.get_majortick
plt.show()
```
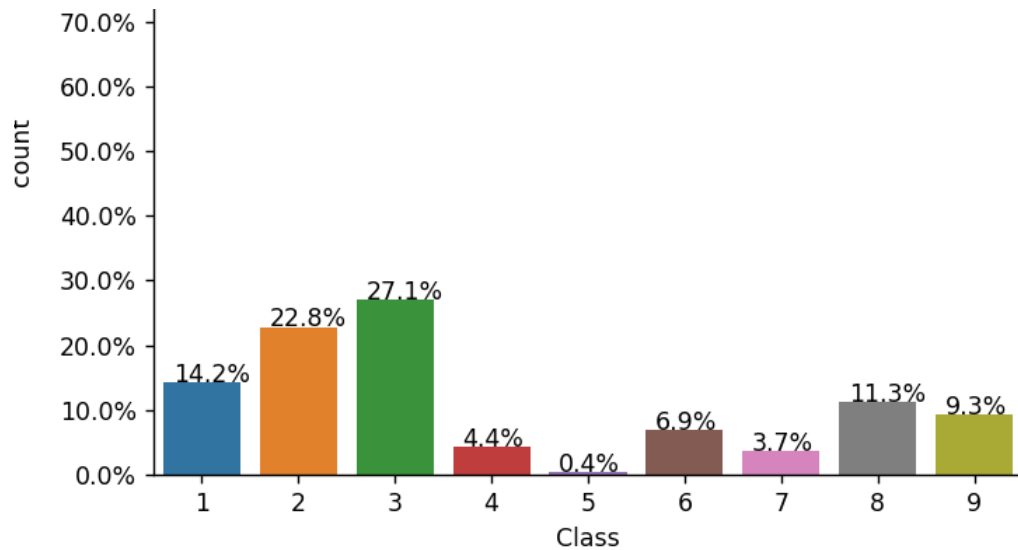
```
<IPython.core.display.Javascript object>
```

```
100.0% ┤

 90.0% ┤

 80.0% ┤
```

## 3.2. Feature extraction

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

```
In [0]: #file sizes of byte files

files=os.listdir('byteFiles')
filenames=Y['Id'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev
    # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st
    # read more about os.stat: here https://www.tutorialspoint.com/
    statinfo=os.stat('byteFiles/'+file)
    # split the file name at '.' and take the first part of it i.e
    file=file.split('.')[0]
    if any(file == filename for filename in filenames):
        i=filenames.index(file)
        class_bytes.append(class_y[i])
        # converting into Mb's
        sizebytes.append(statinfo.st_size/(1024.0*1024.0))
        fnames.append(file)
data_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':c
print (data size byte head())
```
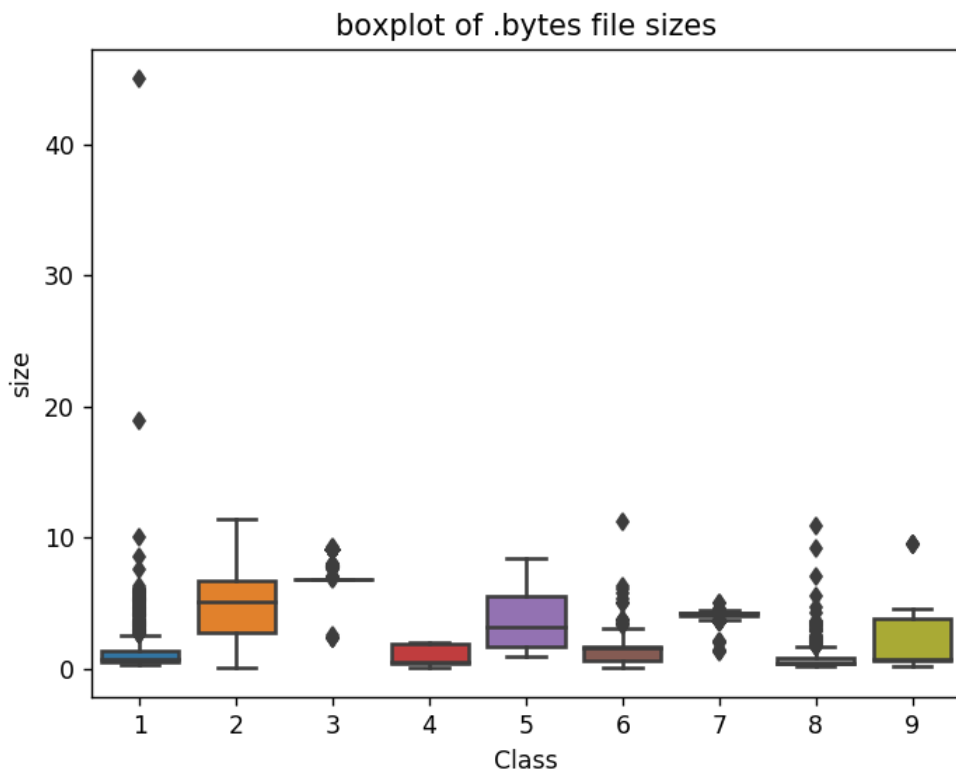
```
    Class              ID      size
0       9  01azqd4InC7m9JpocGv5  4.234863
1       2  01IsoiSMh5gxyDYTl4CB  5.538818
2       9  01jsnpXSAlgw6aPeDxrU  3.887939
3       1  01kcPWA9K2BOxQeS5Rju  0.574219
4       8  01SuzwMJEIXsK7A8dQbl  0.370850
```

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

```
In [0]: #boxplot of byte files
        ax = sns.boxplot(x="Class", y="size", data=data_size_byte)
        plt.title("boxplot of .bytes file sizes")
        plt.show()
```

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



boxplot of .bytes file sizes

### 3.2.3 feature extraction from byte files

```
In [0]: #removal of addres from byte files
        # contents of .byte files
        # ----------------
        #00401000 56 8D 44 24 08 50 8B F1 E8 1C 1B 00 00 C7 06 08
        #-------------------
        #we remove the starting address 00401000

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

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


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

    k += 1

byte_feature_file.close()
```

```python
In [0]: byte_features=pd.read_csv("result.csv")
        print (byte_features.head())
```

```
                                 ID      0      1      2      3      4      5   \
       6      7  \
```

```
In [0]:  result = pd.merge(byte_features, data_size_byte,on='ID', how='left'
         result head()
```

Out[44]:

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

5 rows × 260 columns

```
In [0]:  # https://stackoverflow.com/a/29651514
         def normalize(df):
             result1 = df.copy()
             for feature_name in df.columns:
                 if (str(feature_name) != str('ID') and str(feature_name)!=st
                     max_value = df[feature_name].max()
                     min_value = df[feature_name].min()
                     result1[feature_name] = (df[feature_name] - min_value)
             return result1
         result = normalize(result)
```

```
In [0]:  data_y = result['Class']
         result head()
```
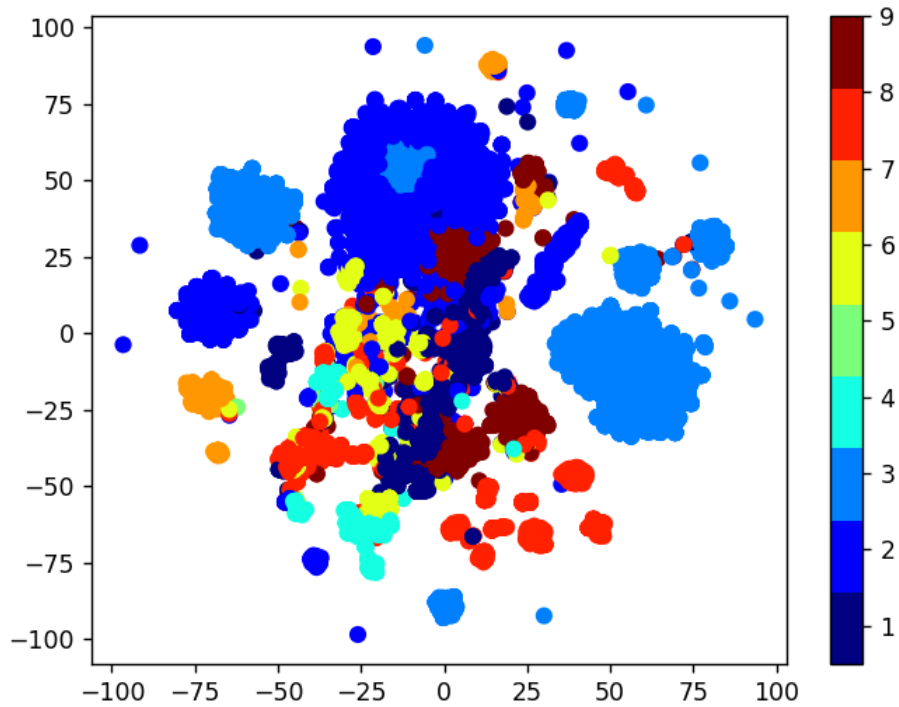
Out[53]:

| | ID | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|
| 0 | 01azqd4InC7m9JpocGv5 | 0.262806 | 0.005498 | 0.001567 | 0.002067 | 0.002048 | 0.001835 | 0.0 |
| 1 | 01IsoiSMh5gxyDYTl4CB | 0.017358 | 0.011737 | 0.004033 | 0.003876 | 0.005303 | 0.003873 | 0.0 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 0.040827 | 0.013434 | 0.001429 | 0.001315 | 0.005464 | 0.005280 | 0.0 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 0.009209 | 0.001708 | 0.000404 | 0.000441 | 0.000770 | 0.000354 | 0.0 |
| 4 | 01SuzwMJEIXsK7A8dQbl | 0.008629 | 0.001000 | 0.000168 | 0.000234 | 0.000342 | 0.000232 | 0.0 |

5 rows × 260 columns

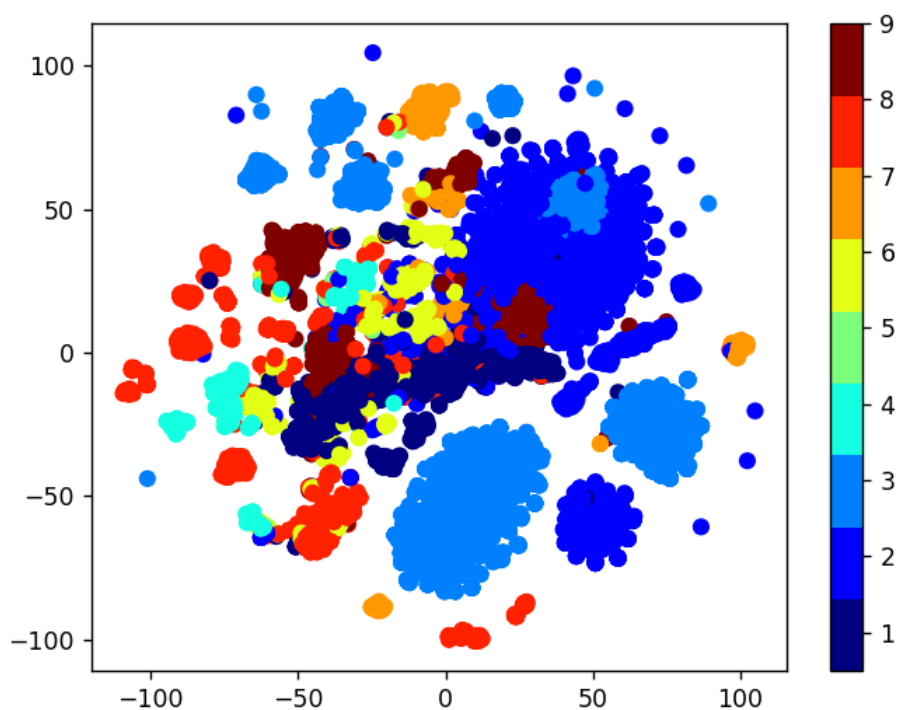### 3.2.4 Multivariate Analysis

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

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

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

<IPython.core.display.Javascript object>

# Train Test split

```
In [0]: data_y = result['Class']
        # split the data into test and train by maintaining same distributi
        X_train, X_test, y_train, y_test = train_test_split(result.drop(['I
        # split the train data into train and cross validation by maintaini
        X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,st
```

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

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

```
In [0]: # it returns a dict, keys as class labels and values as the number
        train_class_distribution = y_train.value_counts().sortlevel()
        test_class_distribution = y_test.value_counts().sortlevel()
        cv_class_distribution = y_cv.value_counts().sortlevel()

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

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


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

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


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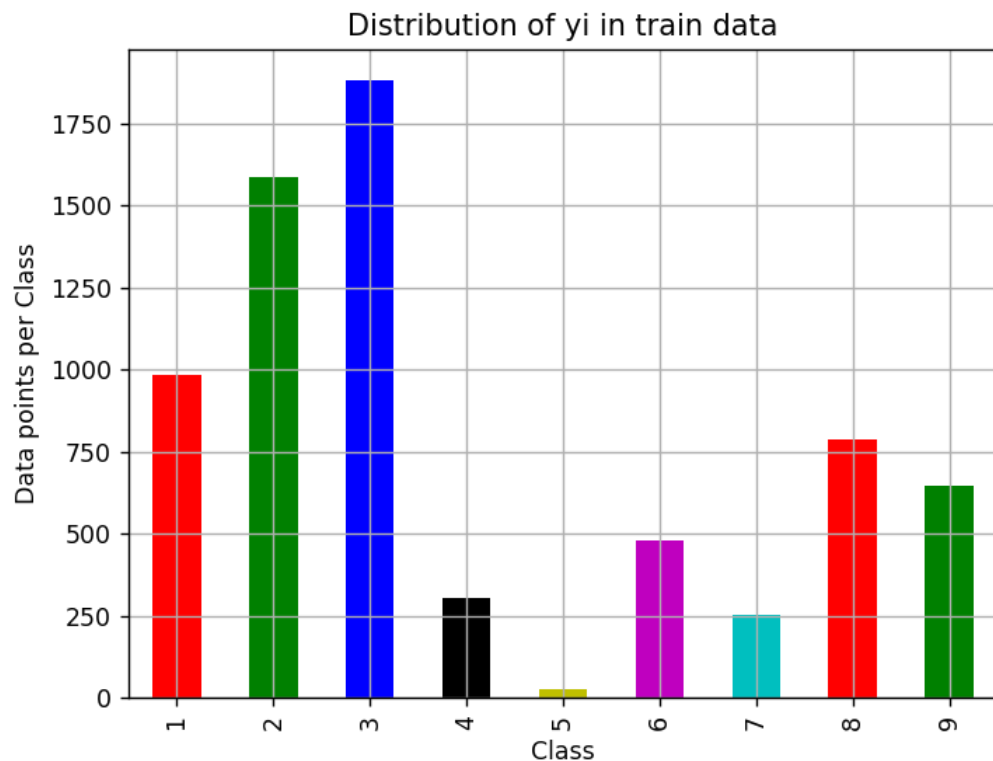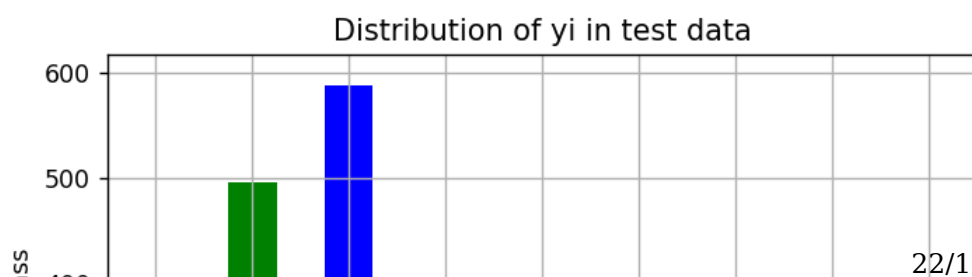
```python
    plt.grid()
    plt.show()

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

<IPython.core.display.Javascript object>

Distribution of yi in train data



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

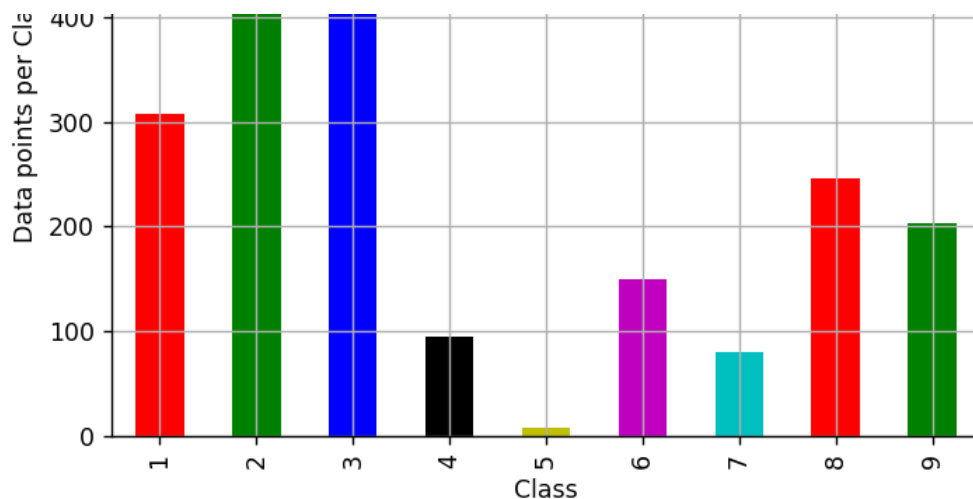<IPython.core.display.Javascript object>

Distribution of yi in test data

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

<IPython.core.display.Javascript object>

### Distribution of yi in cross validation data



```
Number of data points in class 3 : 471 ( 27.085 %)
Number of data points in class 2 : 396 ( 22.772 %)
Number of data points in class 1 : 247 ( 14.204 %)
Number of data points in class 8 : 196 ( 11.271 %)
Number of data points in class 9 : 162 ( 9.316 %)
```

```
            Number of data points in class 6 : 120 ( 6.901 %)
            Number of data points in class 4 : 76 ( 4.37 %)
            Number of data points in class 7 : 64 ( 3.68 %)
            Number of data points in class 5 : 7 ( 0.403 %)
```
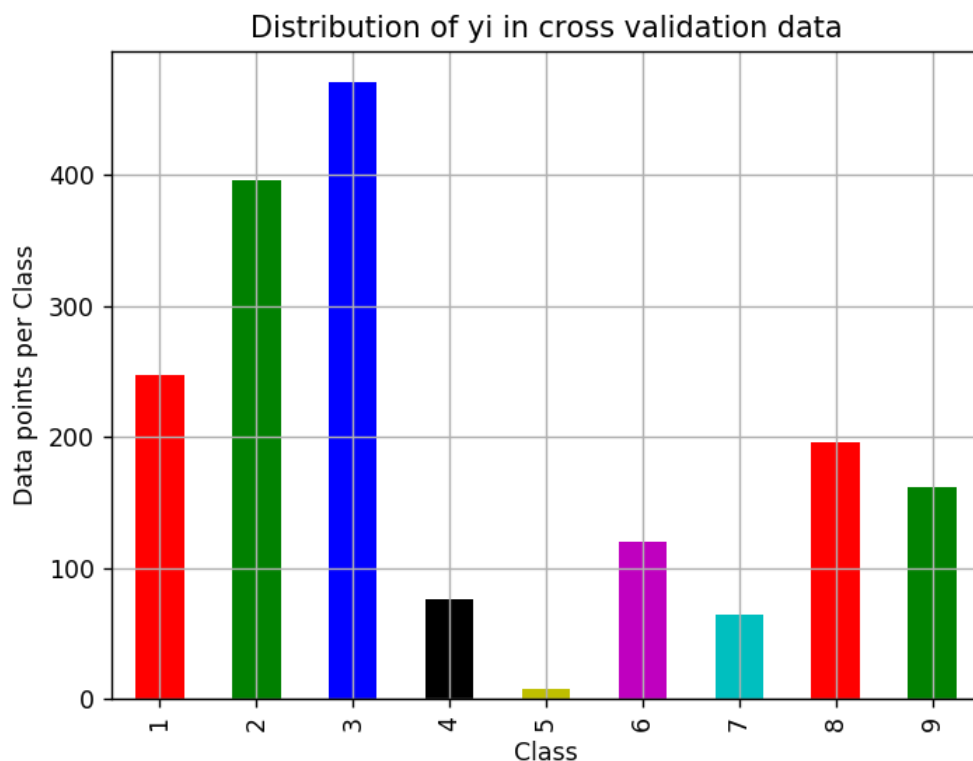
In [0]:
```python
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    print("Number of misclassified points ",(len(test_y)-np.trace(C
    # C = 9,9 matrix, each cell (i,j) represents number of points o

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

    # C = [[1, 2],
    #      [3, 4]]
    # C.T = [[1, 3],
    #         [2, 4]]
    # C.sum(axis = 1)   axis=0 corresonds to columns and axis=1 corr
    # C.sum(axix =1) = [[3, 7]]
    # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
    #                            [2/3, 4/7]]

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

    B =(C/C.sum(axis=0))
    #divid each element of the confusion matrix with the sum of ele
    # C = [[1, 2],
    #      [3, 4]]
    # C.sum(axis = 0)   axis=0 corresonds to columns and axis=1 corr
    # C.sum(axix =0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                      [3/4, 4/6]]

    labels = [1,2,3,4,5,6,7,8,9]
    cmap=sns.light_palette("green")
    # representing A in heatmap format
    print("-"*50, "Confusion matrix", "-"*50)
    plt.figure(figsize=(10,5))
    sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=la
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()

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

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

```
    plt.show()
    print("Sum of rows in precision matrix" A sum(axis=1))
```

# 4. Machine Learning Models

## 4.1. Machine Leaning Models on bytes files

### 4.1.1. Random Model

In [0]:
```python
# we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to genarate 9 numbers and divide each of the numb
# ref: https://stackoverflow.com/a/18662466/4084039

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

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


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

predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)
```
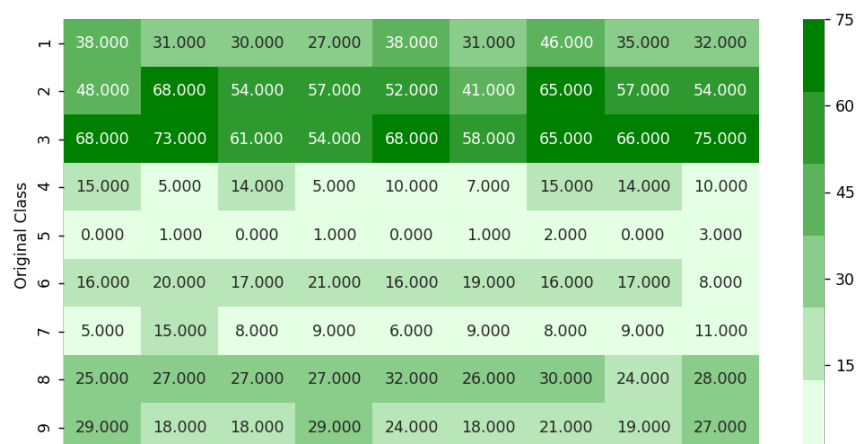
```
Log loss on Cross Validation Data using Random Model 2.45615644965
Log loss on Test Data using Random Model 2.48503905509
Number of misclassified points  88.5004599816
-------------------------------------------------- Confusion matri
x --------------------------------------------------

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

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
|---|---|---|---|---|---|---|---|---|---|---|

Predicted Class

------------------------------------------------ Precision matri
x ------------------------------------------------

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



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

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



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

### 4.1.2. K Nearest Neighbour Classification

```
In [0]: alpha = [x for x in range(1, 15, 2)]
        cv_log_error_array=[]
        for i in alpha:
            k_cfl=KNeighborsClassifier(n_neighbors=i)
            k_cfl.fit(X_train,y_train)
```

```python
        sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
        sig_clf.fit(X_train, y_train)
        predict_y = sig_clf.predict_proba(X_cv)
        cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cf

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

best_alpha = np.argmin(cv_log_error_array)

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

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

predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test l
plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```

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

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



Cross Validation Error for each alpha

For values of best alpha =  1 The train log loss is: 0.07829476692
47
For values of best alpha =  1 The cross validation log loss is: 0.
225386237304
For values of best alpha =  1 The test log loss is: 0.241508604195
Number of misclassified points  4.50781968721
------------------------------------------------- Confusion matri
x -------------------------------------------------

<IPython.core.display.Javascript object>



------------------------------------------------- Precision matri
x -------------------------------------------------

<IPython.core.display.Javascript object>



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

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



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

### 4.1.3. Logistic Regression

In [0]:

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

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

best_alpha = np.argmin(cv_log_error_array)

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

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

```
predict_y = sig_clf.predict_proba(X_train)
print ('log loss for train data',log_loss(y_train, predict_y, label
predict_y = sig_clf.predict_proba(X_cv)
print ('log loss for cv data',log_loss(y_cv, predict_y, labels=logi
predict_y = sig_clf.predict_proba(X_test)
print ('log loss for test data',log_loss(y_test, predict_y, labels=
```

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

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



Cross Validation Error for each alpha

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

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

| Original | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 4.000 | 2.000 | 1.000 | 1.000 |
| 6 | 10.000 | 1.000 | 1.000 | 10.000 | 0.000 | 102.000 | 0.000 | 19.000 | 7.000 |
| 7 | 2.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 76.000 | 0.000 | 1.000 |
| 8 | 18.000 | 0.000 | 0.000 | 7.000 | 0.000 | 3.000 | 0.000 | 217.000 | 1.000 |
| 9 | 27.000 | 5.000 | 0.000 | 2.000 | 0.000 | 8.000 | 0.000 | 18.000 | 143.000 |

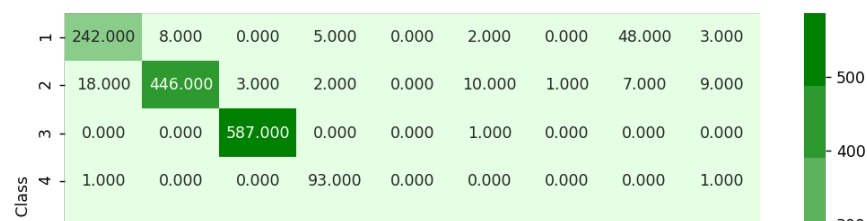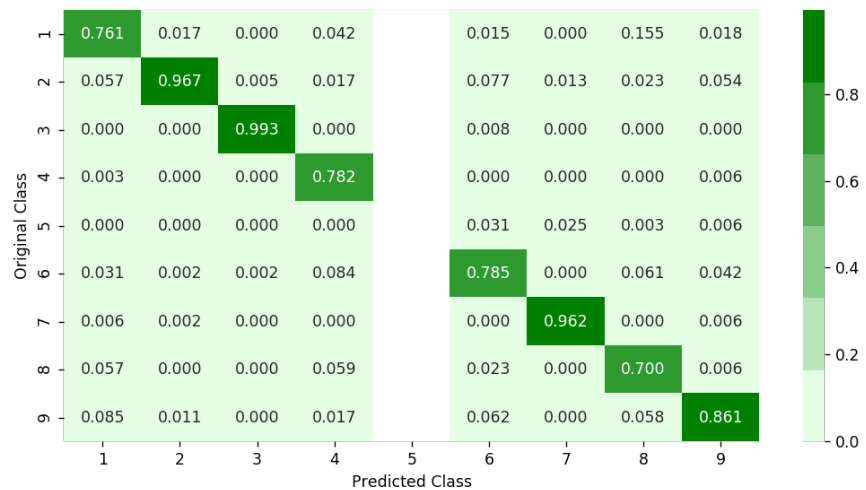Predicted Class

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

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

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

Predicted Class

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

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

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

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

Predicted Class

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

### 4.1.4. Random Forest Classifier

In [0]:

```python
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
train_log_error_array=[]
from sklearn.ensemble import RandomForestClassifier
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_j
    r_cfl.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cf

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


best_alpha = np.argmin(cv_log_error_array)

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


r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_
r_cfl.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train, y_train)

predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test l
plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```
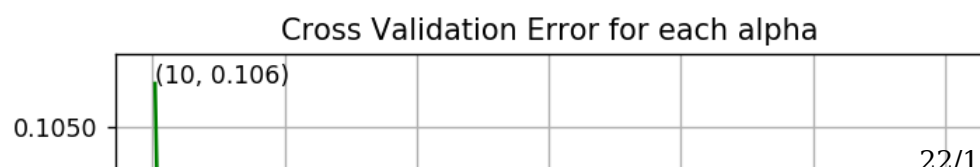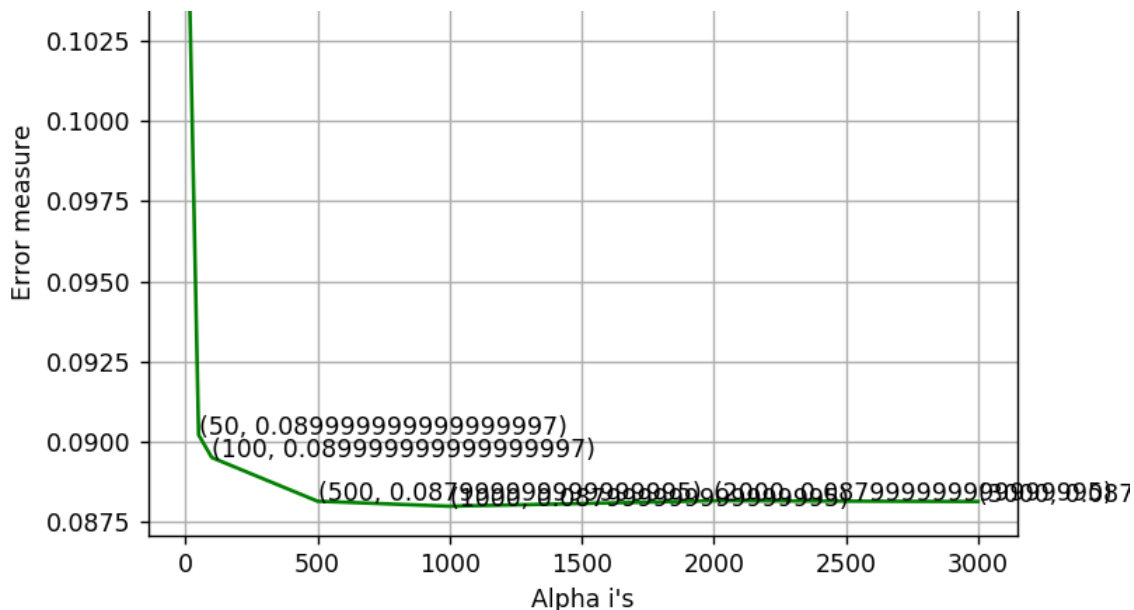
```
log_loss for c =  10 is 0.106357709164
log_loss for c =  50 is 0.0902124124145
log_loss for c =  100 is 0.0895043339776
log_loss for c =  500 is 0.0881420869288
log_loss for c =  1000 is 0.0879849524621
log_loss for c =  2000 is 0.0881566647295
log_loss for c =  3000 is 0.0881318948443

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

Cross Validation Error for each alpha

(10, 0.106)

0.1050

```
For values of best alpha =  1000 The train log loss is: 0.02664762
91801
For values of best alpha =  1000 The cross validation log loss is:
0.0879849524621
For values of best alpha =  1000 The test log loss is: 0.085834696
1407
Number of misclassified points  2.02391904324
------------------------------------------------- Confusion matri
x -------------------------------------------------
```

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



```
------------------------------------------------- Precision matri
x -------------------------------------------------
```

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

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

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



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

### 4.1.5. XgBoost Classification

In [0]:

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

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


best_alpha = np.argmin(cv_log_error_array)

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
```

```python
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

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

predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test l
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>
```



```
For values of best alpha =   500 The train log loss is: 0.022523180
5824
For values of best alpha =   500 The cross validation log loss is:
0.0931035681289
For values of best alpha =   500 The test log loss is: 0.0792067651
731
```

```
Number of misclassified points  1.24195032199
------------------------------------------------- Confusion matri
x -------------------------------------------------
```

<IPython.core.display.Javascript object>

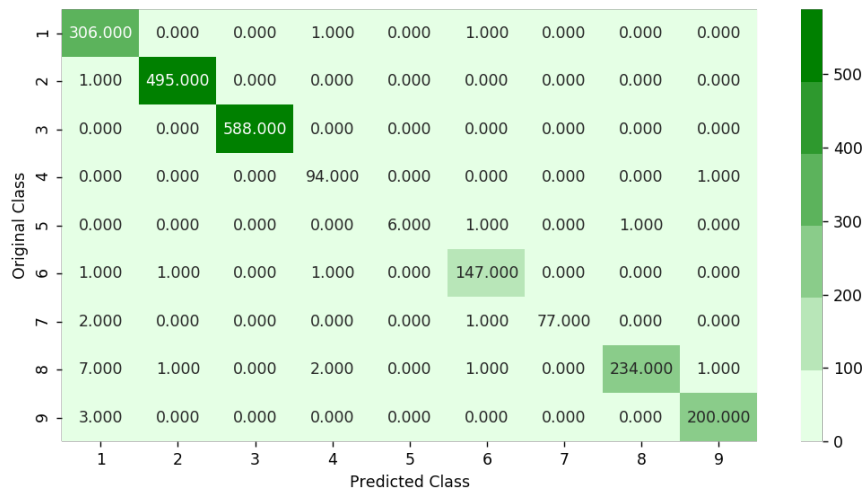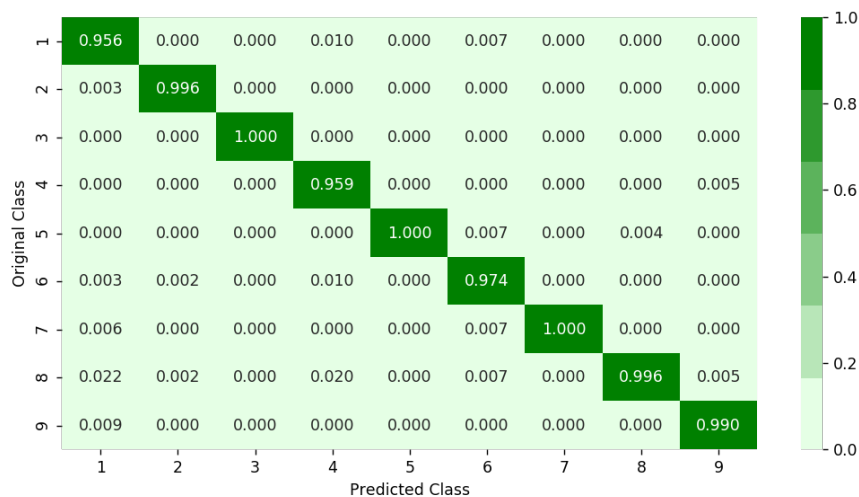| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 306.000 | 0.000 | 0.000 | 1.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 |
| 2 | 1.000 | 495.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 588.000 | 0.000 | 0.000 | 0.000 | 0.000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

Original Class (y-axis), Predicted Class (x-axis)

```
------------------------------------------------- Precision matri
x -------------------------------------------------
```

<IPython.core.display.Javascript object>

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

Original Class (y-axis), Predicted Class (x-axis)

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

<IPython.core.display.Javascript object>

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.994 | 0.000 | 0.000 | 0.003 | 0.000 | 0.003 | 0.000 | 0.000 | 0.000 |
| 2 | 0.002 | 0.998 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 0.000 | 0.000 | 0.000 | 0.989 | 0.000 | 0.000 | 0.000 | 0.000 | 0.011 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.750 | 0.125 | 0.000 | 0.125 | 0.000 |

Predicted Class

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

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

```python
# https://www.analyticsvidhya.com/blog/2016/03/complete-guide-param
x_cfl=XGBClassifier()

prams={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
     'n_estimators':[100,200,500,1000,2000],
     'max_depth':[3,5,10],
     'colsample_bytree':[0.1,0.3,0.5,1],
     'subsample':[0.1,0.3,0.5,1]
}
random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,verb
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 rema
ining:  5.4min
[Parallel(n_jobs=-1)]: Done   23 out of  30 | elapsed: 10.1min rema
ining:  3.1min
[Parallel(n_jobs=-1)]: Done   27 out of  30 | elapsed: 14.0min rema
ining:  1.6min
[Parallel(n_jobs=-1)]: Done   30 out of  30 | elapsed: 14.2min fini
shed
```

Out[75]:
```
RandomizedSearchCV(cv=None, error_score='raise',
          estimator=XGBClassifier(base_score=0.5, colsample_byleve
l=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], 'subs
ample': [0.1, 0.3, 0.5, 1]},
          pre_dispatch='2*n_jobs', random_state=None, refit=True,
          return_train_score=True, scoring=None, verbose=10)
```

In [0]: print (random_cfl1.best_params_)

```
{'subsample': 1, 'n_estimators': 500, 'max_depth': 5, 'learning_ra
```

In [0]:
```python
x_cfl=XGBClassifier(n_estimators=2000, learning_rate=0.05, colsampl
x_cfl.fit(X_train,y_train)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train,y_train)

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

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

# 4.2 Modeling with .asm files

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


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

We read the top solutions and handpicked the features from
those papers/videos/blogs.
 Refer:https://www.kaggle.com/c/malware-classification/disc
ussion
```

### 4.2.1 Feature extraction from asm files

To extract the unigram features from the .asm files we need to process ~150GB of data

**Note: Below two cells will take lot of time (over 48 hours to complete)**

We will provide you the output file of these two cells, which you can directly use it

```
In [0]: #intially create five folders
        #first
        #second
        #thrid
        #fourth
        #fifth
        #this code tells us about random split of files into five folders
        folder_1 ='first'
        folder_2 ='second'
        folder_3 ='third'
        folder_4 ='fourth'
        folder_5 ='fifth'
        folder_6 = 'output'
        for i in [folder_1,folder_2,folder_3,folder_4,folder_5,folder_6]:
            if not os.path.isdir(i):
                os.makedirs(i)

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

```
In [0]: #http://flint.cs.yale.edu/cs421/papers/x86-asm/asm.html

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

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

            opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn',
            #best keywords that are taken from different blogs
            keywords = ['.dll','std::',':dword']
            #Below taken registers are general purpose registers and specia
            #All the registers which are taken are best
            registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip
            file1=open("output\asmsmallfile.txt","w+")
            files = os.listdir('first')
            for f in files:
                #filling the values with zeros into the arrays
                prefixescount=np.zeros(len(prefixes),dtype=int)
```

```python
            opcodescount=np.zeros(len(opcodes),dtype=int)
            keywordcount=np.zeros(len(keywords),dtype=int)
            registerscount=np.zeros(len(registers),dtype=int)
            features=[]
            f2=f.split('.')[0]
            file1.write(f2+",")
            opcodefile.write(f2+" ")
            # https://docs.python.org/3/library/codecs.html#codecs.igno
            # https://docs.python.org/3/library/codecs.html#codecs.Code
            with codecs.open('first/'+f,encoding='cp1252',errors ='repl
                for lines in fli:
                    # https://www.tutorialspoint.com/python3/string_rst
                    line=lines.rstrip().split()
                    l=line[0]
                    #counting the prefixs in each and every line
                    for i in range(len(prefixes)):
                        if prefixes[i] in line[0]:
                            prefixescount[i]+=1
                    line=line[1:]
                    #counting the opcodes in each and every line
                    for i in range(len(opcodes)):
                        if any(opcodes[i]==li for li in line):
                            features.append(opcodes[i])
                            opcodescount[i]+=1
                    #counting registers in the line
                    for i in range(len(registers)):
                        for li in line:
                            # we will use registers only in 'text' and
                            if registers[i] in li and ('text' in l or '
                                registerscount[i]+=1
                    #counting keywords in the line
                    for i in range(len(keywords)):
                        for li in line:
                            if keywords[i] in li:
                                keywordcount[i]+=1
            #pushing the values into the file after reading whole file
            for prefix in prefixescount:
                file1.write(str(prefix)+",")
            for opcode in opcodescount:
                file1.write(str(opcode)+",")
            for register in registerscount:
                file1.write(str(register)+",")
            for key in keywordcount:
                file1.write(str(key)+",")
            file1.write("\n")
        file1.close()


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

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

    # same as smallprocess() functions
    def thirdprocess():
        prefixes = ['HEADER:','.text:','.Pav:','.idata:','.data:','.bss
        opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn',
        keywords = ['.dll','std::',':dword']
        registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip
        file1=open("output\largeasmfile.txt","w+")
        files = os.listdir('thrid')
        for f in files:
            prefixescount=np.zeros(len(prefixes),dtype=int)
            opcodescount=np.zeros(len(opcodes),dtype=int)
            keywordcount=np.zeros(len(keywords),dtype=int)
            registerscount=np.zeros(len(registers),dtype=int)
            features=[]
            f2=f.split('.')[0]
            file1.write(f2+",")
            opcodefile.write(f2+" ")
            with codecs.open('thrid/'+f,encoding='cp1252',errors ='repl
                for lines in fli:
                    line=lines.rstrip().split()
                    l=line[0]
                    for i in range(len(prefixes)):
```

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


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

```python
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
        for prefix in prefixescount:
            file1.write(str(prefix)+",")
        for opcode in opcodescount:
            file1.write(str(opcode)+",")
        for register in registerscount:
            file1.write(str(register)+",")
        for key in keywordcount:
            file1.write(str(key)+",")
        file1.write("\n")
    file1.close()


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

```python
                    file1.write(str(key)+",")
                file1.write("\n")
            file1.close()


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

In [0]:
```python
# asmoutputfile.csv(output genarated from the above two cells) will
if __name__ == "__main__":
#this file will be uploaded in the drive, you can directly use thi
    main()
dfasm=pd.read_csv("asmoutputfile.csv")
Y.columns = ['ID', 'Class']
result_asm = pd.merge(dfasm, Y,on='ID', how='left')
result_asm.head()
```

Out[137]:

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

5 rows × 53 columns

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

In [0]:
```python
#file sizes of byte files

files=os.listdir('asmFiles')
filenames=Y['ID'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
```

```
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev
    # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st
    # read more about os.stat: here https://www.tutorialspoint.com/
    statinfo=os.stat('asmFiles/'+file)
    # split the file name at '.' and take the first part of it i.e
    file=file.split('.')[0]
    if any(file == filename for filename in filenames):
        i=filenames.index(file)
        class_bytes.append(class_y[i])
        # converting into Mb's
        sizebytes.append(statinfo.st_size/(1024.0*1024.0))
        fnames.append(file)
asm_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':cl
print (asm size byte head())
```

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

**4.2.1.2 Distribution of .asm file sizes**

In [0]:
```
#boxplot of asm files
ax = sns.boxplot(x="Class", y="size", data=asm_size_byte)
plt.title("boxplot of .bytes file sizes")
plt show()
```

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



In [0]:
```
# add the file size feature to previous extracted features
print(result_asm.shape)
print(asm_size_byte.shape)
```

```
result_asm = pd.merge(result_asm, asm_size_byte.drop(['Class'], axi
(10868, 53)head()
```
(10868, 3)

Out[140]:

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

5 rows × 54 columns

```
In [0]:   # we normalize the data each column
          result_asm = normalize(result_asm)
          result_asm.head()
```

Out[145]:

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

5 rows × 54 columns
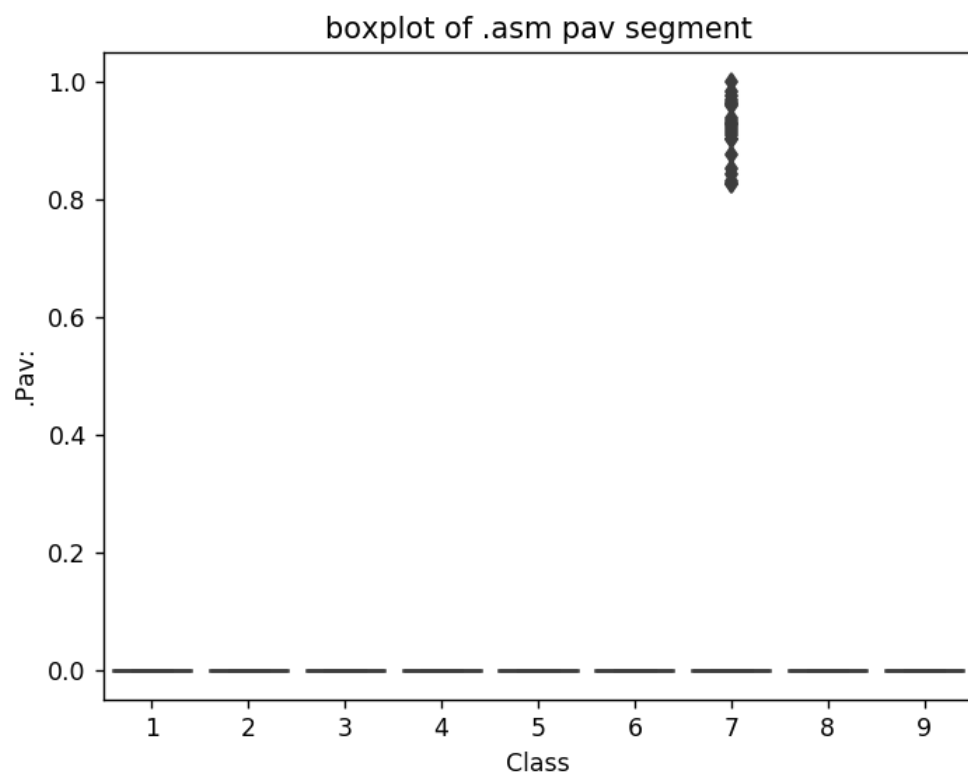
### 4.2.2 Univariate analysis on asm file features

```
In [0]:   ax = sns.boxplot(x="Class", y=".text:", data=result_asm)
          plt.title("boxplot of .asm text segment")
          plt.show()
```
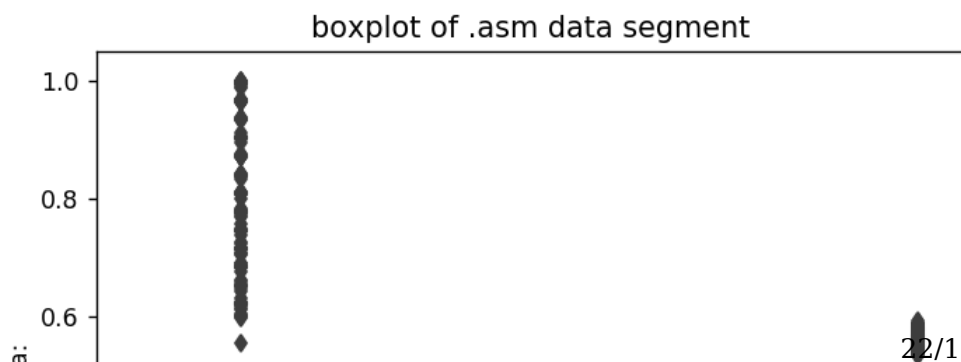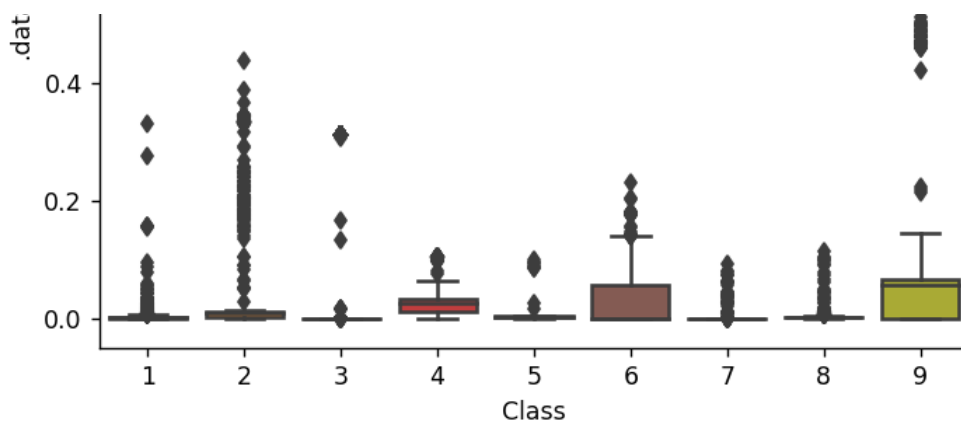<IPython.core.display.Javascript object>



boxplot of .asm text segment

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

In [0]: 
```
ax = sns.boxplot(x="Class", y=".Pav:", data=result_asm)
plt.title("boxplot of .asm pav segment")
plt show()
```
<IPython.core.display.Javascript object>
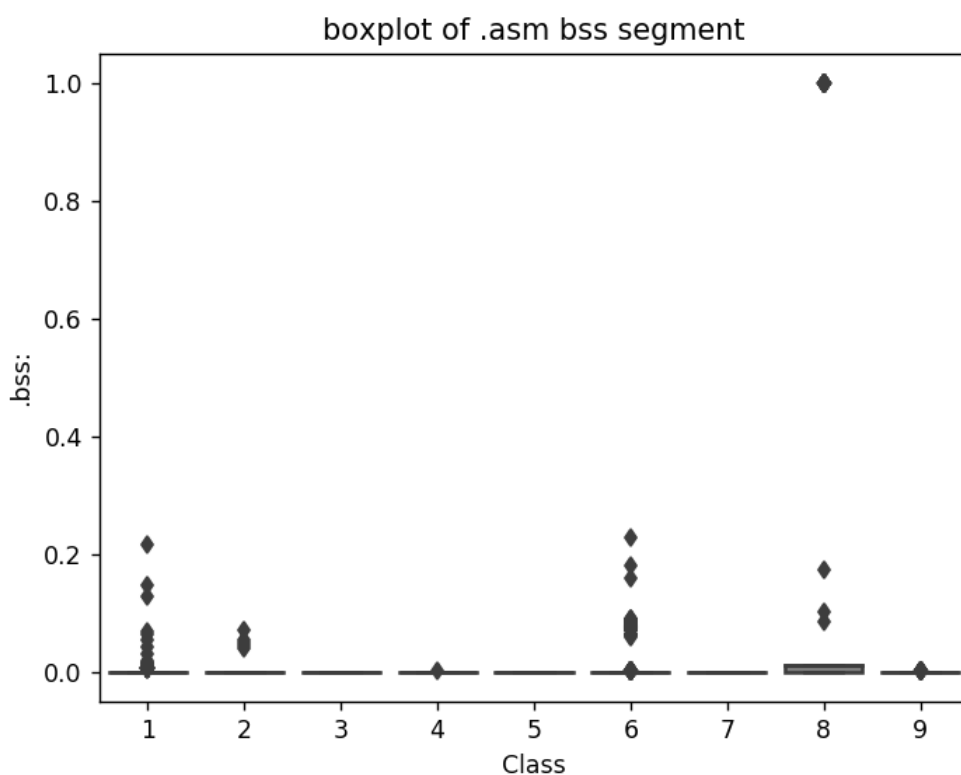


In [0]: 
```
ax = sns.boxplot(x="Class", y=".data:", data=result_asm)
plt.title("boxplot of .asm data segment")
plt show()
```
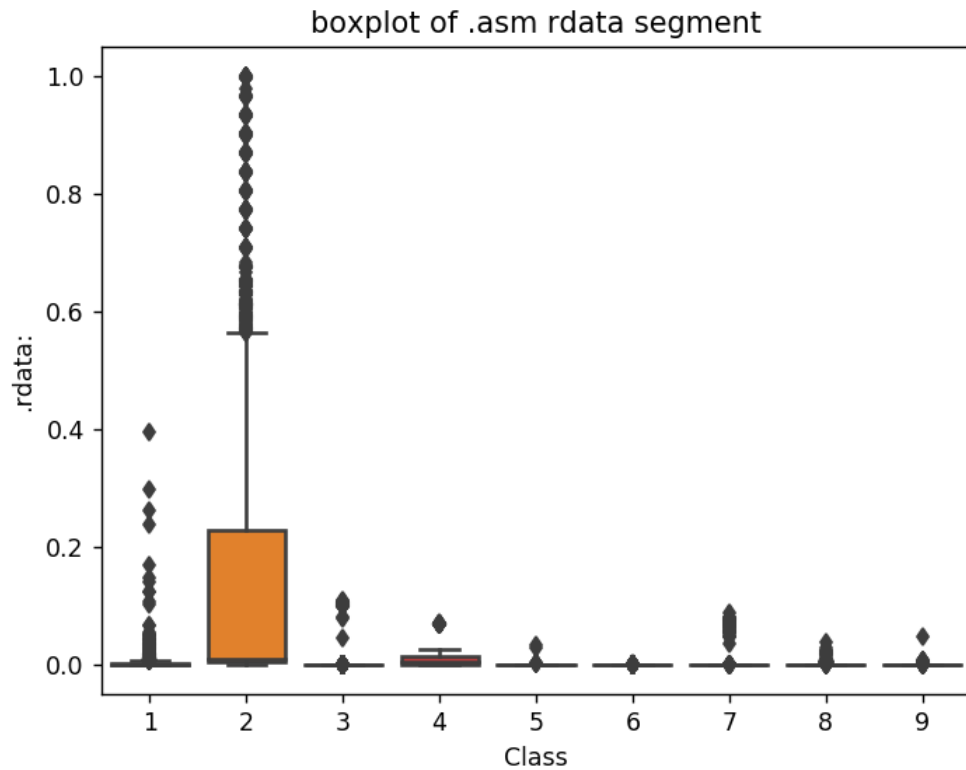<IPython.core.display.Javascript object>

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

```
In [0]: ax = sns.boxplot(x="Class", y=".bss:", data=result_asm)
        plt.title("boxplot of .asm bss segment")
        plt.show()
```
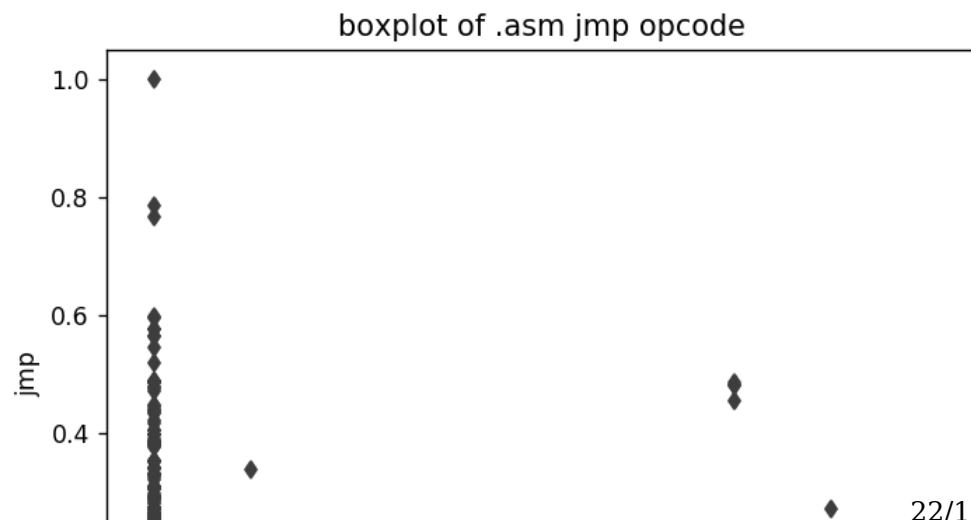
<IPython.core.display.Javascript object>



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

In [0]:
```
ax = sns.boxplot(x="Class", y=".rdata:", data=result_asm)
plt.title("boxplot of .asm rdata segment")
plt.show()
```
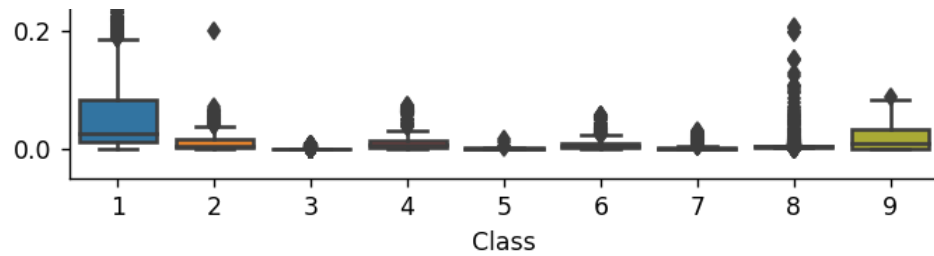
<IPython.core.display.Javascript object>

boxplot of .asm rdata segment

Plot between rdata segment and Class segment
Class 2 can be easily separated 75 pecentile files are havi
ng 1M rdata lines

In [0]:
```
ax = sns.boxplot(x="Class", y="jmp", data=result_asm)
plt.title("boxplot of .asm jmp opcode")
plt.show()
```

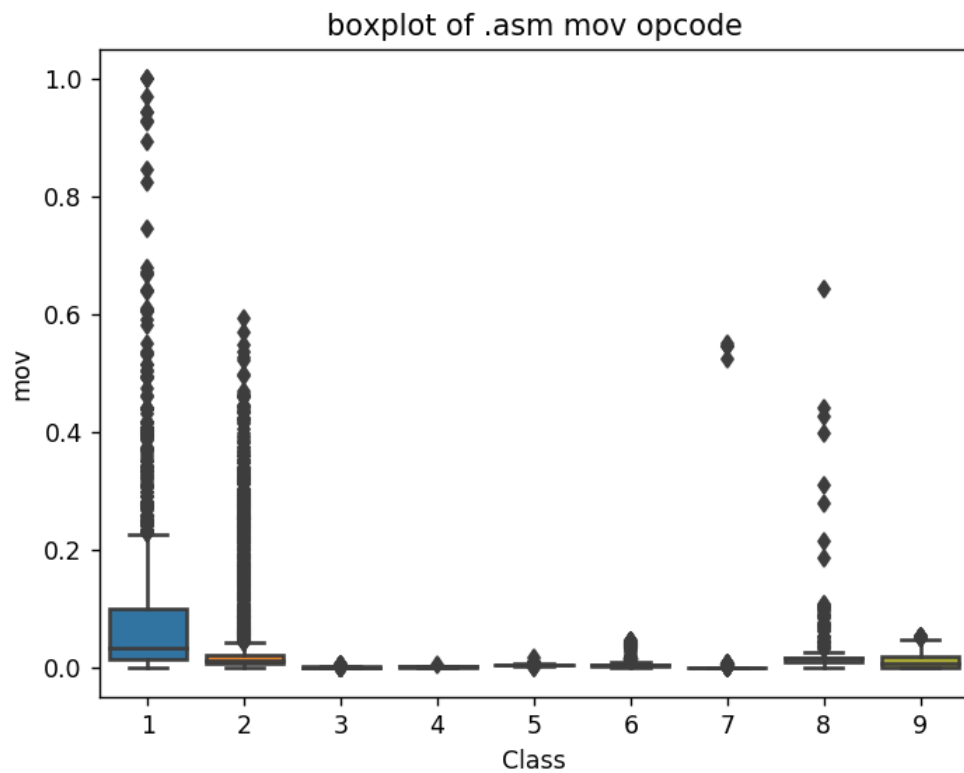<IPython.core.display.Javascript object>

boxplot of .asm jmp opcode

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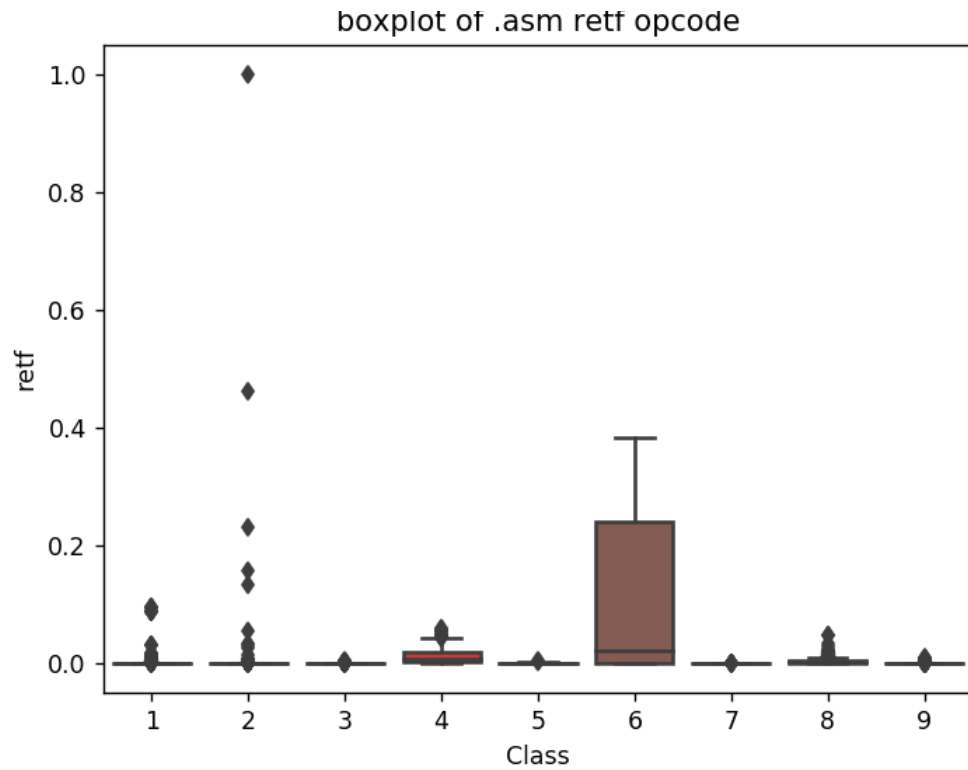of files

In [0]: 
```
ax = sns.boxplot(x="Class", y="mov", data=result_asm)
plt.title("boxplot of .asm mov opcode")
plt.show()
```
<IPython.core.display.Javascript object>
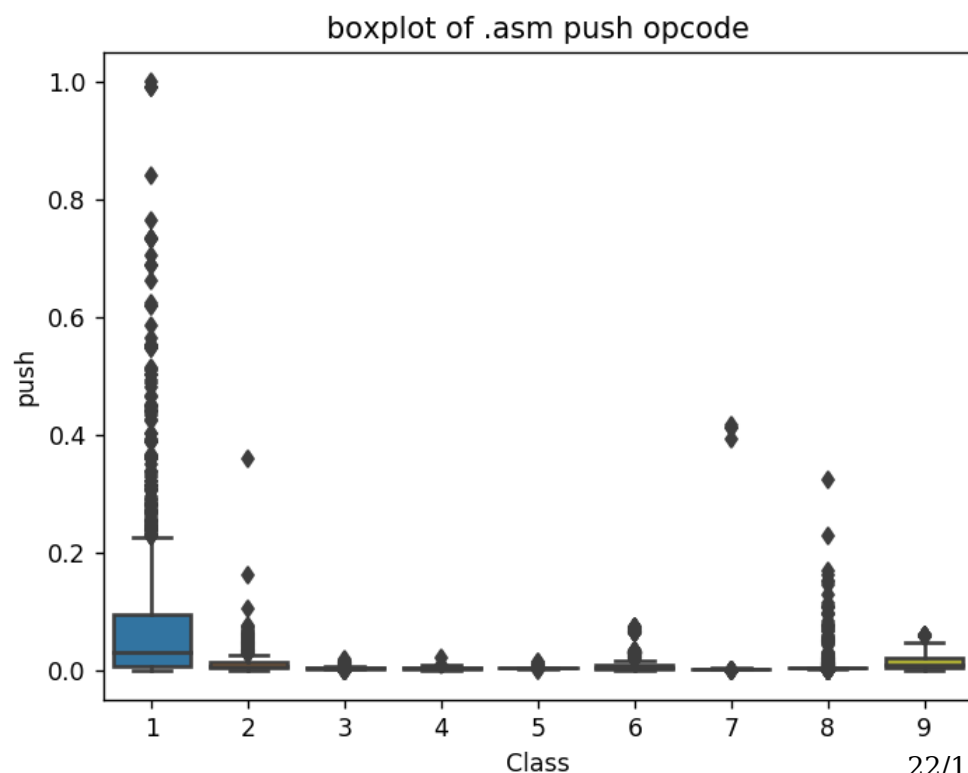


plot between Class label and mov opcode
Class 1 is having frequency of 2000 approx in 75 perentile
of files

In [0]: 
```
ax = sns.boxplot(x="Class", y="retf", data=result_asm)
plt.title("boxplot of .asm retf opcode")
plt.show()
```
<IPython.core.display.Javascript object>

### boxplot of .asm retf opcode



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

In [0]:
```
ax = sns.boxplot(x="Class", y="push", data=result_asm)
plt.title("boxplot of .asm push opcode")
plt show()
```
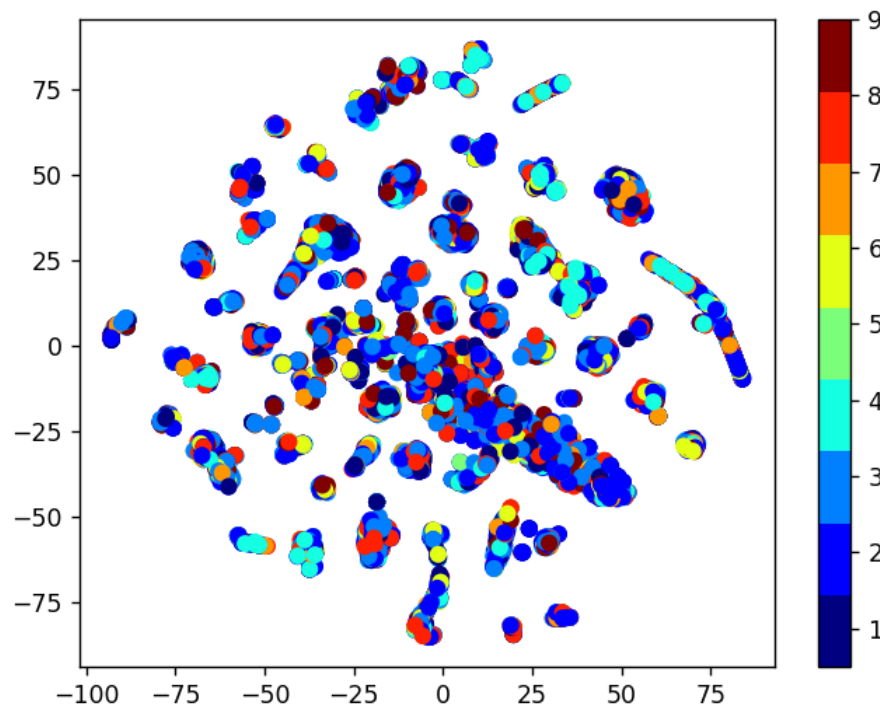
<IPython.core.display.Javascript object>

### boxplot of .asm push opcode

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

### 4.2.2 Multivariate Analysis on .asm file features

In [0]:
```
xtsne=TSNE(perplexity=50)
results=xtsne.fit_transform(result_asm.drop(['ID','Class'], axis=1)
vis_x = results[:, 0]
vis_y = results[:, 1   ]
plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
plt.colorbar(ticks=range(10))
plt.clim(0.5, 9)
plt show()
```
```
<IPython.core.display.Javascript object>
```



In [0]:
```
# by univariate analysis on the .asm file features we are getting v
# 'rtn', '.BSS:' '.CODE' features, so heare we are trying multivari
# the plot looks very messy

xtsne=TSNE(perplexity=30)
results=xtsne.fit_transform(result_asm.drop(['ID','Class', 'rtn', '
vis_x = results[:, 0]
vis_y = results[:, 1]
plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
plt.colorbar(ticks=range(10))
plt.clim(0.5, 9)
```
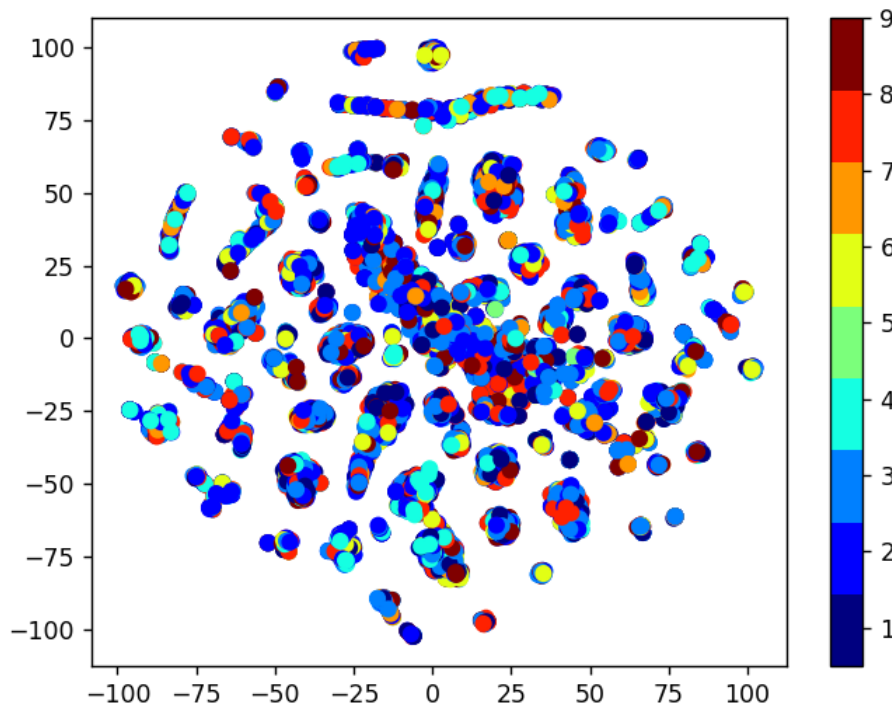
```
plt show()
<IPython.core.display.Javascript object>
```



```
TSNE for asm data with perplexity 50
```

### 4.2.3 Conclusion on EDA

We have taken only 52 features from asm files (after reading through many blogs and research papers)
The univariate analysis was done only on few important features.
Take-aways
  - 1. Class 3 can be easily separated because of the frequency of segments,opcodes and keywords being less
  - 2. Each feature has its unique importance in separating the Class labels.

## 4.3 Train and test split

```
In [0]: asm_y = result_asm['Class']
        asm x = result asm drop(['ID' 'Class' ' BSS·' 'rtn' ' CODE'l  axis=
```

```
In [0]: X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split
        X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_t
```

```
In [0]: print( X cv asm isnull() all())
```

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

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

### 4.4.1 K-Nearest Neigbors

In [0]:

```python
alpha = [x for x in range(1, 21,2)]
```

22/12/20, 9:17 pm

```python
cv_log_error_array=[]
for i in alpha:
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=

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

best_alpha = np.argmin(cv_log_error_array)

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

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


predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',log_loss(y_train_asm, predict_y))
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',log_loss(y_cv_asm, predict_y))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',log_loss(y_test_asm, predict_y))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```
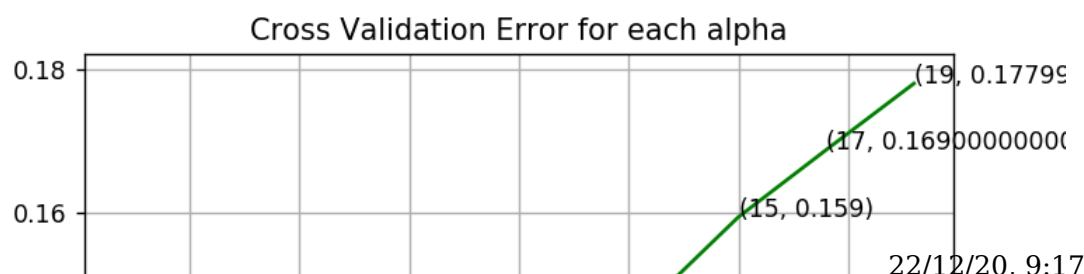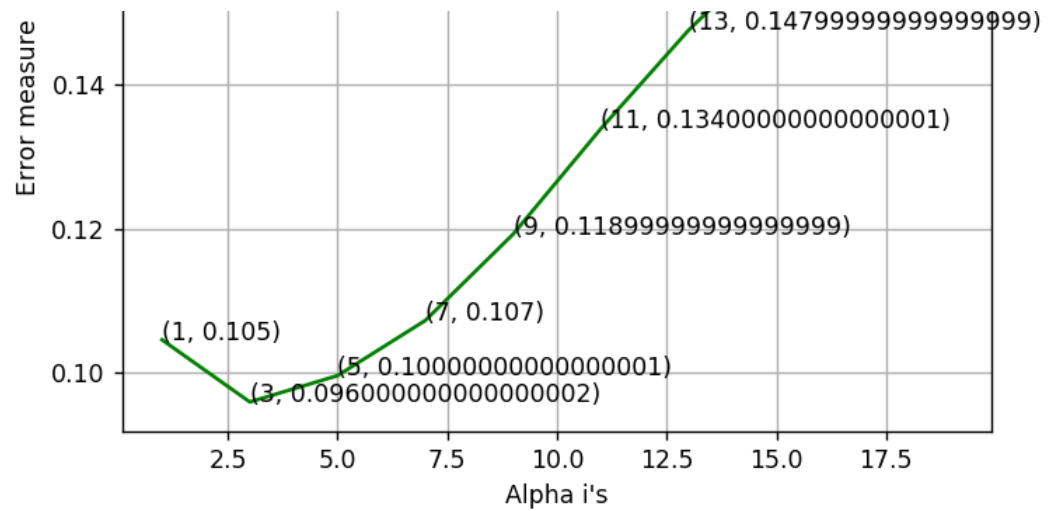
```
log_loss for k =   1 is 0.104531321344
log_loss for k =   3 is 0.0958800580948
log_loss for k =   5 is 0.0995466557335
log_loss for k =   7 is 0.107227274345
log_loss for k =   9 is 0.119239543547
log_loss for k =   11 is 0.133926642781
log_loss for k =   13 is 0.147643793967
log_loss for k =   15 is 0.159439699615
log_loss for k =   17 is 0.16878376444
log_loss for k =   19 is 0.178020728839

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

Cross Validation Error for each alpha

0.18

(19, 0.17799

(17, 0.1690000000(

0.16

(15, 0.159)

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

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



```
------------------------------------------------- Precision matri
x -------------------------------------------------

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

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

<IPython.core.display.Javascript object>



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

### 4.4.2 Logistic Regression

In [0]:

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

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

best_alpha = np.argmin(cv_log_error_array)

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

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

predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y,
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',(log_loss(y_test_asm, predict_y, la
plot confusion matrix(y test asm sig clf predict(X test asm))
```
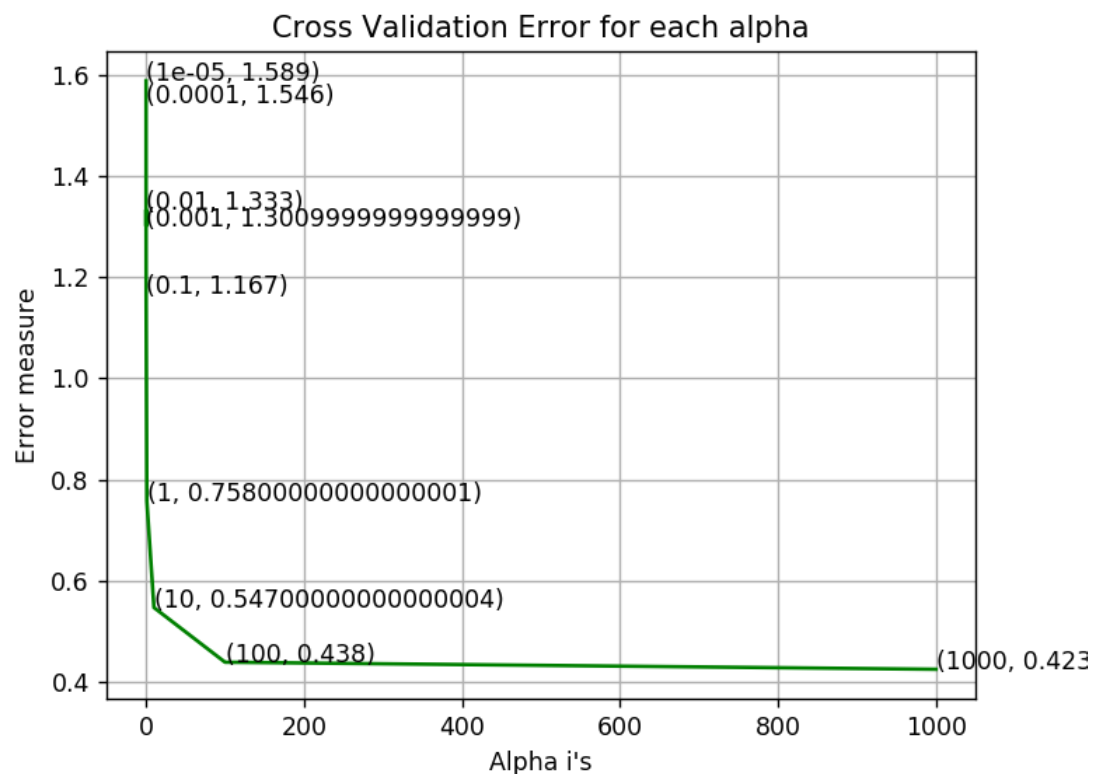
```
log_loss for c =  1e-05 is 1.58867274165
log_loss for c =  0.0001 is 1.54560797884
log_loss for c =  0.001 is 1.30137786807
log_loss for c =  0.01 is 1.33317456931
log_loss for c =  0.1 is 1.16705751378
log_loss for c =  1 is 0.757667807779
log_loss for c =  10 is 0.546533939819
log_loss for c =  100 is 0.438414998062
log_loss for c =  1000 is 0.424423536526
```

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

Cross Validation Error for each alpha

(1e-05, 1.589)
(0.0001, 1.546)
(0.01, 1.333)
(0.001, 1.3009999999999999)
(0.1, 1.167)
(1, 0.75800000000000001)
(10, 0.54700000000000004)
(100, 0.438)
(1000, 0.423

Error measure / Alpha i's

```
log loss for train data 0.396219394701
log loss for cv data 0.424423536526
log loss for test data 0.415685592517
Number of misclassified points  9.61361545538
--------------------------------------------------- Confusion matri
x -------------------------------------------------
```
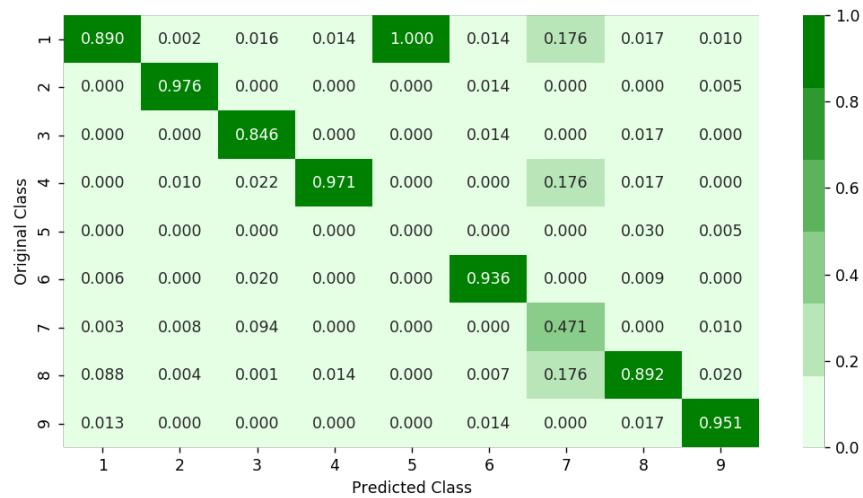
```
<IPython.core.display.Javascript object>
```

283.000  1.000  11.000  1.000  1.000  2.000  3.000  4.000  2.000

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| **2** | 0.000 | 493.000 | 0.000 | 0.000 | 0.000 | 2.000 | 0.000 | 0.000 | 1.000 |
| **3** | 0.000 | 0.000 | 582.000 | 0.000 | 0.000 | 2.000 | 0.000 | 4.000 | 0.000 |
| **4** | 0.000 | 5.000 | 15.000 | 68.000 | 0.000 | 0.000 | 3.000 | 4.000 | 0.000 |
| **5** | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 7.000 | 1.000 |
| **6** | 2.000 | 0.000 | 14.000 | 0.000 | 0.000 | 132.000 | 0.000 | 2.000 | 0.000 |
| **7** | 1.000 | 4.000 | 65.000 | 0.000 | 0.000 | 0.000 | 8.000 | 0.000 | 2.000 |
| **8** | 28.000 | 2.000 | 1.000 | 1.000 | 0.000 | 1.000 | 3.000 | 206.000 | 4.000 |
| **9** | 4.000 | 0.000 | 0.000 | 0.000 | 0.000 | 2.000 | 0.000 | 4.000 | 193.000 |

Original Class (y-axis), Predicted Class (x-axis)

------------------------------------------------- Precision matri
x -------------------------------------------------

<IPython.core.display.Javascript object>

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| **1** | 0.890 | 0.002 | 0.016 | 0.014 | 1.000 | 0.014 | 0.176 | 0.017 | 0.010 |
| **2** | 0.000 | 0.976 | 0.000 | 0.000 | 0.000 | 0.014 | 0.000 | 0.000 | 0.005 |
| **3** | 0.000 | 0.000 | 0.846 | 0.000 | 0.000 | 0.014 | 0.000 | 0.017 | 0.000 |
| **4** | 0.000 | 0.010 | 0.022 | 0.971 | 0.000 | 0.000 | 0.176 | 0.017 | 0.000 |
| **5** | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.030 | 0.005 |
| **6** | 0.006 | 0.000 | 0.020 | 0.000 | 0.000 | 0.936 | 0.000 | 0.009 | 0.000 |
| **7** | 0.003 | 0.008 | 0.094 | 0.000 | 0.000 | 0.000 | 0.471 | 0.000 | 0.010 |
| **8** | 0.088 | 0.004 | 0.001 | 0.014 | 0.000 | 0.007 | 0.176 | 0.892 | 0.020 |
| **9** | 0.013 | 0.000 | 0.000 | 0.000 | 0.000 | 0.014 | 0.000 | 0.017 | 0.951 |

Original Class (y-axis), Predicted Class (x-axis)

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

<IPython.core.display.Javascript object>

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

Original Class (y-axis), Predicted Class (x-axis)

1.]

### 4.4.3 Random Forest Classifier

In [0]:
```python
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_j
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=

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


best_alpha = np.argmin(cv_log_error_array)

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

r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y,
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',(log_loss(y_test_asm, predict_y, la
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```
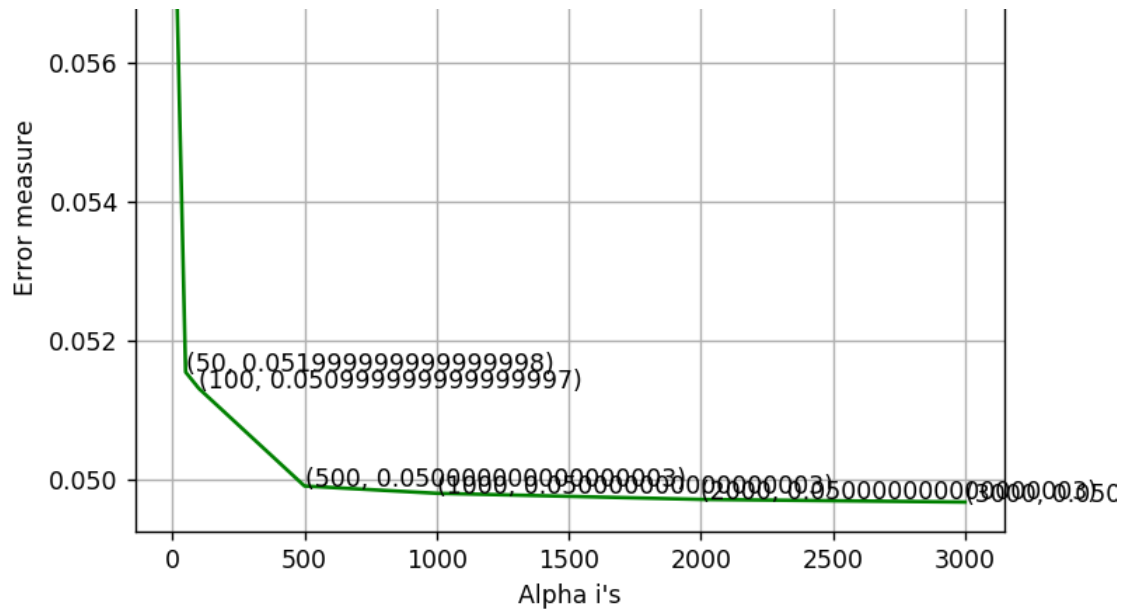
```
log_loss for c =   10 is 0.0581657906023
log_loss for c =   50 is 0.0515443148419
log_loss for c =   100 is 0.0513084973231
log_loss for c =   500 is 0.0499021761479
log_loss for c =   1000 is 0.0497972474298
log_loss for c =   2000 is 0.0497091690815
log_loss for c =   3000 is 0.0496706817633

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

Cross Validation Error for each alpha

(10, 0.058000000000000003)

0.058

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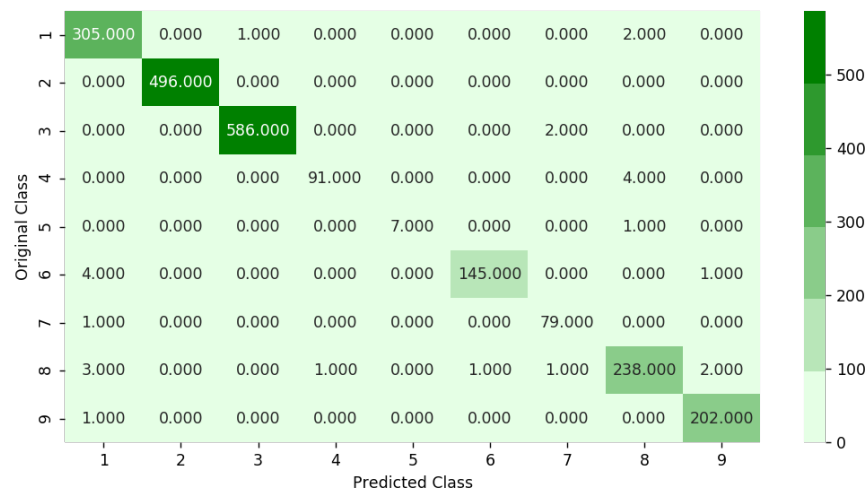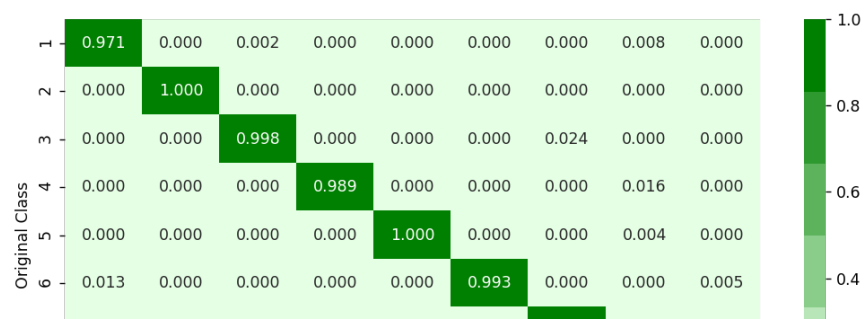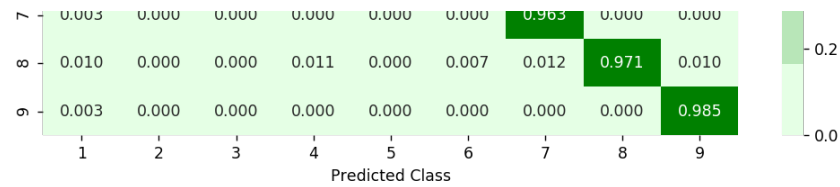
<IPython.core.display.Javascript object>



```
------------------------------------------------- Precision matri
x -------------------------------------------------
```

<IPython.core.display.Javascript object>

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 7 | 0.003 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.963 | 0.000 | 0.000 |
| 8 | 0.010 | 0.000 | 0.000 | 0.011 | 0.000 | 0.007 | 0.012 | 0.971 | 0.010 |
| 9 | 0.003 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.985 |

Predicted Class

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

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

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.990 | 0.000 | 0.003 | 0.000 | 0.000 | 0.000 | 0.000 | 0.006 | 0.000 |
| 2 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.000 | 0.997 | 0.000 | 0.000 | 0.000 | 0.003 | 0.000 | 0.000 |
| 4 | 0.000 | 0.000 | 0.000 | 0.958 | 0.000 | 0.000 | 0.000 | 0.042 | 0.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.875 | 0.000 | 0.000 | 0.125 | 0.000 |
| 6 | 0.027 | 0.000 | 0.000 | 0.000 | 0.000 | 0.967 | 0.000 | 0.000 | 0.007 |
| 7 | 0.013 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.988 | 0.000 | 0.000 |
| 8 | 0.012 | 0.000 | 0.000 | 0.004 | 0.000 | 0.004 | 0.004 | 0.967 | 0.008 |
| 9 | 0.005 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.995 |

Original Class / Predicted Class

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

### 4.4.4 XgBoost Classifier

In [0]:

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

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


best_alpha = np.argmin(cv_log_error_array)

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

```
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

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

predict_y = sig_clf.predict_proba(X_train_asm)

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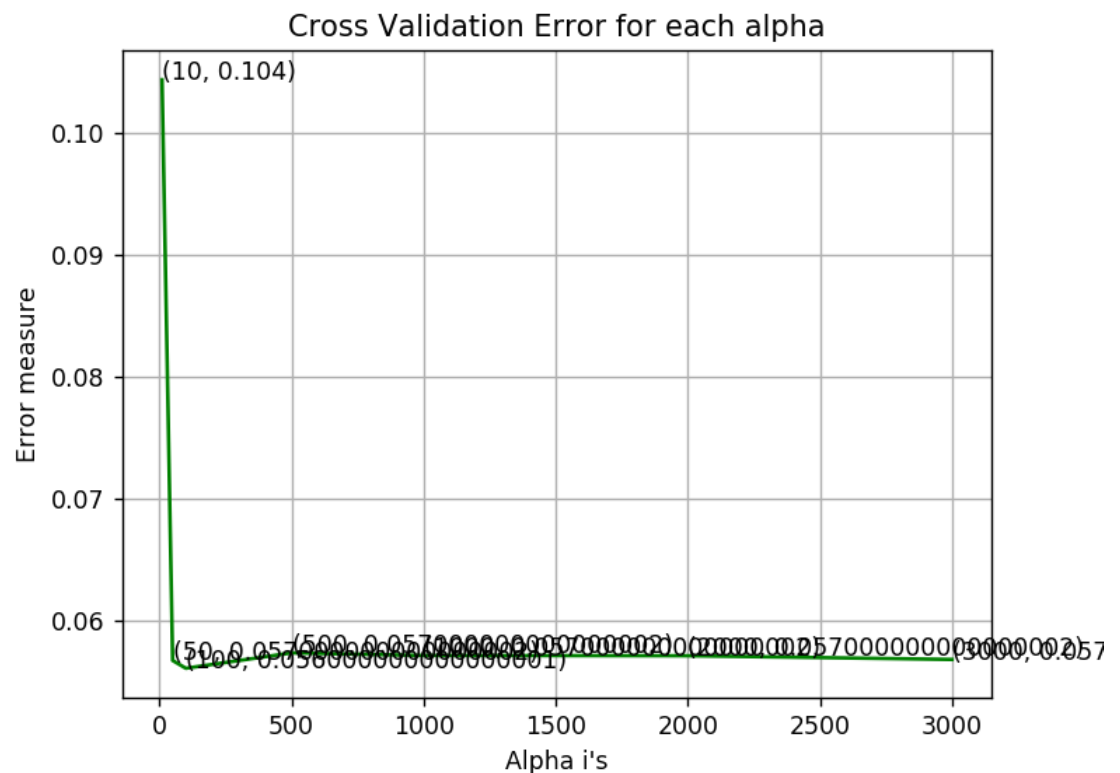
```
log_loss for c =  10 is 0.104344888454
log_loss for c =  50 is 0.0567190635611
log_loss for c =  100 is 0.056075038646
log_loss for c =  500 is 0.057336051683
log_loss for c =  1000 is 0.0571265109903
log_loss for c =  2000 is 0.057103406781
log_loss for c =  3000 is 0.0567993215778

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



Cross Validation Error for each alpha

```
For values of best alpha =  100 The train log loss is: 0.011788374
2574
For values of best alpha =  100 The cross validation log loss is:
0.056075038646
For values of best alpha =  100 The test log loss is: 0.0491647763
845
Number of misclassified points  0.873965041398
```
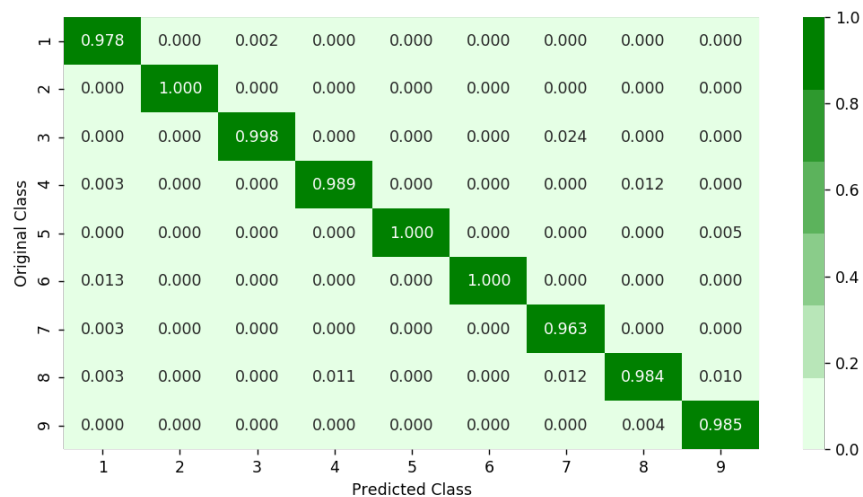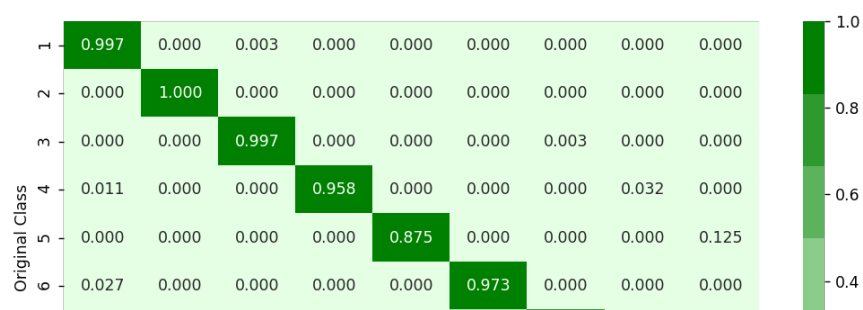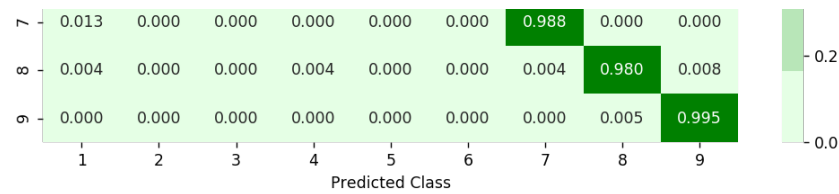
```
------------------------------------------------- Confusion matri
x -------------------------------------------------
<IPython.core.display.Javascript object>
```



```
------------------------------------------------- Precision matri
x -------------------------------------------------
<IPython.core.display.Javascript object>
```



```
Sum of columns in precision matrix [ 1.   1.   1.   1.   1.   1.   1.
1.   1.]
------------------------------------------------- Recall matrix
-------------------------------------------------
<IPython.core.display.Javascript object>
```

| | 0.013 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | **0.988** | 0.000 | 0.000 |
| 8 | 0.004 | 0.000 | 0.000 | 0.004 | 0.000 | 0.000 | 0.004 | **0.980** | 0.008 |
| 9 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.005 | **0.995** |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

Predicted Class

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

### 4.4.5 Xgboost Classifier with best hyperparameters

In [0]:
```python
x_cfl=XGBClassifier()

prams={
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
     'n_estimators':[100,200,500,1000,2000],
     'max_depth':[3,5,10],
     'colsample_bytree':[0.1,0.3,0.5,1],
     'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbo
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 rema
ining:   39.3s
[Parallel(n_jobs=-1)]: Done   23 out of  30 | elapsed:  1.3min rema
ining:   23.0s
[Parallel(n_jobs=-1)]: Done   27 out of  30 | elapsed:  1.4min rema
ining:    9.2s
[Parallel(n_jobs=-1)]: Done   30 out of  30 | elapsed:  2.3min fini
shed
```

Out[163]:
```
RandomizedSearchCV(cv=None, error_score='raise',
          estimator=XGBClassifier(base_score=0.5, colsample_byleve
l=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], 'subs
ample': [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]:
```python
print (random_cfl.best_params_)
```

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

In [0]:

```
x_cfl=XGBClassifier(n_estimators=200,subsample=0.5,learning_rate=0.
x_cfl.fit(X_train_asm,y_train_asm)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train_asm,y_train_asm)

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

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

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

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

In [0]: `result.head()`

Out[171]:

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

5 rows × 260 columns

In [0]: `result_asm.head()`

Out[174]:

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

5 rows × 54 columns

In [0]:
```
print(result.shape)
print(result_asm.shape)
```
```
(10868, 260)
(10868, 54)
```

In [0]:
```
result_x = pd.merge(result,result_asm.drop(['Class'], axis=1),on='I
result_y = result_x['Class']
```

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

Out[182]:

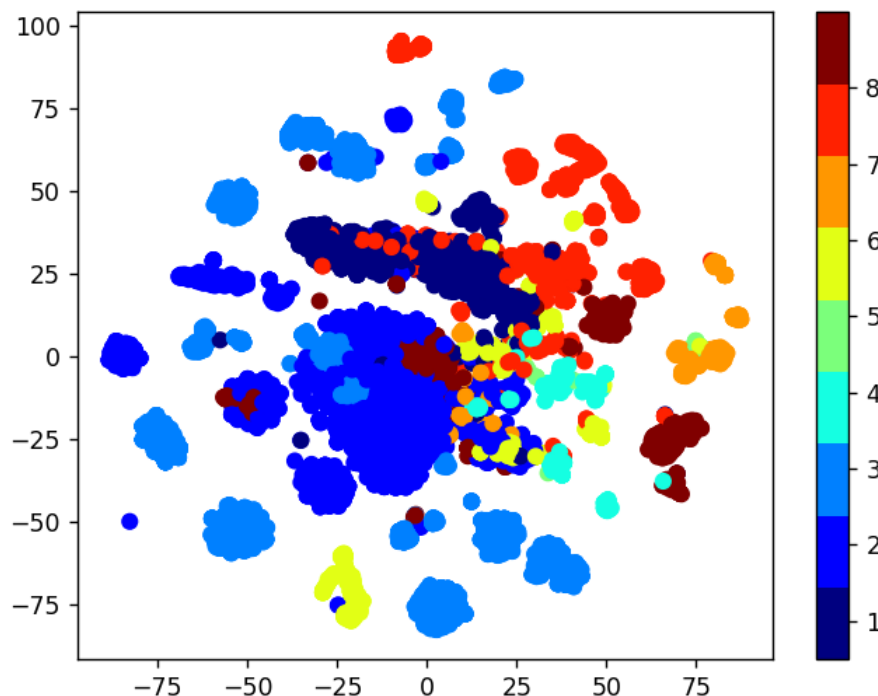| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.262806 | 0.005498 | 0.001567 | 0.002067 | 0.002048 | 0.001835 | 0.002058 | 0.002946 | 0.002638 |
| 1 | 0.017358 | 0.011737 | 0.004033 | 0.003876 | 0.005303 | 0.003873 | 0.004747 | 0.006984 | 0.008267 |
| 2 | 0.040827 | 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 | 0.000955 |
| 4 | 0.008629 | 0.001000 | 0.000168 | 0.000234 | 0.000342 | 0.000232 | 0.000148 | 0.000229 | 0.000370 |

5 rows × 307 columns

### 4.5.2. Multivariate Analysis on final fearures

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

<IPython.core.display.Javascript object>



### 4.5.3. Train and Test split

```
In [0]: X_train, X_test_merge, y_train, y_test_merge = train_test_split(res
X_train_merge, X_cv_merge, y_train_merge, y_cv_merge = train_test_s
```

### 4.5.4. Random Forest Classifier on final features

In [0]:

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

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


best_alpha = np.argmin(cv_log_error_array)

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


r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_
r_cfl.fit(X_train_merge,y_train_merge)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

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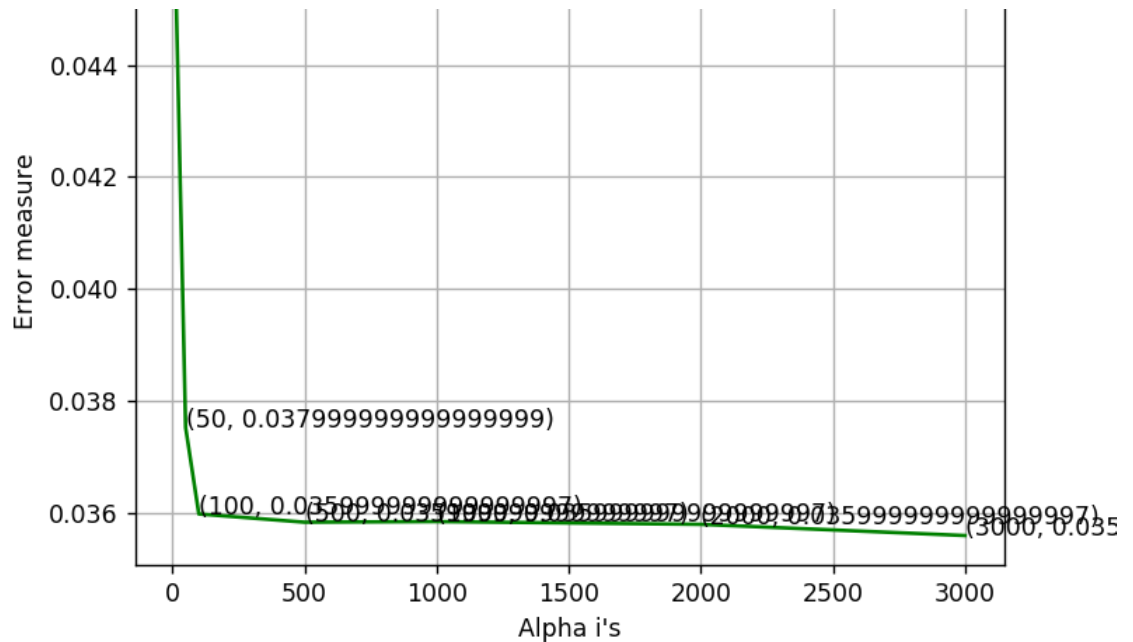
```
log_loss for c =   10 is 0.0461221662017
log_loss for c =   50 is 0.0375229563452
log_loss for c =   100 is 0.0359765822455
log_loss for c =   500 is 0.0358291883873
log_loss for c =   1000 is 0.0358403093496
log_loss for c =   2000 is 0.0357908022178
log_loss for c =   3000 is 0.0355909487962

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

Cross Validation Error for each alpha

(10, 0.045999999999999999)

0.046

```
For values of best alpha =  3000 The train log loss is: 0.01662676
14753
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: 0.040114130
3589
```

### 4.5.5. XgBoost Classifier on final features

In [0]:

```python
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=i)
    x_cfl.fit(X_train_merge,y_train_merge)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(X_train_merge, y_train_merge)
    predict_y = sig_clf.predict_proba(X_cv_merge)
    cv_log_error_array.append(log_loss(y_cv_merge, predict_y, label

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


best_alpha = np.argmin(cv_log_error_array)

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

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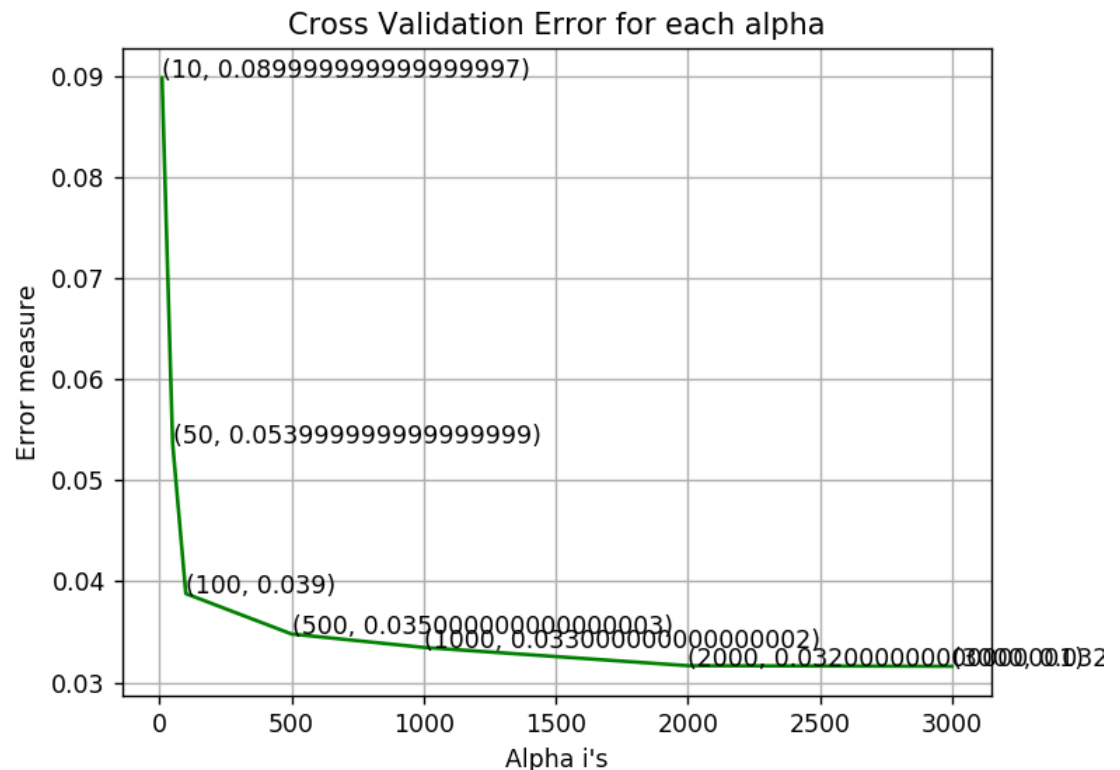
```
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

predict_y = sig_clf.predict_proba(X_train_merge)
print ('For values of best alpha = ', alpha[best_alpha], "The train
predict_y = sig_clf.predict_proba(X_cv_merge)
print('For values of best alpha = ', alpha[best_alpha], "The cross
predict_y = sig_clf.predict_proba(X_test_merge)
print('For values of best alpha = ', alpha[best_alpha], "The test l
```

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

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

Cross Validation Error for each alpha



```
For values of best alpha =  3000 The train log loss is: 0.01119188
09342
For values of best alpha =  3000 The cross validation log loss is:
0.0315972694477
For values of best alpha =  3000 The test log loss is: 0.032397851
5915
```

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

```
In [0]: x_cfl=XGBClassifier()

prams={
```

```python
    'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
     'n_estimators':[100,200,500,1000,2000],
     'max_depth':[3,5,10],
    'colsample_bytree':[0.1,0.3,0.5,1],
    'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbo
random_cfl.fit(X_train,meta_v_train.meta)
```

```
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 rema
ining:  2.6min
[Parallel(n_jobs=-1)]: Done   23 out of   30 | elapsed:  5.8min rema
ining:  1.8min
[Parallel(n_jobs=-1)]: Done   27 out of   30 | elapsed:  6.7min rema
ining:   44.5s
[Parallel(n_jobs=-1)]: Done   30 out of   30 | elapsed:  7.4min fini
shed
```

Out[187]:
```
RandomizedSearchCV(cv=None, error_score='raise',
          estimator=XGBClassifier(base_score=0.5, colsample_byleve
l=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], 'subs
ample': [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]:
```python
print (random_cfl.best_params_)
```

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

In [0]:
```python
# find more about XGBClassifier function here http://xgboost.readth
# -------------------------
# default paramters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_est
# objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=
# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_byle
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, mi

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

```python
x_cfl=XGBClassifier(n_estimators=1000,max_depth=10,learning_rate=0.
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

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

```
For values of best alpha =  3000 The train log loss is: 0.01219228
32297
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: 0.031704113
2442
```

# 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