

▼ 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 causes 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 protection, forcing the anti-malware groups/communities to build more robust softwares to detect an of protecting a computer system from a malware attack is to **identify whether a given piece of file/sc**

1.3 Source/Useful Links

Microsoft has been very active in building anti-malware products over the years and it runs its anti-malware system around the world. This generates tens of millions of daily data points to be analyzed as potential malware. And classifying such large amounts of data, we need to be able to group them into groups and identify them.

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

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

1.4. Real-world/Business objectives and constraints.

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

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
.text:00401000 56
.text:00401001 8D 44 24      08
.text:00401005 50
.text:00401006 8B F1
.text:00401008 E8 1C 1B      00 00
.text:0040100D C7 06 08      BB 42 00
                                assume es:nothing, ss:nothing, ds:_data
                                push    esi
                                lea     eax, [esp+8]
                                push    eax
                                mov    esi, ecx
                                call   ??0exception@std@@QAE
                                mov    dword ptr [esi], of
```

```

.text:00401013 8B C6          mov     eax, esi
.text:00401015 5E          pop     esi
.text:00401016 C2 04 00      retn    4
.text:00401016                 ; -----
.text:00401019 CC CC CC CC    CC CC CC CC      align 10h
.text:00401020 C7 01 08      BB 42 00      mov     dword ptr [ecx], of
.text:00401026 E9 26 1C      00 00      jmp     sub_402C51
.text:00401026                 ; -----
.text:0040102B CC 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], of
.text:00401039 E8 13 1C      00 00      call    sub_402C51
.text:0040103E F6 44 24      08 01      test    byte ptr [esp+8],
.text:00401043 74 09          jz     short loc_40104E
.text:00401045 56          push    esi
.text:00401046 E8 6C 1E      00 00      call    ??3@YAXPAX@Z      ; ope
.text:0040104B 83 C4 04      add     esp, 4
.text:0040104E
.text:0040104E 8B C6          loc_40104E:      ; CODE XREF: .tex
.text:0040104E                 mov     eax, esi
.text:00401050 5E          pop     esi
.text:00401051 C2 04 00      retn    4
.text:00401051                 ; -----

```

.bytes file

```

00401000 00 00 80 40 40 28 00 1C 02 42 00 C4 00 20 04 20
00401010 00 00 20 09 2A 02 00 00 00 00 8E 10 41 0A 21 01
00401020 40 00 02 01 00 90 21 00 32 40 00 1C 01 40 C8 18
00401030 40 82 02 63 20 00 00 09 10 01 02 21 00 82 00 04
00401040 82 20 08 83 00 08 00 00 00 02 00 60 80 10 80
00401050 18 00 00 20 A9 00 00 00 00 04 04 78 01 02 70 90
00401060 00 02 00 08 20 12 00 00 00 40 10 00 80 00 40 19
00401070 00 00 00 00 11 20 80 04 80 10 00 20 00 00 25 00
00401080 00 00 01 00 00 04 00 10 02 C1 80 80 00 20 20 00
00401090 08 A0 01 01 44 28 00 00 08 10 20 00 02 08 00 00
004010A0 00 40 00 00 00 34 40 40 00 04 00 08 80 08 00 08
004010B0 10 00 40 00 68 02 40 04 E1 00 28 14 00 08 20 0A
004010C0 06 01 02 00 40 00 00 00 00 00 20 00 02 00 04
004010D0 80 18 90 00 00 10 A0 00 45 09 00 10 04 40 44 82
004010E0 90 00 26 10 00 00 04 00 82 00 00 00 20 40 00 00
004010F0 B4 00 00 40 00 02 20 25 08 00 00 00 00 00 00 00

```

```
00401100 08 00 00 50 00 08 40 50 00 02 06 22 08 85 30 00
00401110 00 80 00 80 60 00 09 00 04 20 00 00 00 00 00 00
00401120 00 82 40 02 00 11 46 01 4A 01 8C 01 E6 00 86 10
00401130 4C 01 22 00 64 00 AE 01 EA 01 2A 11 E8 10 26 11
00401140 4E 11 8E 11 C2 00 6C 00 0C 11 60 01 CA 00 62 10
00401150 6C 01 A0 11 CE 10 2C 11 4E 10 8C 00 CE 01 AE 01
00401160 6C 10 6C 11 A2 01 AE 00 46 11 EE 10 22 00 A8 00
00401170 EC 01 08 11 A2 01 AE 10 6C 00 6E 00 AC 11 8C 00
00401180 EC 01 2A 10 2A 01 AE 00 40 00 C8 10 48 01 4E 11
00401190 0E 00 EC 11 24 10 4A 10 04 01 C8 11 E6 01 C2 00
```

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

2.2.1. Type of Machine Learning Problem

There are nine different classes of malware that we need to classify a given a data point => M

2.2.2. Performance Metric

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

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learing Objectives and Constraints

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

Constraints:

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

2.3. Train and Test Dataset

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

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

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

source = 'train'
```

```
destination = 'byteFiles'

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

# if we have folder called 'train' (train folder contains both .asm files and .bytes files) w
# for every file that we have in our 'asmFiles' directory we check if it is ending with .byte
# 'byteFiles' folder

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

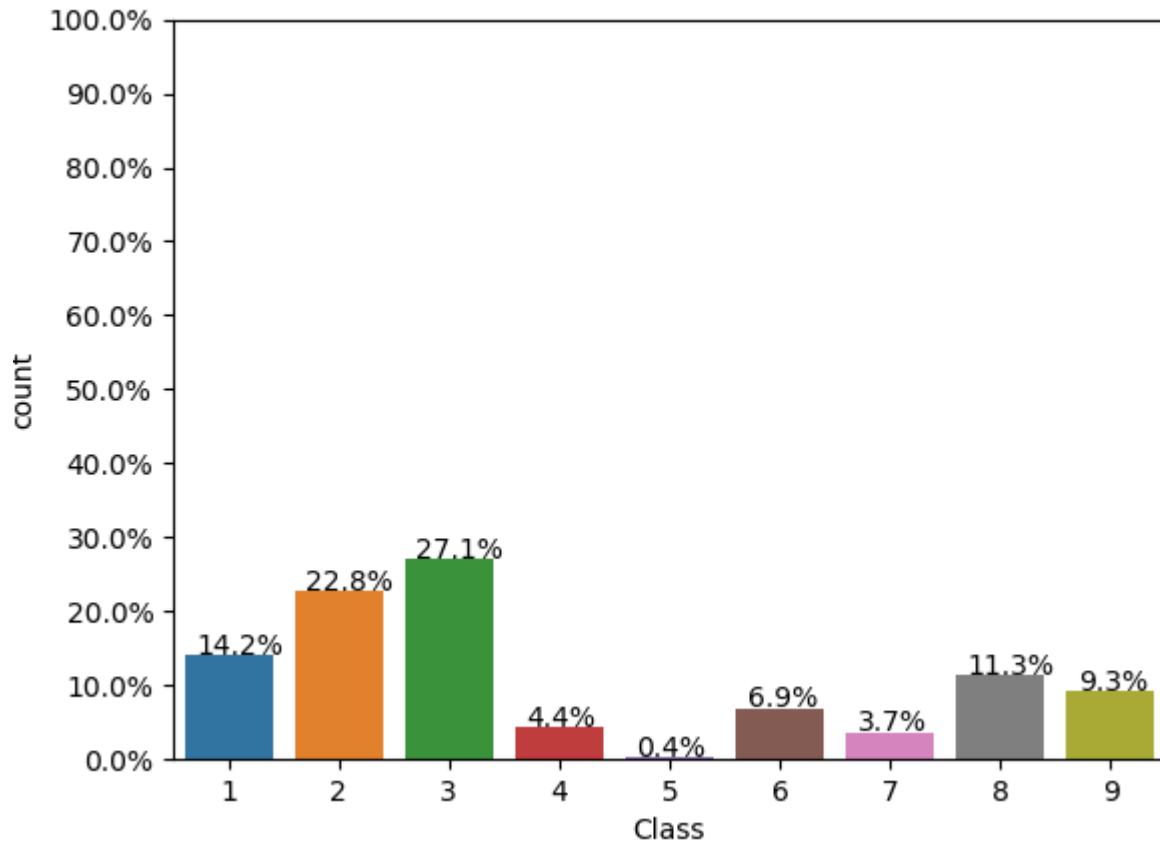
3.1. Distribution of malware classes in whole data set

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

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

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





3.2. Feature extraction

3.2.1 File size of byte files as a feature

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

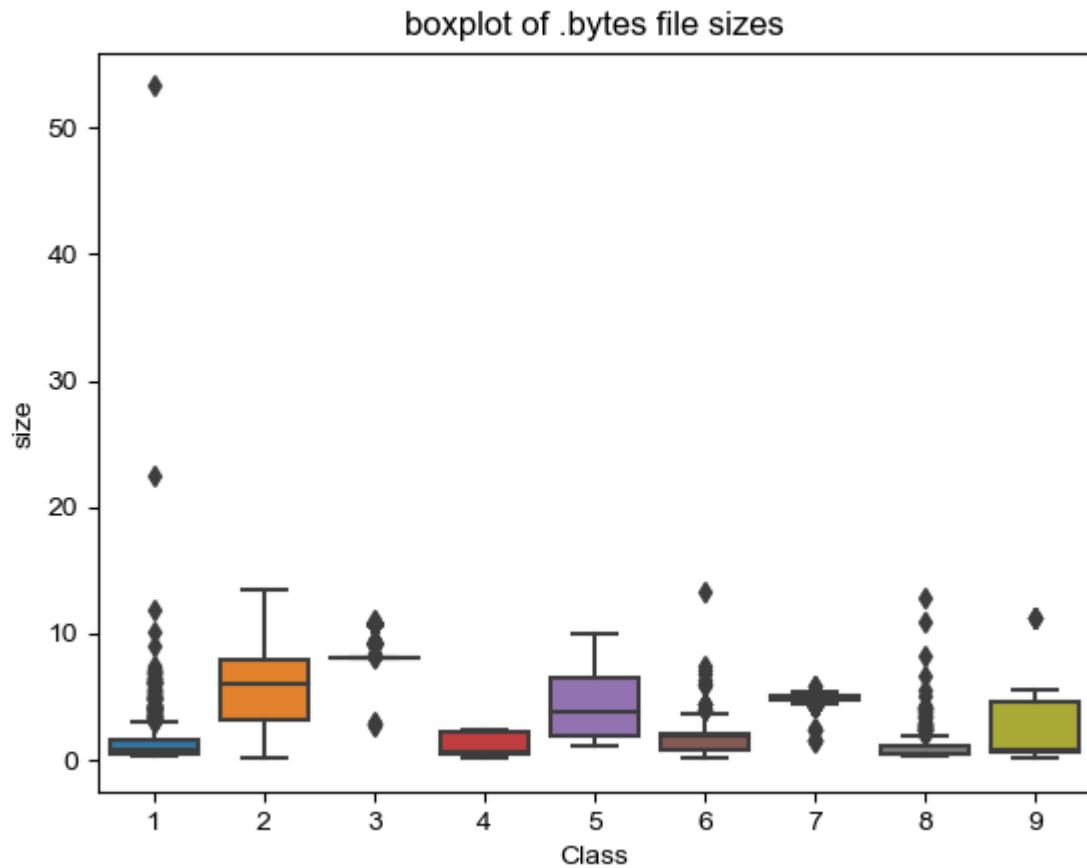
```
i=filenames.index(file)
class_bytes.append(class_y[i])
# converting into Mb's
sizebytes.append(statinfo.st_size/(1024.0*1024.0))
fnames.append(file)
data_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes})
print (data_size_byte.head())
```



	ID	size	Class
0	01azqd4InC7m9JpocGv5	5.012695	9
1	01IsoiSMh5gxyDYT14CB	6.556152	2
2	01jsnpXSAlgw6aPeDxrU	4.602051	9
3	01kcPWA9K2B0xQeS5Rju	0.679688	1
4	01SuzwMJEIXsK7A8dQb1	0.438965	8

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

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



3.2.3 feature extraction from byte files

```

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

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

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

#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,1,2,3,4,5,6,7,8,9,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,17,18,19
for file in files:
    filenames2.append(file)
    byte_feature_file.write(file+",")
    if(file.endswith("txt")):
        with open('byteFiles/'+file,"r") as byte_file:
            for lines in byte_file:
                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_file.close()
for i in feature_matrix[k]:
    byte_feature_file.write(str(i)+",")
byte_feature_file.write("\n")

k += 1

byte_feature_file.close()

```

mk

```

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

```

	Unnamed: 0	ID	0	1	2	3	4	5	6	7
0	0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201
1	1	01lsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589
2	2	01jsnpXSAlg6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342
3	3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523
4	4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249

5 rows × 263 columns

```

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

```

```

      Unnamed: 0          ID      0      1      2      3      4      5  \
0          0  01azqd4InC7m9JpocGv5  601905  3905  2816  3832  3345  3242
1          1  01IsoiSMh5gxyDYT14CB  39755  8337  7249  7186  8663  6844
2          2  01jsnpXSAlgw6aPeDxrU  93506  9542  2568  2438  8925  9330
3          3  01kcPWA9K2B0xQeS5Rju  21091  1213  726   817  1257  625
4          4  01SuzwMJEIXsK7A8dQbl  19764  710   302  433   559  410

      6      7  ...      f9      fa      fb      fc      fd      fe      ff  ??  \
0  3650  3201  ...  3101  3211  3097  2758  3099  2759  5753  1824
1  8420  7589  ...  439   281   302   7639  518   17001  54902  8588
2  9007  2342  ...  2242  2885  2863  2471  2786  2680  49144  468
3  550   523   ...  485   462   516   1133  471   761   7998  13940
4  262   249   ...  350   209   239   653   221   242   2199  9008

      size  Class
0  5.012695      9
1  6.556152      2
2  4.602051      9
3  0.679688      1
4  0.438965      8

[5 rows x 261 columns]

```

```
result=byte_features
```

```

# 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)!=str('Class')):
            max_value = df[feature_name].max()
            min_value = df[feature_name].min()
            result1[feature_name] = (df[feature_name] - min_value) / (max_value - min_value)
    return result1
result = normalize(result)

```

```
from sklearn.externals import joblib
```

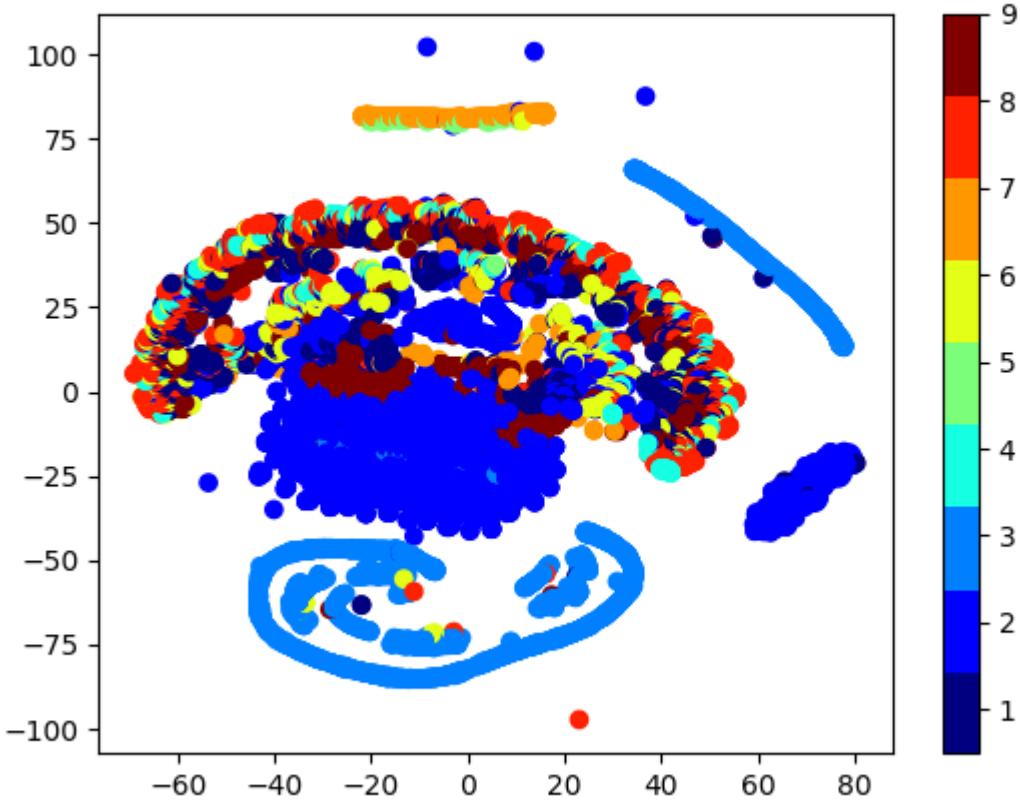
```
joblib.dump(result, 'result.pkl')
```

```
👤 ['result.pkl']
```

3.2.4 Multivariate Analysis

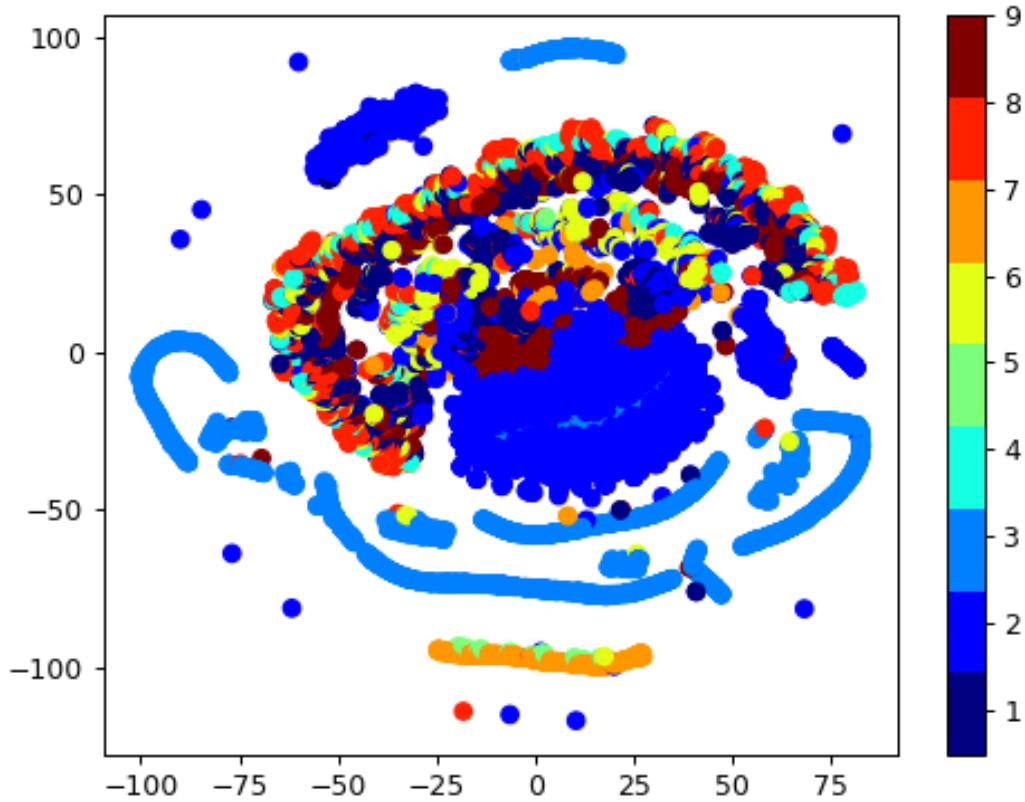
```
#multivariate analysis on byte files
```

```
#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()
```



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





▼ Train Test split

```
data_y = result['Class']
# split the data into test and train by maintaining same distribution of output variable 'y_t'
X_train, X_test, y_train, y_test = train_test_split(result.drop(['ID','Class'], axis=1), data_y, test_size=0.2)
# split the train data into train and cross validation by maintaining same distribution of output variable 'y_t'
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)

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

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

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

```
cv_class_distribution = y_cv.value_counts().sortlevel()

my_colors = ['#b23850', '#3b8beb', '#e7e3d4', '#c4dbf6', '#8590aa', '#0d19a3', '#15db95', '#0
train_class_distribution.plot(kind='bar', color=my_colors)
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in train data')
plt.grid()
plt.show()

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

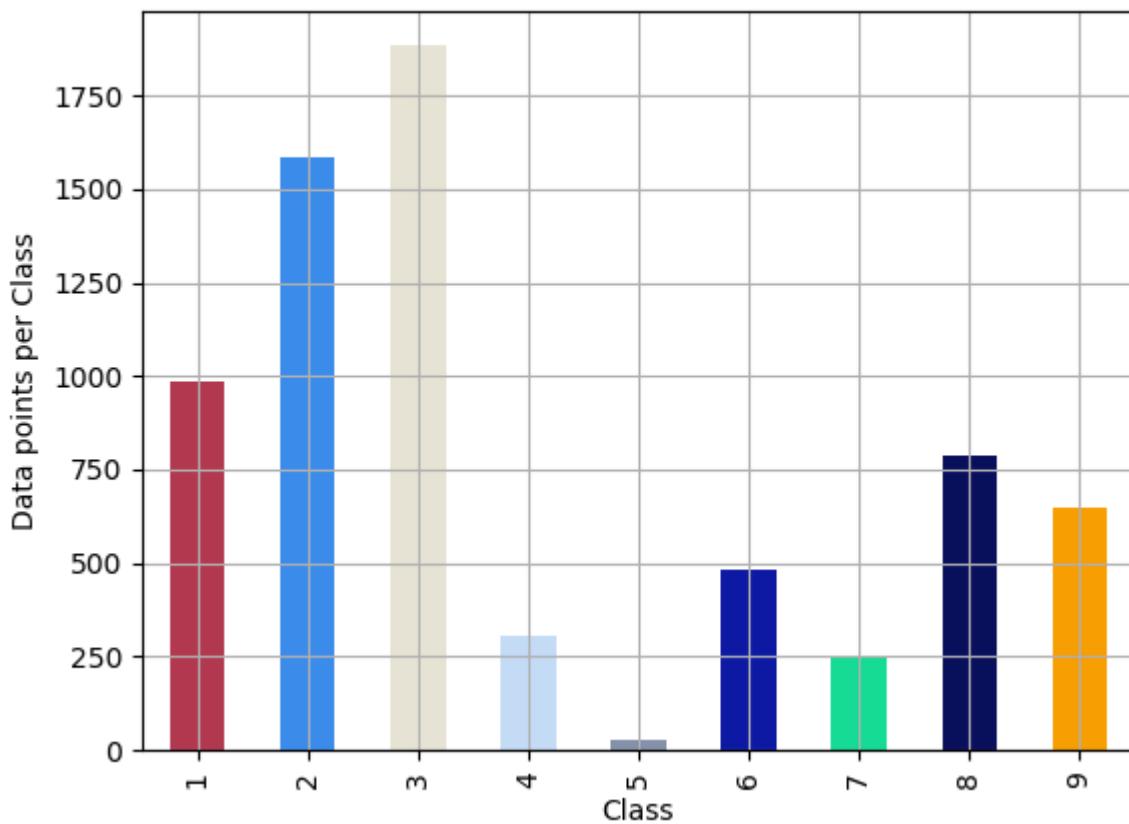
print('*'*80)
my_colors = ['#b23850', '#3b8beb', '#e7e3d4', '#c4dbf6', '#8590aa', '#0d19a3', '#15db95', '#0
test_class_distribution.plot(kind='bar', color=my_colors)
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

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

print('*'*80)
my_colors = ['#b23850', '#3b8beb', '#e7e3d4', '#c4dbf6', '#8590aa', '#0d19a3', '#15db95', '#0
cv_class_distribution.plot(kind='bar', color=my_colors)
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in cross validation data')
plt.grid()
plt.show()

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



Distribution of y_i 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 %)

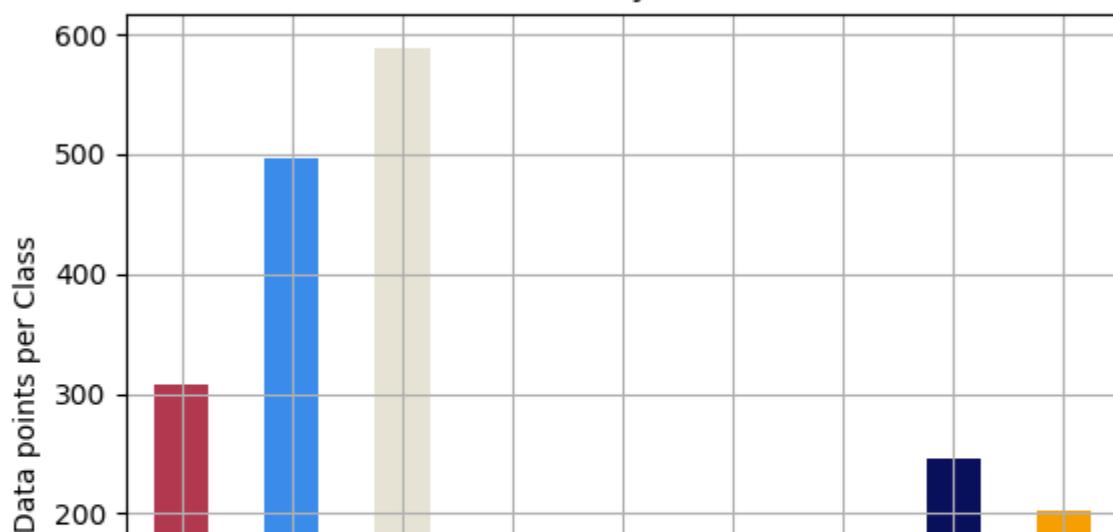
Number of data points in class 9 : 648 (9.317 %)

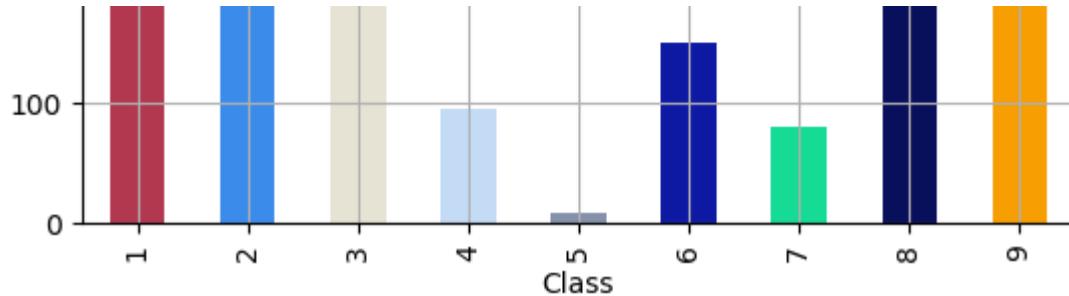
Number of data points in class 6 : 481 (6.916 %)

Number of data points in class 4 : 304 (4.371 %)

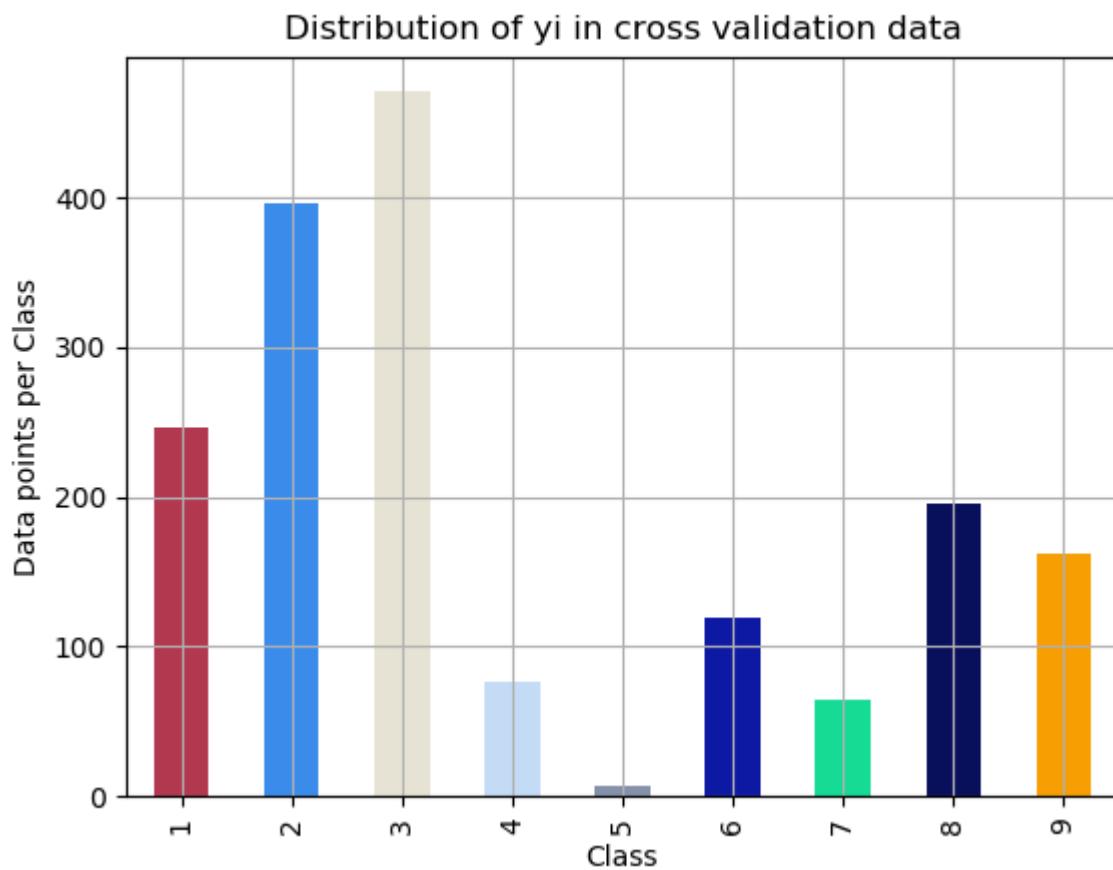
Number of data points in class 7 : 254 (3.652 %)

Number of data points in class 5 : 27 (0.388 %)

Distribution of y_i 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 %)



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 %)

```

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

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

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

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

    B =(C/C.sum(axis=0))
    #divid each element of the confusion matrix with the sum of elements in that row

    # C = [[1, 2],
    #       [3, 4]]
    # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to rows in two dia
    # C.sum(axix =0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                      [3/4, 4/6]]

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

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

```

```

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

```

4. Machine Learning Models

4.1. Machine Learning Models on bytes files

4.1.1. Random Model

```

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

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

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

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

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

```



Log loss on Cross Validation Data using Random Model 2.4987116946656167

Log loss on Test Data using Random Model 2.4553327958473936

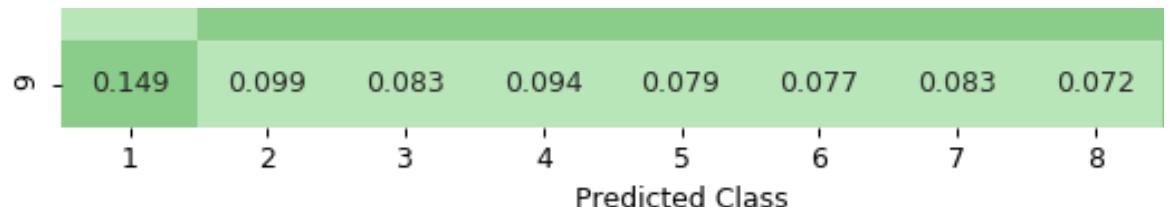
Number of misclassified points 88.45446182152715

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

Original Class	Predicted Class							
	1	2	3	4	5	6	7	8
1	35.000	40.000	24.000	30.000	37.000	32.000	44.000	40.000
2	47.000	55.000	61.000	61.000	59.000	56.000	46.000	56.000
3	61.000	64.000	69.000	64.000	60.000	74.000	69.000	74.000
4	7.000	14.000	11.000	14.000	11.000	11.000	10.000	6.000
5	1.000	2.000	1.000	1.000	0.000	0.000	0.000	3.000
6	16.000	15.000	23.000	18.000	10.000	18.000	13.000	21.000
7	7.000	12.000	6.000	7.000	9.000	10.000	8.000	17.000
8	20.000	34.000	25.000	26.000	24.000	28.000	30.000	28.000
9	34.000	26.000	20.000	23.000	18.000	19.000	20.000	19.000

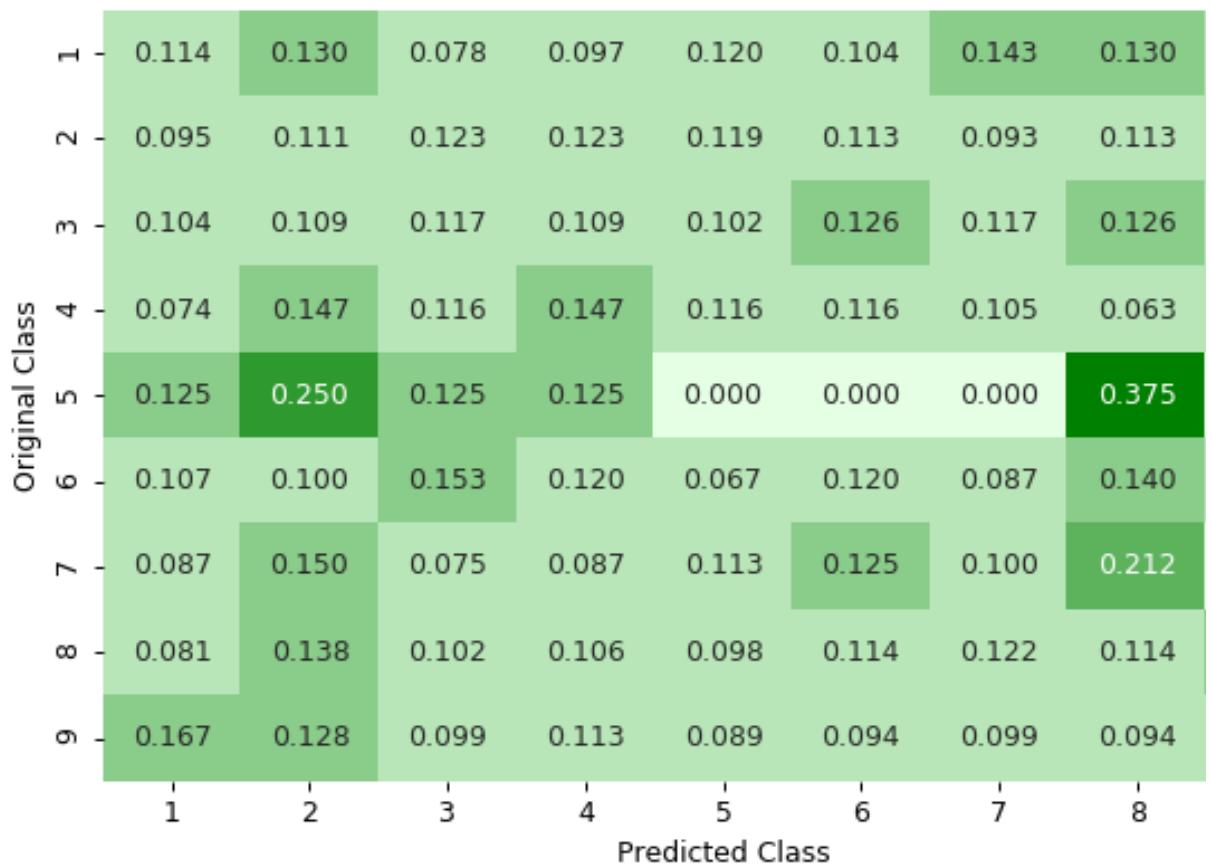
----- Precision matrix -----

Original Class	Predicted Class							
	1	2	3	4	5	6	7	8
1	0.154	0.153	0.100	0.123	0.162	0.129	0.183	0.152
2	0.206	0.210	0.254	0.250	0.259	0.226	0.192	0.212
3	0.268	0.244	0.287	0.262	0.263	0.298	0.287	0.280
4	0.031	0.053	0.046	0.057	0.048	0.044	0.042	0.023
5	0.004	0.008	0.004	0.004	0.000	0.000	0.000	0.011
6	0.070	0.057	0.096	0.074	0.044	0.073	0.054	0.080
7	0.031	0.046	0.025	0.029	0.039	0.040	0.033	0.064
8	0.088	0.130	0.104	0.107	0.105	0.113	0.125	0.106



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

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



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

4.1.2. K Nearest Neighbour Classification

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

# methods of
# fit(X, y) : Fit the model using X as training data and y as target values
# predict(X):Predict the class labels for the provided data
```

```

# predict_proba(X):Return probability estimates for the test data X.
#-----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/k-neare
#-----


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


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

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

best_alpha = np.argmin(cv_log_error_array)

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

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

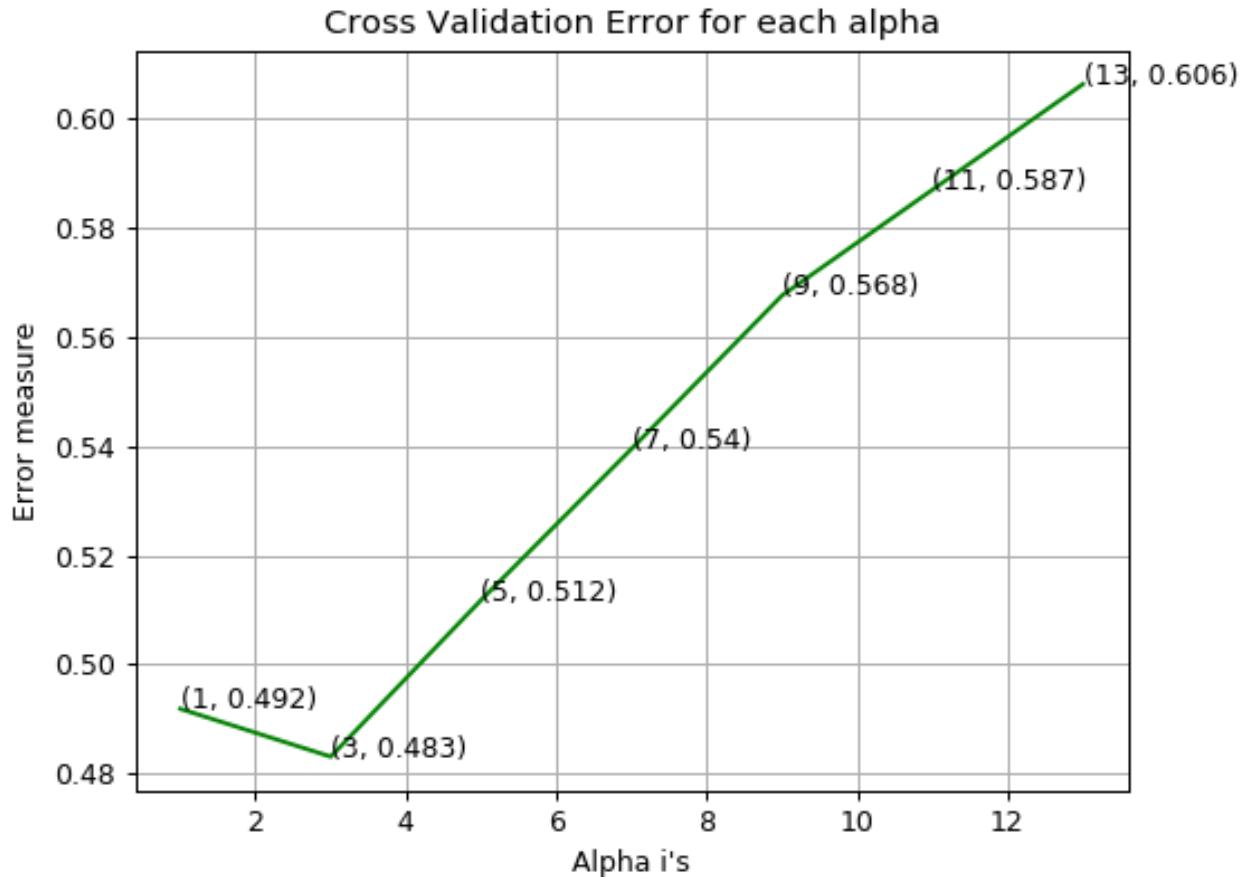
predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_cv, predict_y))

```

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



```
log_loss for k = 1 is 0.49188045368463196
log_loss for k = 3 is 0.483116902642161
log_loss for k = 5 is 0.5118350087441232
log_loss for k = 7 is 0.5395490778512431
log_loss for k = 9 is 0.5676371813660702
log_loss for k = 11 is 0.5870170308367498
log_loss for k = 13 is 0.606375118318671
```



For values of best alpha = 3 The train log loss is: 0.29320594139515405
 For values of best alpha = 3 The cross validation log loss is: 0.483116902642161
 For values of best alpha = 3 The test log loss is: 0.4851954874286108
 Number of misclassified points 12.97148114075437

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

Original Class	1	2	3	4	5	6	7	8
1	271.000	0.000	1.000	4.000	0.000	8.000	1.000	15.000
2	21.000	417.000	3.000	3.000	0.000	16.000	1.000	3.000
3	0.000	0.000	588.000	0.000	0.000	0.000	0.000	0.000
4	3.000	1.000	0.000	86.000	0.000	3.000	0.000	2.000
5	1.000	0.000	0.000	3.000	0.000	0.000	4.000	0.000

Original Class	6	30.000	1.000	3.000	11.000	0.000	98.000	0.000	2.000
7	-	4.000	1.000	1.000	0.000	0.000	5.000	66.000	1.000
8	-	14.000	0.000	0.000	5.000	1.000	5.000	1.000	211.000
9	-	26.000	4.000	1.000	1.000	0.000	5.000	0.000	11.000
	1	2	3	4	5	6	7	8	
	Predicted Class								

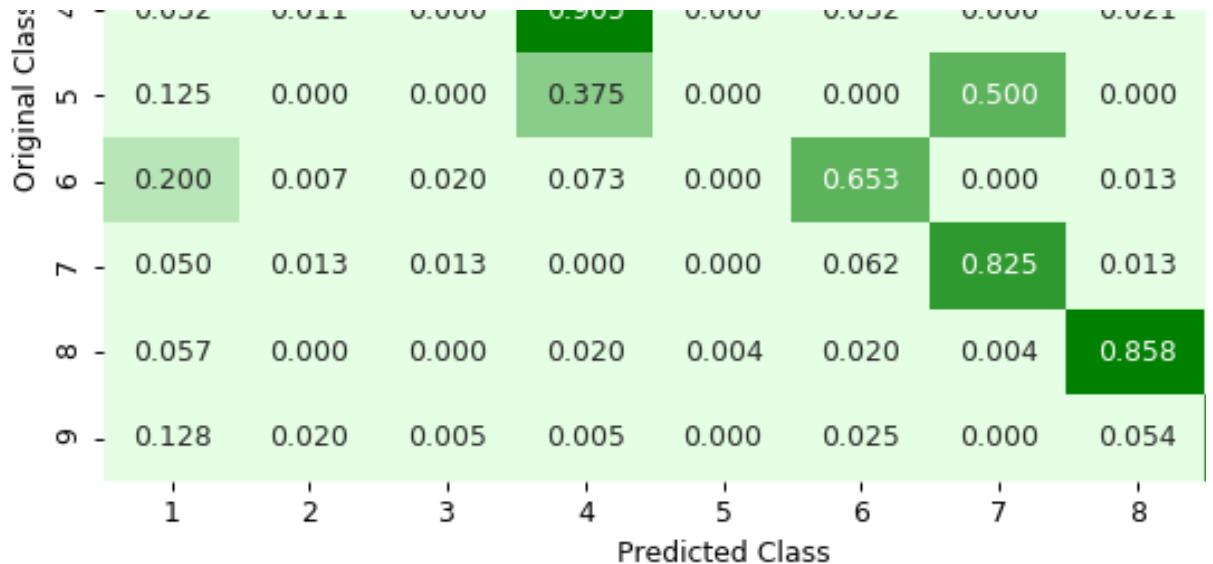
----- Precision matrix -----

Original Class	1	0.732	0.000	0.002	0.035	0.000	0.057	0.014	0.061
2	-	0.057	0.983	0.005	0.027	0.000	0.114	0.014	0.012
3	-	0.000	0.000	0.985	0.000	0.000	0.000	0.000	0.000
4	-	0.008	0.002	0.000	0.761	0.000	0.021	0.000	0.008
5	-	0.003	0.000	0.000	0.027	0.000	0.000	0.055	0.000
6	-	0.081	0.002	0.005	0.097	0.000	0.700	0.000	0.008
7	-	0.011	0.002	0.002	0.000	0.000	0.036	0.904	0.004
8	-	0.038	0.000	0.000	0.044	1.000	0.036	0.014	0.861
9	-	0.070	0.009	0.002	0.009	0.000	0.036	0.000	0.045
	1	2	3	4	5	6	7	8	
	Predicted Class								

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

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

Original Class	1	0.880	0.000	0.003	0.013	0.000	0.026	0.003	0.049
2	-	0.042	0.841	0.006	0.006	0.000	0.032	0.002	0.006
3	-	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000
4	-	0.032	0.011	0.000	0.005	0.000	0.032	0.000	0.021



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

4.1.3. Logistic Regression

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

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

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

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

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

```
best_alpha = np.argmin(cv_log_error_array)

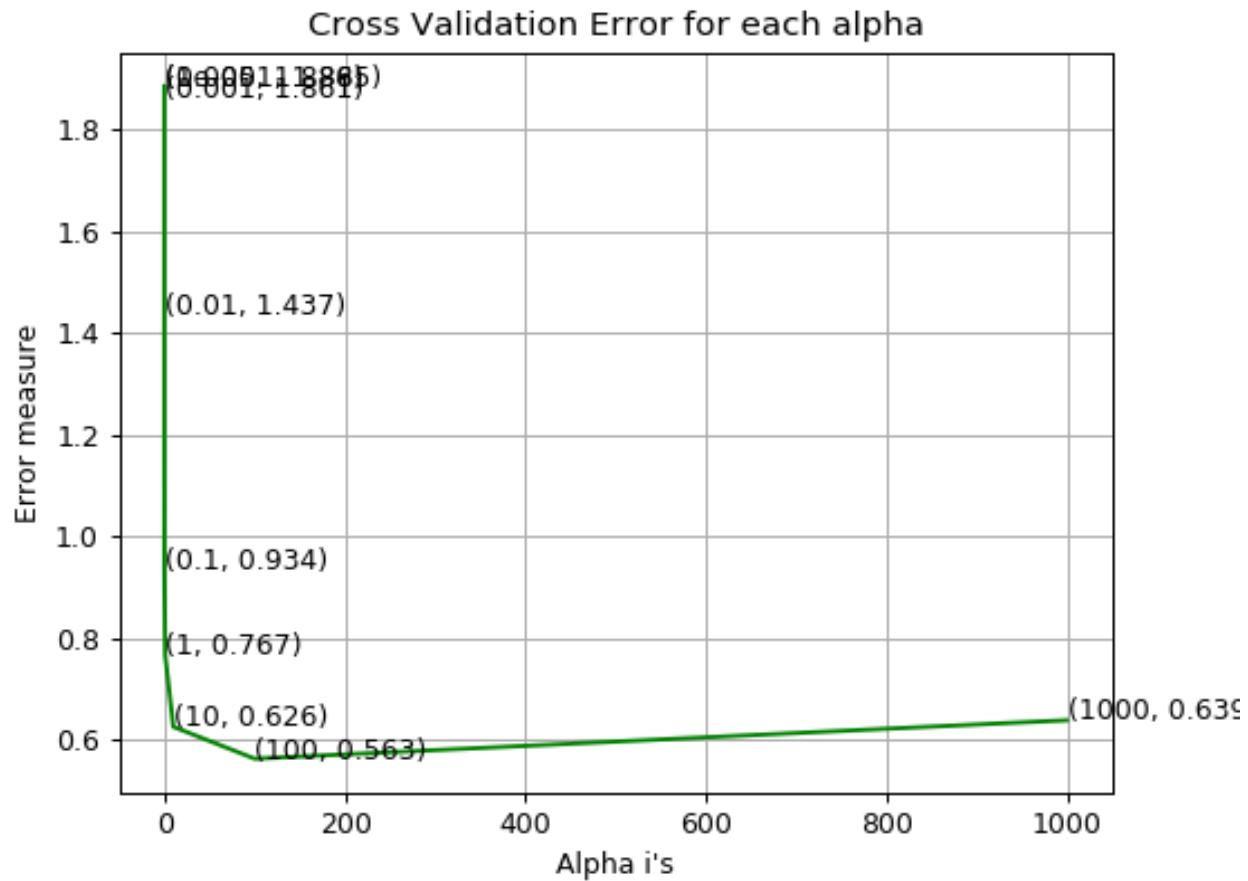
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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



```
log_loss for c = 1e-05 is 1.8861109311791922
log_loss for c = 0.0001 is 1.8845880471197165
log_loss for c = 0.001 is 1.8614285515198798
log_loss for c = 0.01 is 1.437434379095915
log_loss for c = 0.1 is 0.9337959204695321
log_loss for c = 1 is 0.7667190017910965
log_loss for c = 10 is 0.6257185173536978
log_loss for c = 100 is 0.56294262675526
log_loss for c = 1000 is 0.6385231825855628
```



```
log loss for train data 0.4871437301878761
log loss for cv data 0.56294262675526
log loss for test data 0.5294168099375685
Number of misclassified points 12.603495860165593
```

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

		Confusion matrix									
		Actual Class		Predicted Class 1		Predicted Class 2		Predicted Class 3		Predicted Class 4	
Predicted Class 1	Predicted Class 2	Predicted Class 3	Predicted Class 4	1	2	3	4	5	6	7	8
				247.000	2.000	0.000	3.000	0.000	3.000	1.000	46.000
1	30.000	433.000	5.000	4.000	0.000	7.000	0.000	9.000			
2	0.000	0.000	588.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
3	1.000	0.000	0.000	92.000	0.000	0.000	0.000	0.000	0.000	0.000	2.000
4	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Original Class	1	2	3	4	5	6	7	8
1	6.000	1.000	6.000	19.000	0.000	82.000	0.000	26.000
2	3.000	1.000	0.000	0.000	0.000	2.000	72.000	0.000
3	10.000	0.000	1.000	3.000	0.000	6.000	3.000	221.000
4	13.000	1.000	1.000	1.000	0.000	5.000	0.000	17.000

Predicted Class

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

Original Class	1	2	3	4	5	6	7	8
1	0.797	0.005	0.000	0.024		0.028	0.012	0.143
2	0.097	0.989	0.008	0.033		0.065	0.000	0.028
3	0.000	0.000	0.978	0.000		0.000	0.000	0.000
4	0.003	0.000	0.000	0.748		0.000	0.000	0.006
5	0.000	0.000	0.000	0.008		0.019	0.062	0.000
6	0.019	0.002	0.010	0.154		0.766	0.000	0.081
7	0.010	0.002	0.000	0.000		0.019	0.889	0.000
8	0.032	0.000	0.002	0.024		0.056	0.037	0.688
9	0.042	0.002	0.002	0.008		0.047	0.000	0.053

Predicted Class

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

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

Original Class	1	2	3	4	5	6	7	8
1	0.802	0.006	0.000	0.010	0.000	0.010	0.003	0.149
2	0.060	0.873	0.010	0.008	0.000	0.014	0.000	0.018
3	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000



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

4.1.4. Random Forest Classifier

```

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

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

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

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-classifier
# -----
```

```

alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
train_log_error_array=[]
from sklearn.ensemble import RandomForestClassifier
for i in alpha:
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    ...
```

```
r_cfl.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train, y_train)
predict_y = sig_clf.predict_proba(X_cv)
cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

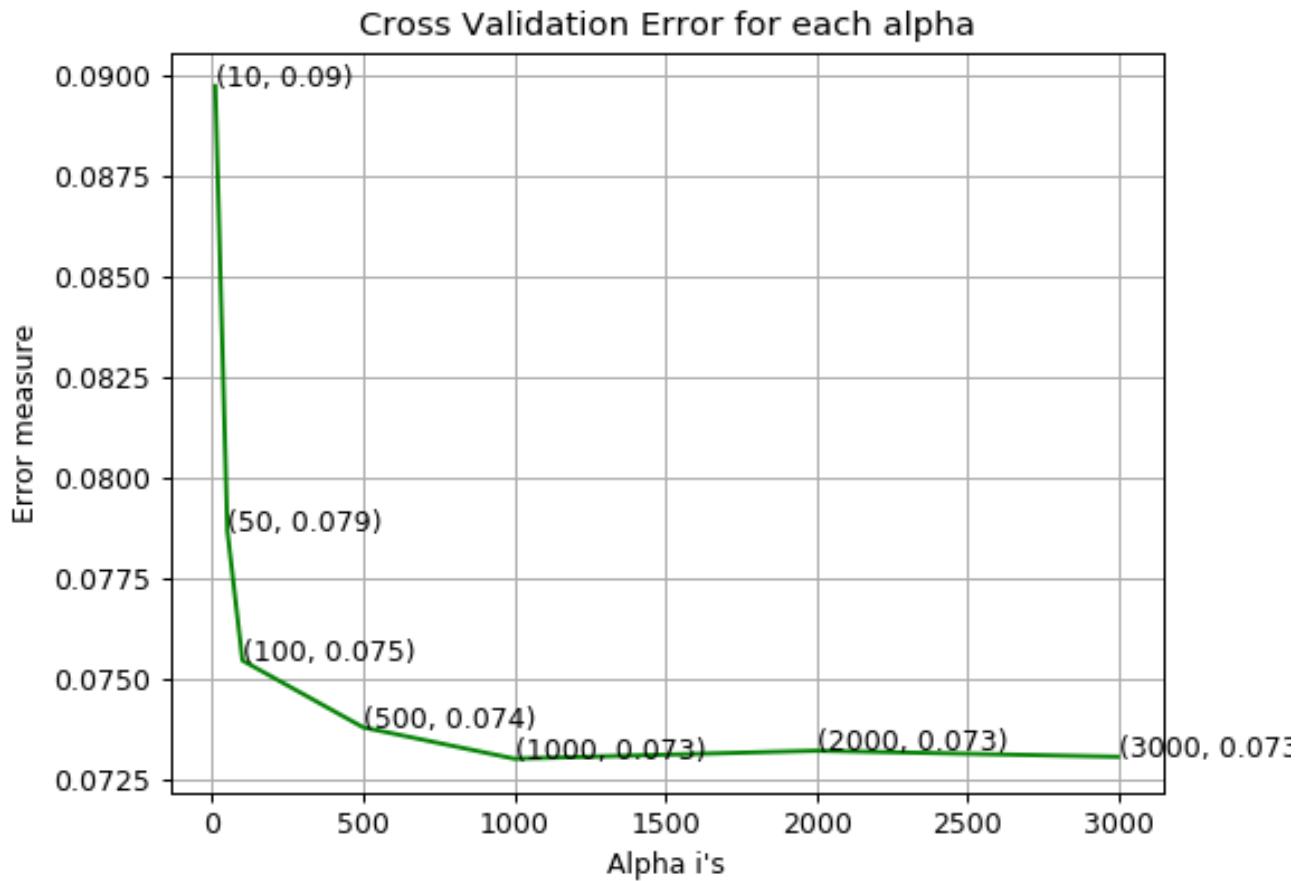
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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



```
log_loss for c = 10 is 0.0897527296356877
log_loss for c = 50 is 0.07869374681057625
log_loss for c = 100 is 0.07546292664044586
log_loss for c = 500 is 0.07379883728342362
log_loss for c = 1000 is 0.07302077078516724
log_loss for c = 2000 is 0.07321813574020479
log_loss for c = 3000 is 0.073068132864414
```



For values of best alpha = 1000 The train log loss is: 0.02941015229514485

For values of best alpha = 1000 The cross validation log loss is: 0.07302077078516724

For values of best alpha = 1000 The test log loss is: 0.06623869322452512

Number of misclassified points 1.2419503219871204

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

Original Class	1	2	3	4	5	6	7	8
1	302.000	3.000	0.000	0.000	0.000	2.000	0.000	0.000
2	0.000	495.000	0.000	0.000	0.000	1.000	0.000	0.000
3	0.000	0.000	587.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	92.000	0.000	0.000	1.000	2.000
5	0.000	0.000	0.000	0.000	8.000	0.000	0.000	0.000

Original Class	6	0.000	0.000	0.000	0.000	0.000	148.000	0.000	2.000
7	-	0.000	0.000	0.000	0.000	0.000	0.000	79.000	1.000
8	-	3.000	1.000	0.000	0.000	0.000	2.000	0.000	238.000
9	-	0.000	0.000	1.000	0.000	0.000	1.000	0.000	3.000
	1	2	3	4	5	6	7	8	
	Predicted Class								

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

Original Class	1	0.990	0.006	0.000	0.000	0.000	0.013	0.000	0.000
2	-	0.000	0.992	0.000	0.000	0.000	0.006	0.000	0.000
3	-	0.000	0.000	0.998	0.000	0.000	0.000	0.000	0.000
4	-	0.000	0.000	0.000	1.000	0.000	0.000	0.013	0.008
5	-	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000
6	-	0.000	0.000	0.000	0.000	0.000	0.961	0.000	0.008
7	-	0.000	0.000	0.000	0.000	0.000	0.000	0.988	0.004
8	-	0.010	0.002	0.000	0.000	0.000	0.013	0.000	0.967
9	-	0.000	0.000	0.002	0.000	0.000	0.006	0.000	0.012
	1	2	3	4	5	6	7	8	
	Predicted Class								

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

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

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

4.1.5. XgBoost Classification

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

```
# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/
```

```

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

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None,
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regres
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-a
# -----
```

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

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

best_alpha = np.argmin(cv_log_error_array)

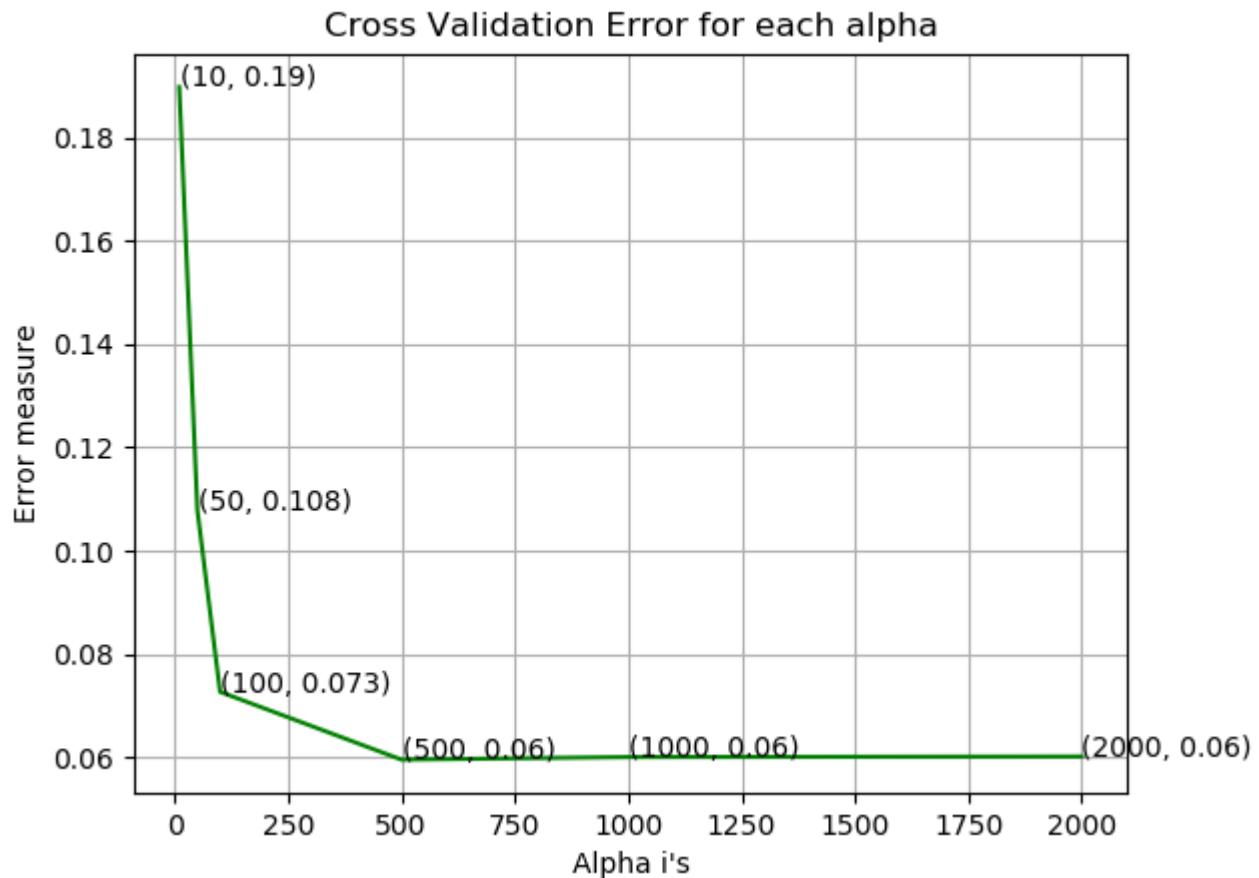
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
 ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:" 1

```
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_te  
predict_y = sig_clf.predict_proba(X_test)  
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_te
```

log_loss for c = 10 is 0.18972194520294353
log_loss for c = 50 is 0.10788109266220497
log_loss for c = 100 is 0.0727450829972164
log_loss for c = 500 is 0.059635927131908094
log_loss for c = 1000 is 0.06014538270053144
log_loss for c = 2000 is 0.060249389062255305



For values of best alpha = 500 The train log loss is: 0.02468009568654092
For values of best alpha = 500 The cross validation log loss is: 0.059635927131908094
For values of best alpha = 500 The test log loss is: 0.07847700799402009

```
plot_confusion_matrix(y_test, sig_clf.predict(X_test))
```



Number of misclassified points 1.6559337626494939

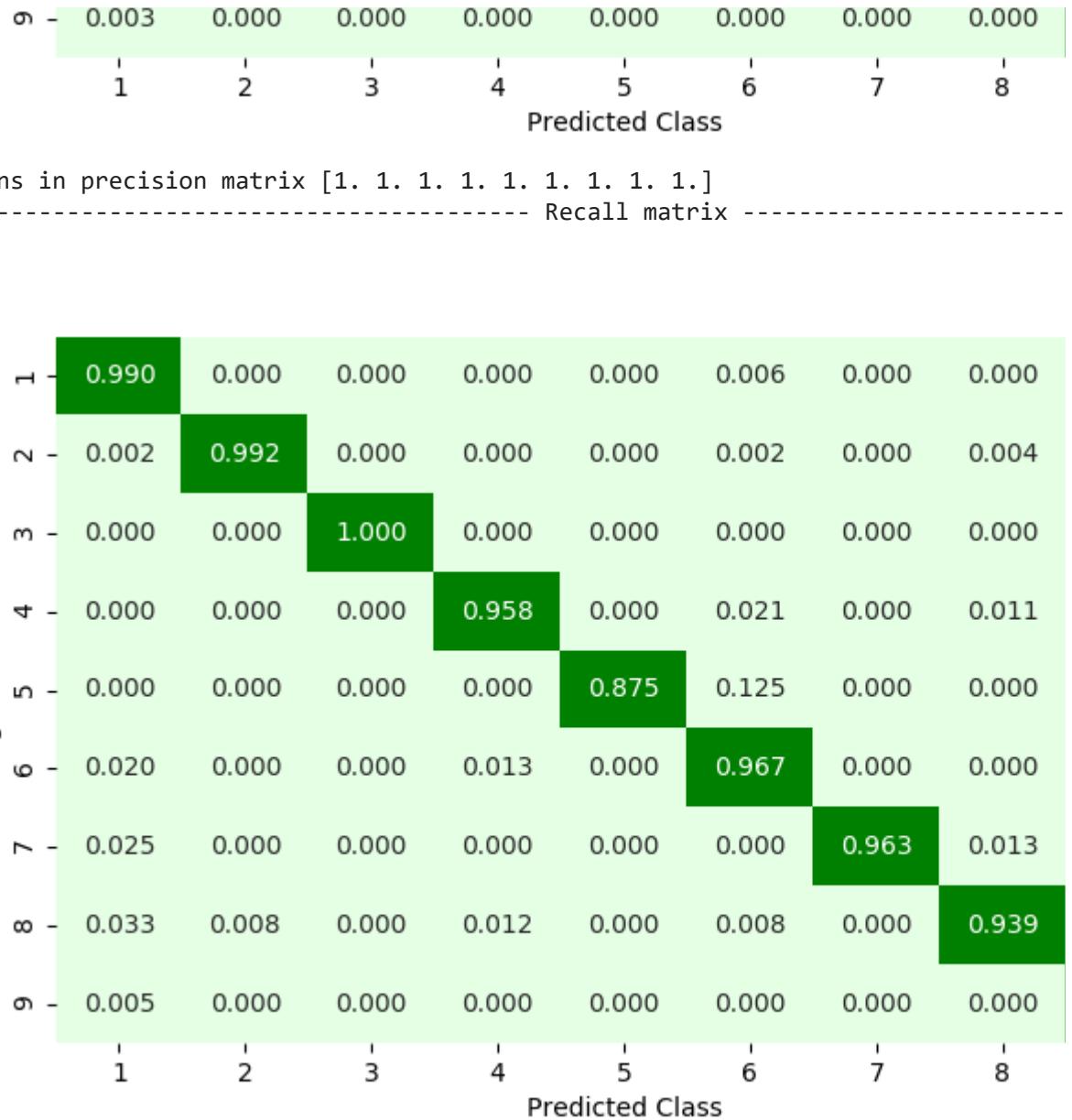
----- Confusion matrix -----

Original Class	1	2	3	4	5	6	7	8
	1	2	3	4	5	6	7	8
1	305.000	0.000	0.000	0.000	0.000	2.000	0.000	0.000
2	1.000	492.000	0.000	0.000	0.000	1.000	0.000	2.000
3	0.000	0.000	588.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	91.000	0.000	2.000	0.000	1.000
5	0.000	0.000	0.000	0.000	7.000	1.000	0.000	0.000
6	3.000	0.000	0.000	2.000	0.000	145.000	0.000	0.000
7	2.000	0.000	0.000	0.000	0.000	0.000	77.000	1.000
8	8.000	2.000	0.000	3.000	0.000	2.000	0.000	231.000
9	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Predicted Class

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

Original Class	1	2	3	4	5	6	7	8
	1	2	3	4	5	6	7	8
1	0.953	0.000	0.000	0.000	0.000	0.013	0.000	0.000
2	0.003	0.996	0.000	0.000	0.000	0.007	0.000	0.009
3	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	0.948	0.000	0.013	0.000	0.004
5	0.000	0.000	0.000	0.000	1.000	0.007	0.000	0.000
6	0.009	0.000	0.000	0.021	0.000	0.948	0.000	0.000
7	0.006	0.000	0.000	0.000	0.000	0.000	1.000	0.004
8	0.025	0.004	0.000	0.031	0.000	0.013	0.000	0.983



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

```
# 

```

```
random_cfl1.fit(X_train,y_train)
```

👤 Fitting 3 folds for each of 10 candidates, totalling 30 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 5 tasks | elapsed: 4.1min
[Parallel(n_jobs=-1)]: Done 10 tasks | elapsed: 9.8min
[Parallel(n_jobs=-1)]: Done 17 tasks | elapsed: 31.8min
[Parallel(n_jobs=-1)]: Done 27 out of 30 | elapsed: 48.9min remaining: 5.4min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 54.2min finished
RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
 estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
 colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
 max_depth=3, min_child_weight=1, missing=None, n_estimators=100,
 n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
 reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
 silent=True, subsample=1),
 fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
 param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15, 0.2], 'n_e
 pre_dispatch='2*n_jobs', random_state=None, refit=True,
 return_train_score='warn', scoring=None, verbose=10)

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

👤 {'subsample': 1, 'n_estimators': 2000, 'max_depth': 5, 'learning_rate': 0.01, 'colsample

```
x_cfl=XGBClassifier(n_estimators=2000, learning_rate=0.01, colsample_bytree=0.5, max_depth=5,  
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 0.024348854759529453
cv loss 0.06210692336009718
test loss 0.07717065282799755

4.2 Modeling with .asm files

There are 10868 files of asm

All the files make up about 150 GB

The asm files contains :

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

With the help of parallel processing we extracted all the features. In parallel we can use al

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

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

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

4.2.1 Feature extraction from asm files

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

```
#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')

```

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

```

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

    prefixes = ['HEADER','.text','.Pav','.idata','.data','.bss','.rdata','.edata','.r'
    #this are opcodes that are used to get best results
    #https://en.wikipedia.org/wiki/X86\_instruction\_listings

    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec'
    #best keywords that are taken from different blogs
    keywords = ['.dll','std::',':dword']
    #Below taken registers are general purpose registers and special registers
    #All the registers which are taken are best
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\asmsmallfile.txt","w+")
    files = os.listdir('first')
    for f in files:
        #filling the values with zeros into the arrays
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        # https://docs.python.org/3/library/codecs.html#codecs.ignore\_errors
        # https://docs.python.org/3/library/codecs.html#codecs.Codec.encode
        with codecs.open('first/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                # https://www.tutorialspoint.com/python3/string\_rstrip.htm
                line=lines.rstrip().split()

```

```

l=line[0]
#counting the prefixes in each and every line
for i in range(len(prefixes)):
    if prefixes[i] in line[0]:
        prefixescount[i]+=1
line=line[1:]
#counting the opcodes in each and every line
for i in range(len(opcodes)):
    if any(opcodes[i]==li for li in line):
        features.append(opcodes[i])
        opcodescount[i]+=1
#counting registers in the line
for i in range(len(registers)):
    for li in line:
        # we will use registers only in 'text' and 'CODE' segments
        if registers[i] in li and ('text' in l or 'CODE' in l):
            registerscount[i]+=1
#counting keywords in the line
for i in range(len(keywords)):
    for li in line:
        if keywords[i] in li:
            keywordcount[i]+=1
#pushing the values into the file after reading whole file
for prefix in prefixescount:
    file1.write(str(prefix)+",")
for opcode in opcodescount:
    file1.write(str(opcode)+",")
for register in registerscount:
    file1.write(str(register)+",")
for key in keywordcount:
    file1.write(str(key)+",")
file1.write("\n")
file1.close()

#same as above
def secondprocess():
    prefixes = ['HEADER','.text','.Pav','.idata','.data','.bss','.rdata','.edata','.r'
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec'
    keywords = ['.dll','std::','::dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\mediumasmfile.txt","w+")
    files = os.listdir('second')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")

```

```

opcoaet1ie.write(tz+ " ")
with codecs.open('second/'+f,encoding='cp1252',errors ='replace') as fli:
    for lines in fli:
        line=lines.rstrip().split()
        l=line[0]
        for i in range(len(prefixes)):
            if prefixes[i] in line[0]:
                prefixescount[i]+=1
        line=line[1:]
        for i in range(len(opcodes)):
            if any(opcodes[i]==li for li in line):
                features.append(opcodes[i])
                opcodescount[i]+=1
        for i in range(len(registers)):
            for li in line:
                if registers[i] in li and ('text' in l or 'CODE' in l):
                    registerscount[i]+=1
        for i in range(len(keywords)):
            for li in line:
                if keywords[i] in li:
                    keywordcount[i]+=1
    for prefix in prefixescount:
        file1.write(str(prefix)+",")
    for opcode in opcodescount:
        file1.write(str(opcode)+",")
    for register in registerscount:
        file1.write(str(register)+",")
    for key in keywordcount:
        file1.write(str(key)+",")
    file1.write("\n")
file1.close()

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

```

```

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

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

```

```

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

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

```

```

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

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

if __name__=="__main__":
    main()

# asmoutputfile.csv(output generated from the above two cells) will contain all the extracted
# this file will be uploaded in the drive, you can directly use this
dfasm=pd.read_csv("asmoutputfile.csv")

```

```

Y.columns = ['ID', 'Class']
result_asm = pd.merge(dfasm, Y, on='ID', how='left')
result_asm.head()

```



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

5 rows × 53 columns

4.2.1.1 Files sizes of each .asm file

```

#file sizes of byte files

files=os.listdir('asmFiles')
filenames=Y['ID'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, st_nlink=1, s
    # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=1519638522)
    # read more about os.stat: here https://www.tutorialspoint.com/python/os\_stat.htm
    statinfo=os.stat('asmFiles/'+file)
    # split the file name at '.' and take the first part of it i.e the file name
    file=file.split('.')[0]
    if any(file == filename for filename in filenames):
        i=filenames.index(file)
        class_bytes.append(class_y[i])
        # converting into Mb's
        sizebytes.append(statinfo.st_size/(1024.0*1024.0))
        fnames.append(file)
asm_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes})
print (asm_size_byte.head())

```



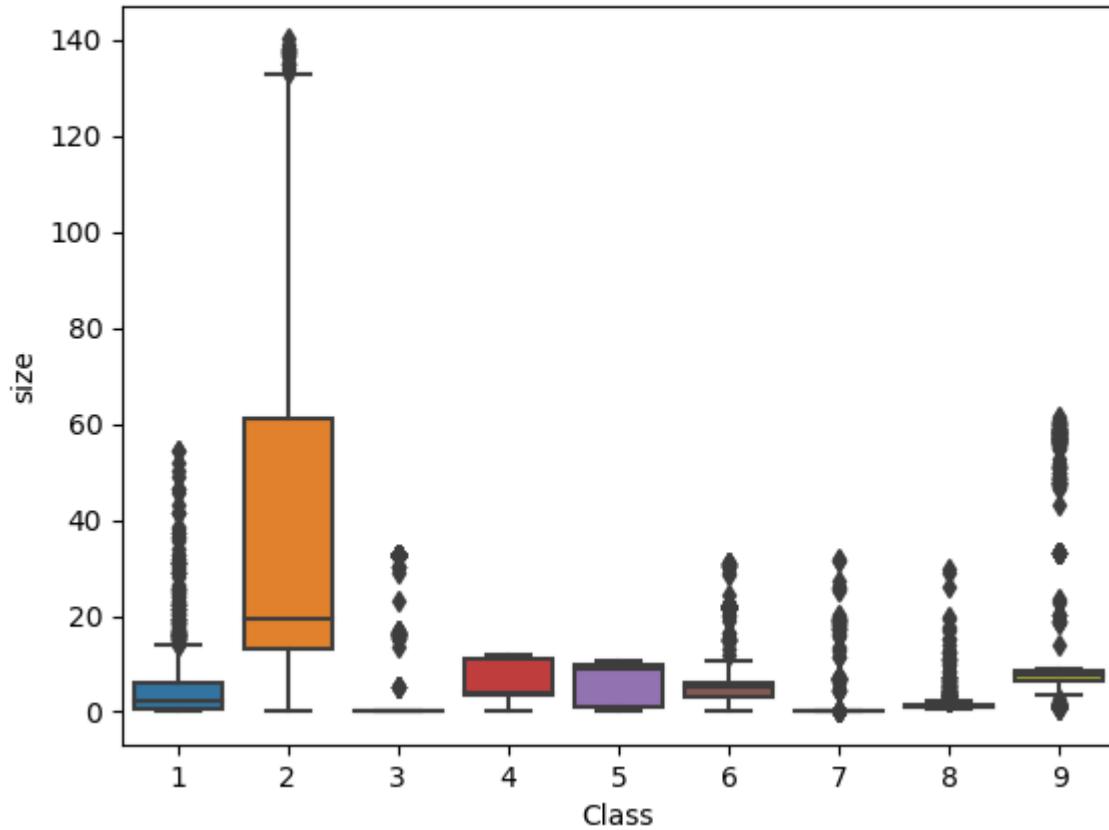
TD size Class

4.2.1.2 Distribution of .asm file sizes

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



boxplot of .bytes file sizes



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



```
(10868, 53)
(10868, 3)
```

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

5 rows × 54 columns

```
# we normalize the data each column
result_asm.head()
```



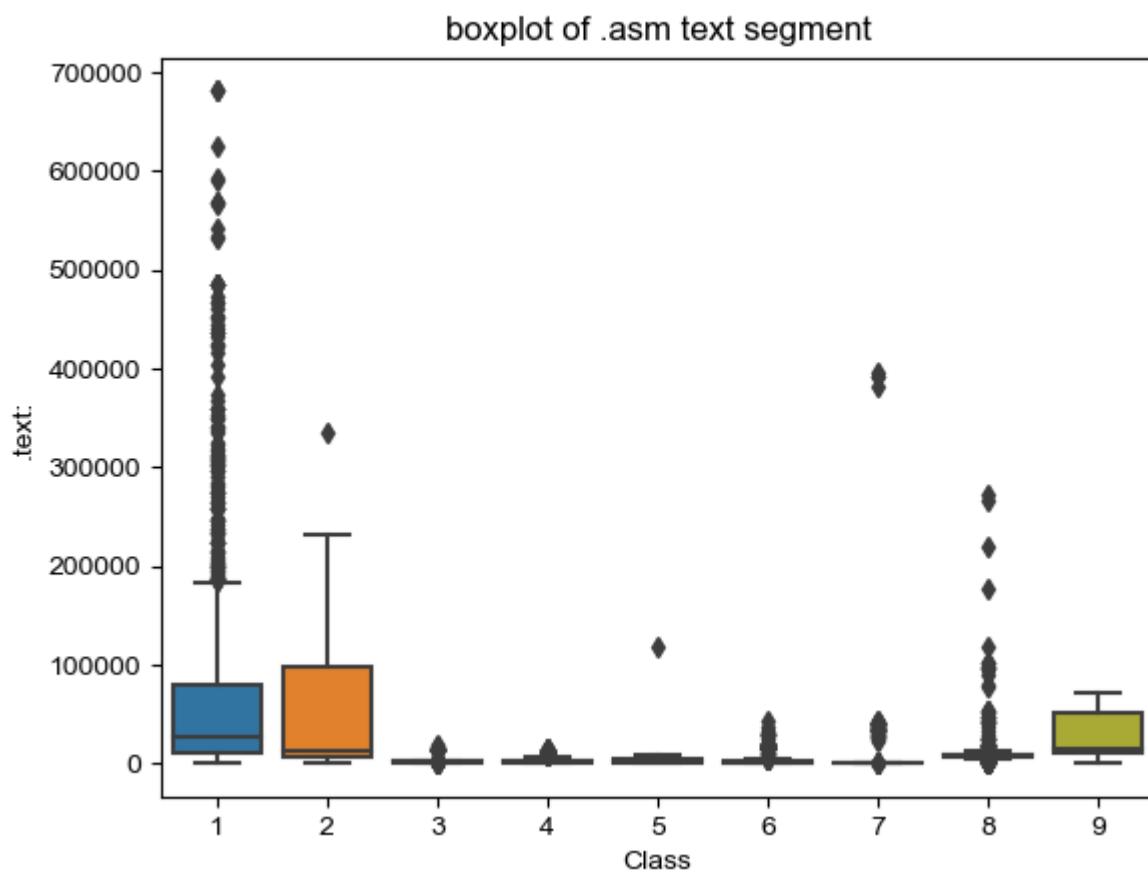
	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.eda
0	01kcPWA9K2BOxQeS5Rju	19	744	0	127	57	0	323	
1	1E93CpP60RHFNiT5Qfvn	17	838	0	103	49	0	0	
2	3ekVow2ajZHbTnBcsDfX	17	427	0	50	43	0	145	
3	3X2nY7iQaPBIWDrAZqJe	17	227	0	43	19	0	0	
4	46OZzdsSKDCFV8h7XWxf	17	402	0	59	170	0	0	

5 rows × 54 columns

4.2.2 Univariate analysis on asm file features

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

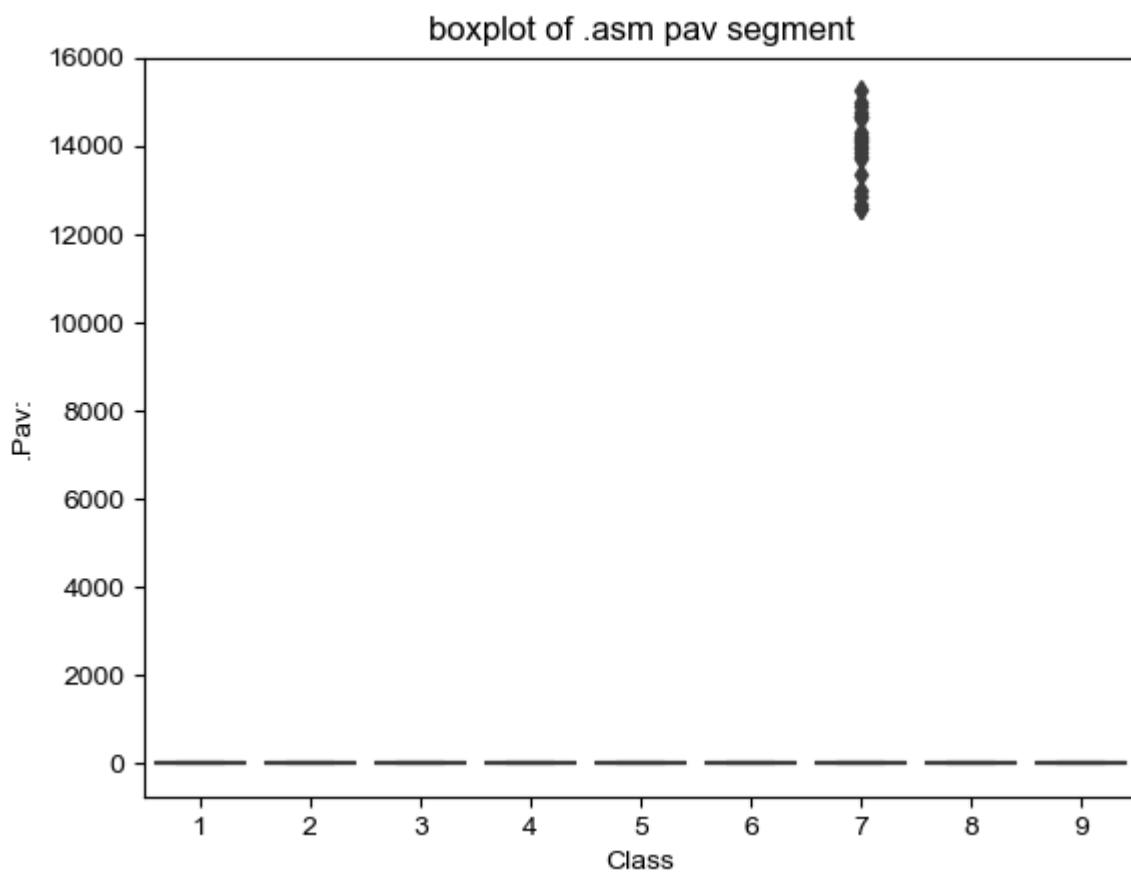




The plot is between Text and class
Class 1,2 and 9 can be easily separated

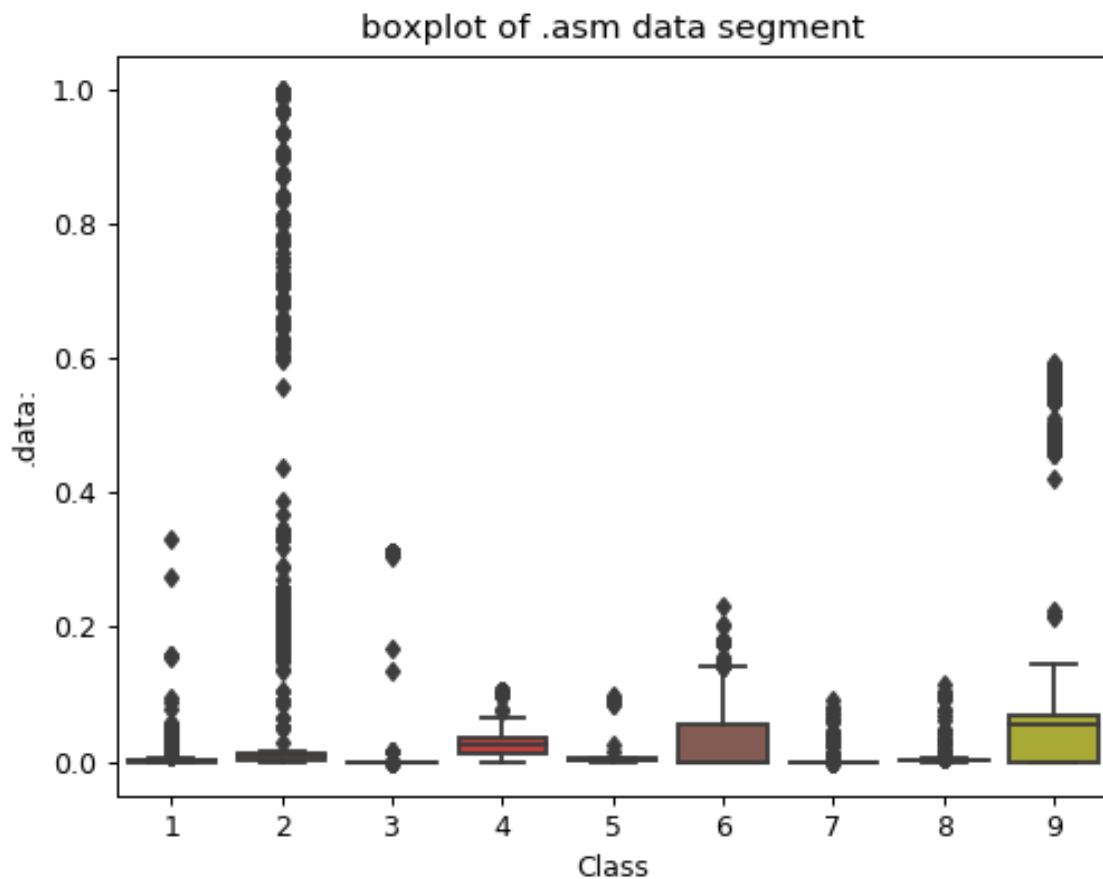
```
ax = sns.boxplot(x="Class", y=".Pav:", data=result_asm)
plt.title("boxplot of .asm pav segment")
plt.show()
```





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

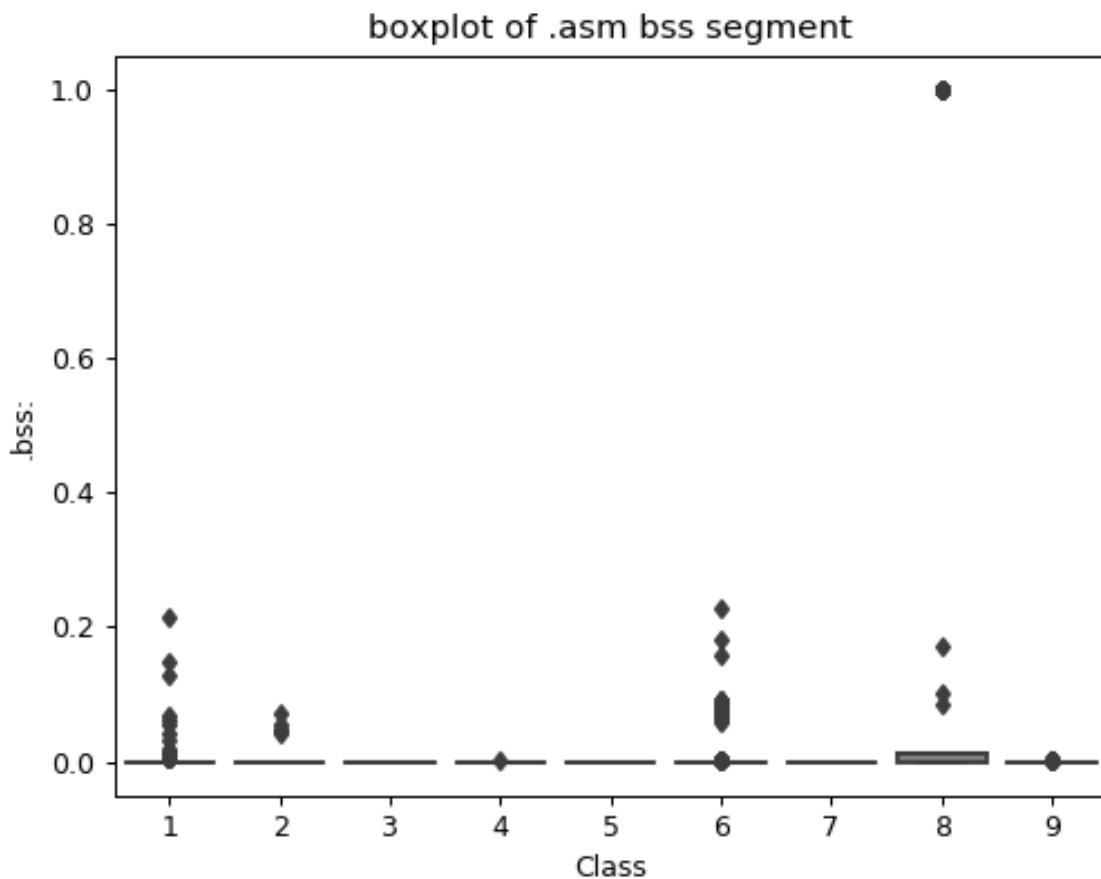




The plot is between data segment and class label
class 6 and class 9 can be easily separated from given points

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

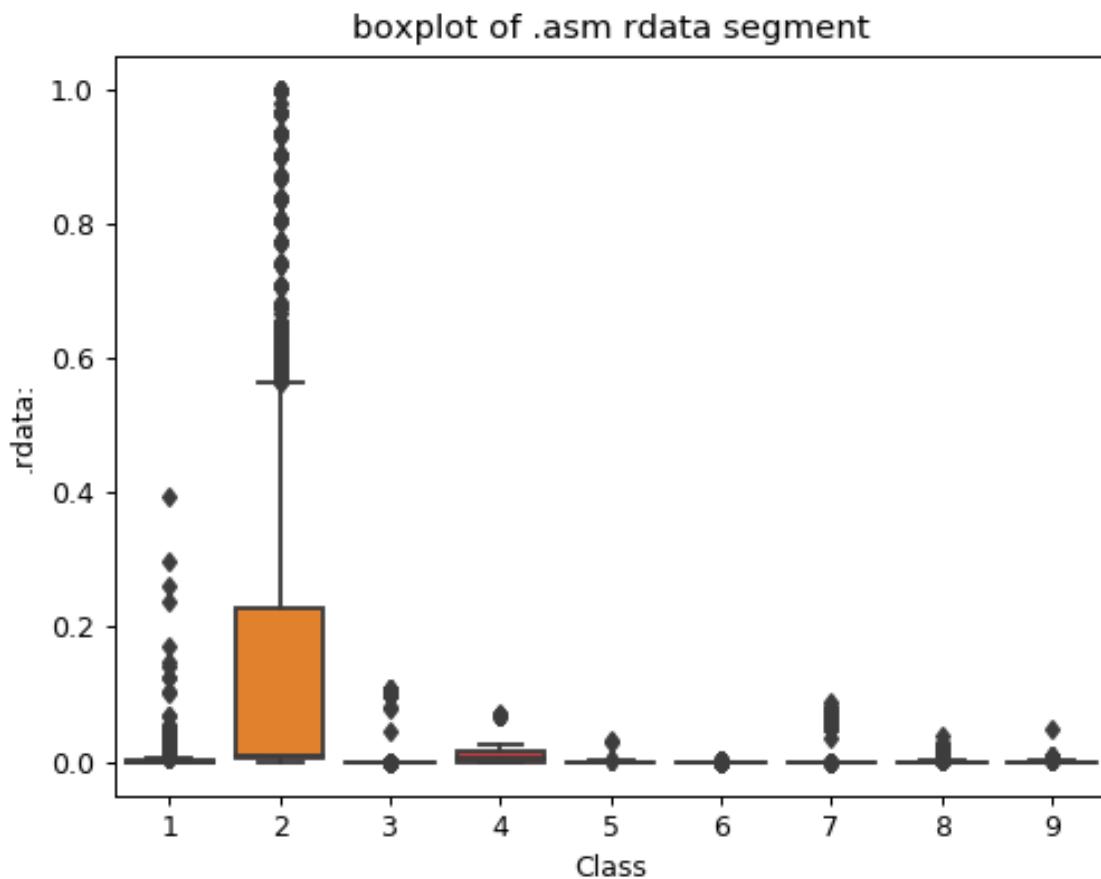




plot between bss segment and class label
very less number of files are having bss segment

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



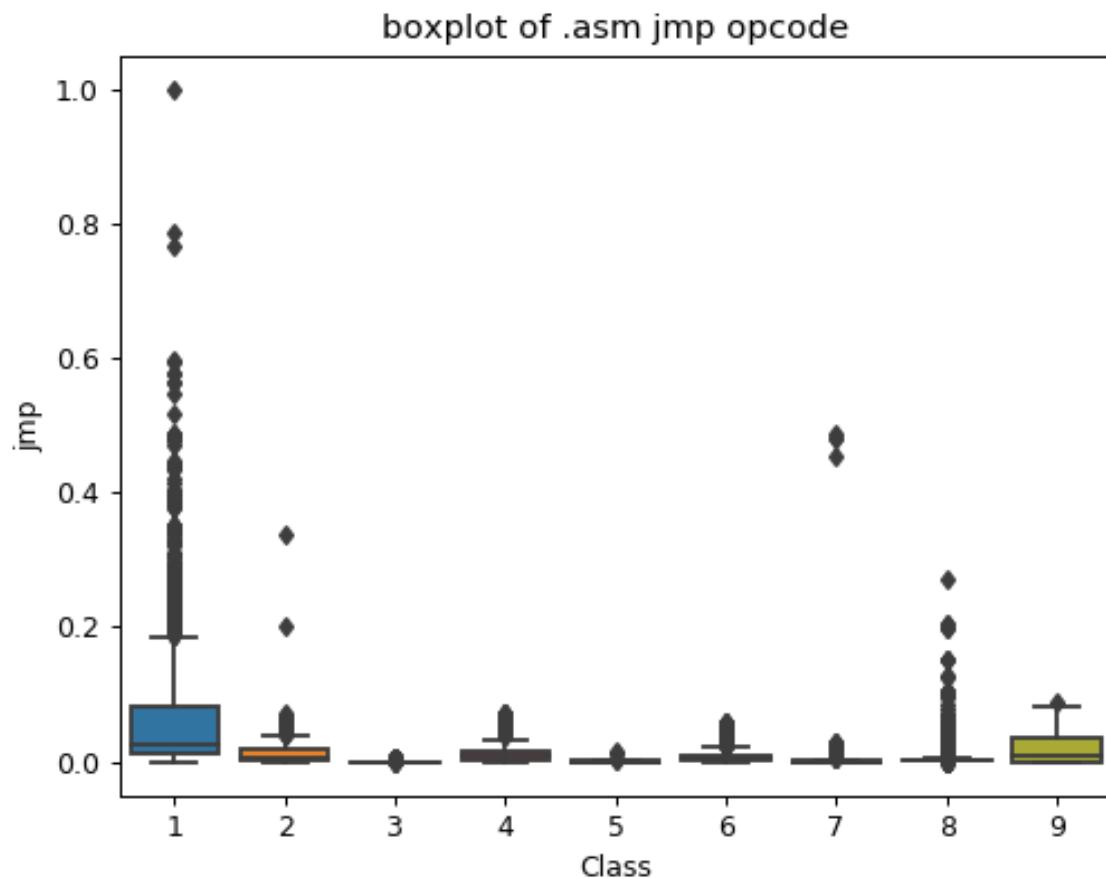


Plot between rdata segment and Class segment

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

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



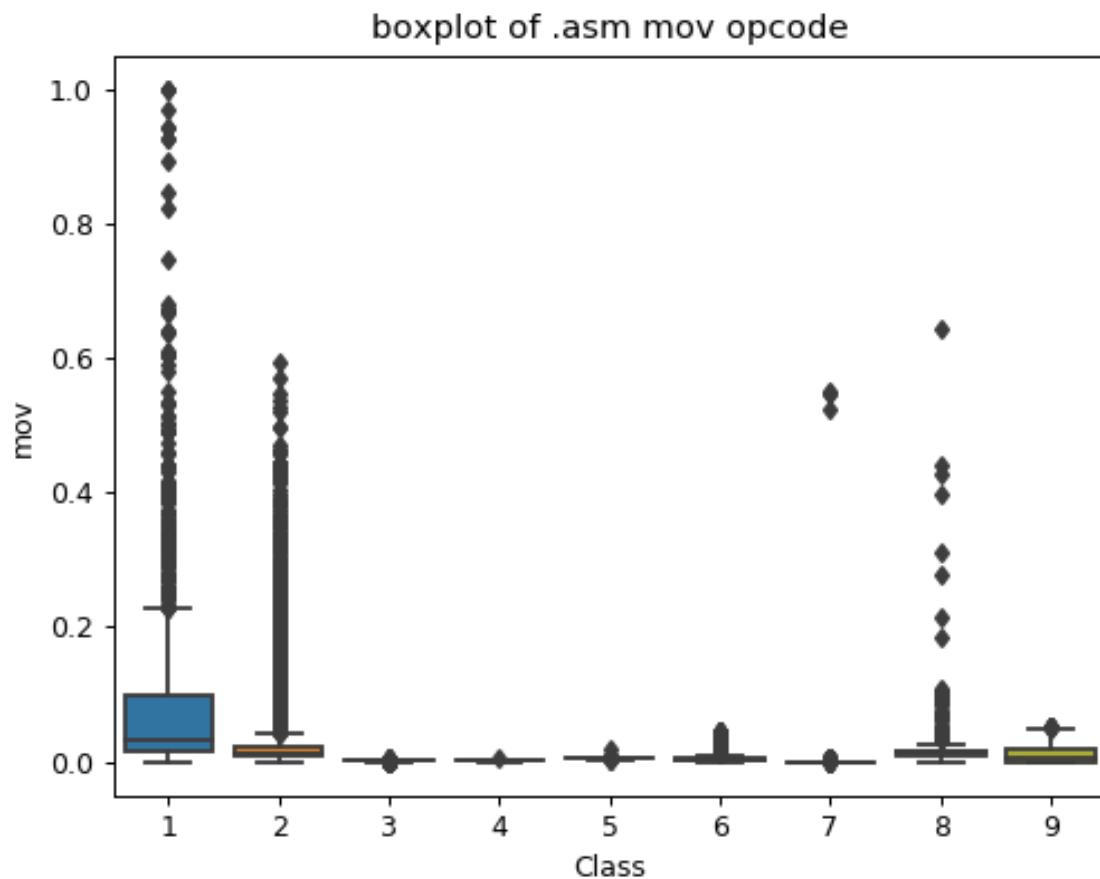


plot between jmp and Class label

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

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



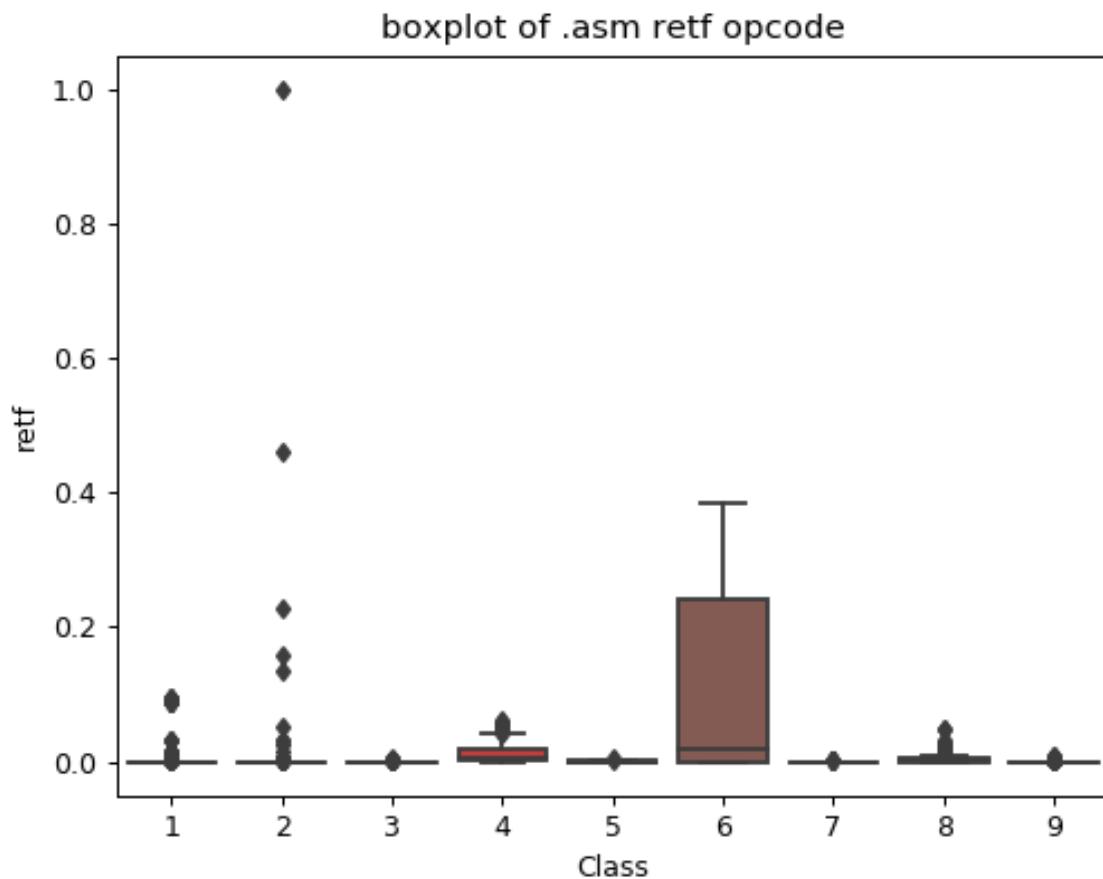


plot between Class label and mov opcode

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

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

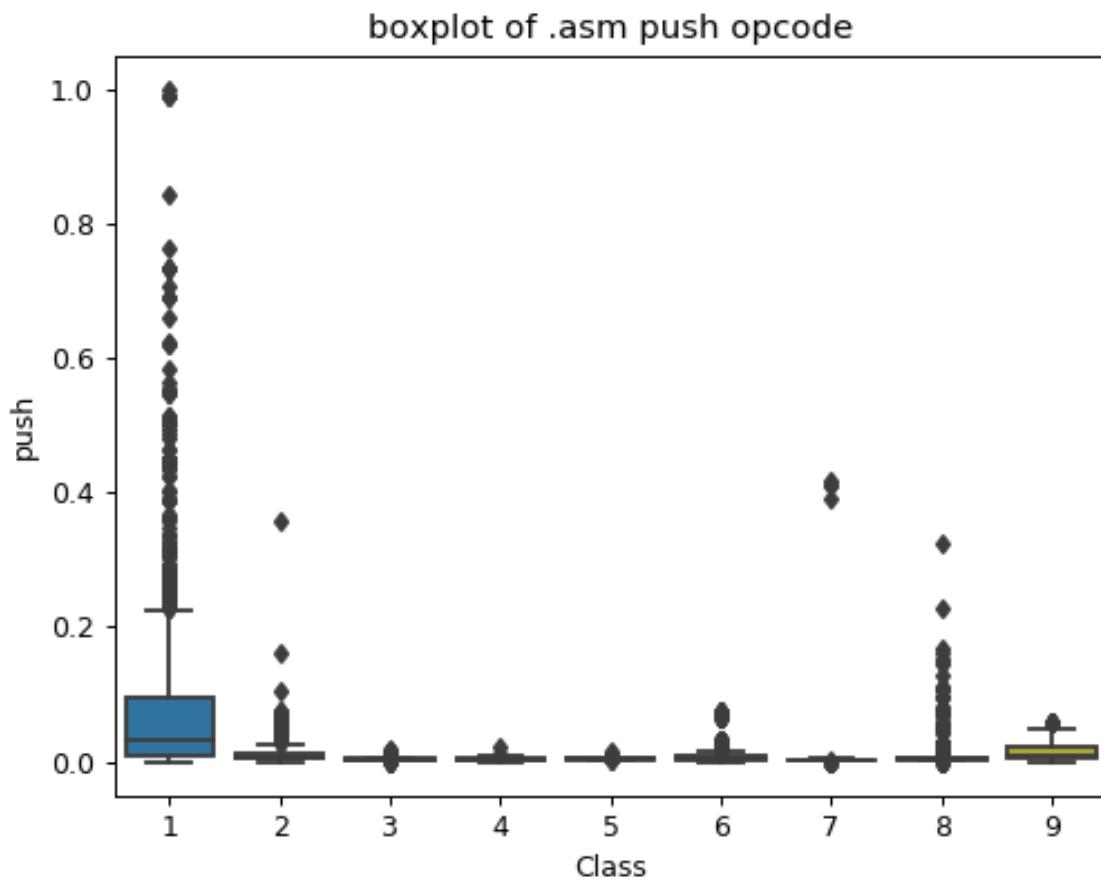




plot between Class label and retf
Class 6 can be easily separated with opcode retf
The frequency of retf is approx of 250.

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





plot between push opcode and Class label

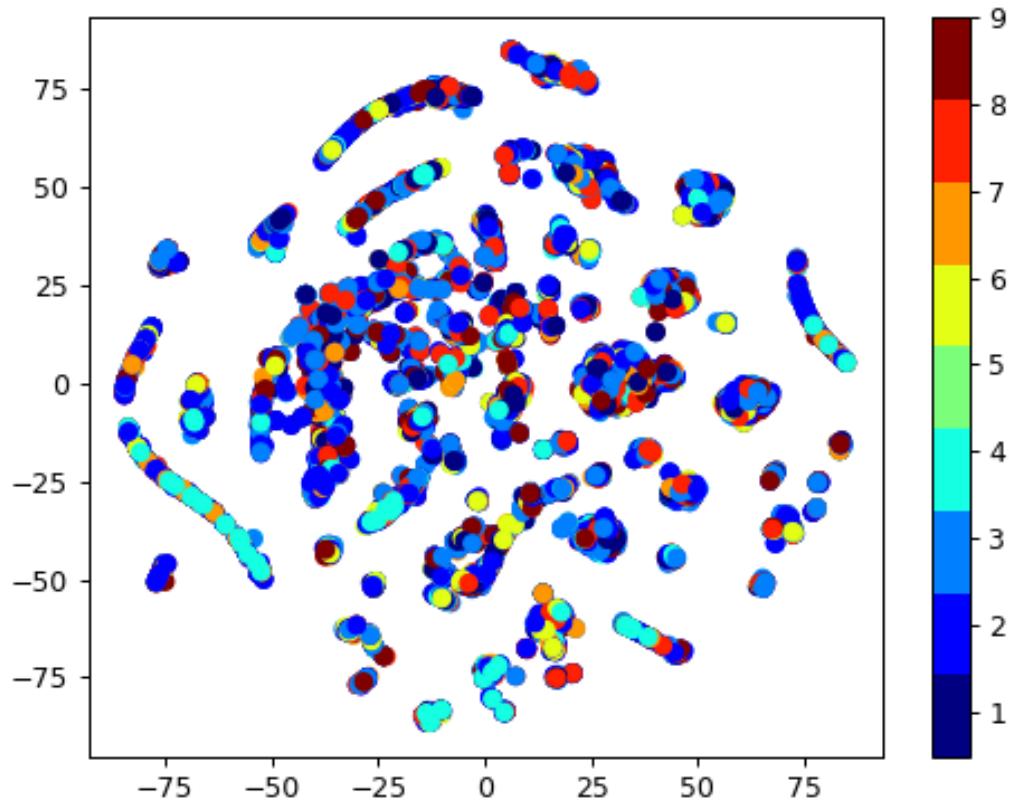
Class 1 is having 75 percentile files with push opcodes of frequency 1000

4.2.2 Multivariate Analysis on .asm file features

```
# check out the course content for more explanation on tsne algorithm
# https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/t-distributed-stochastic-nearest-neighbor

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

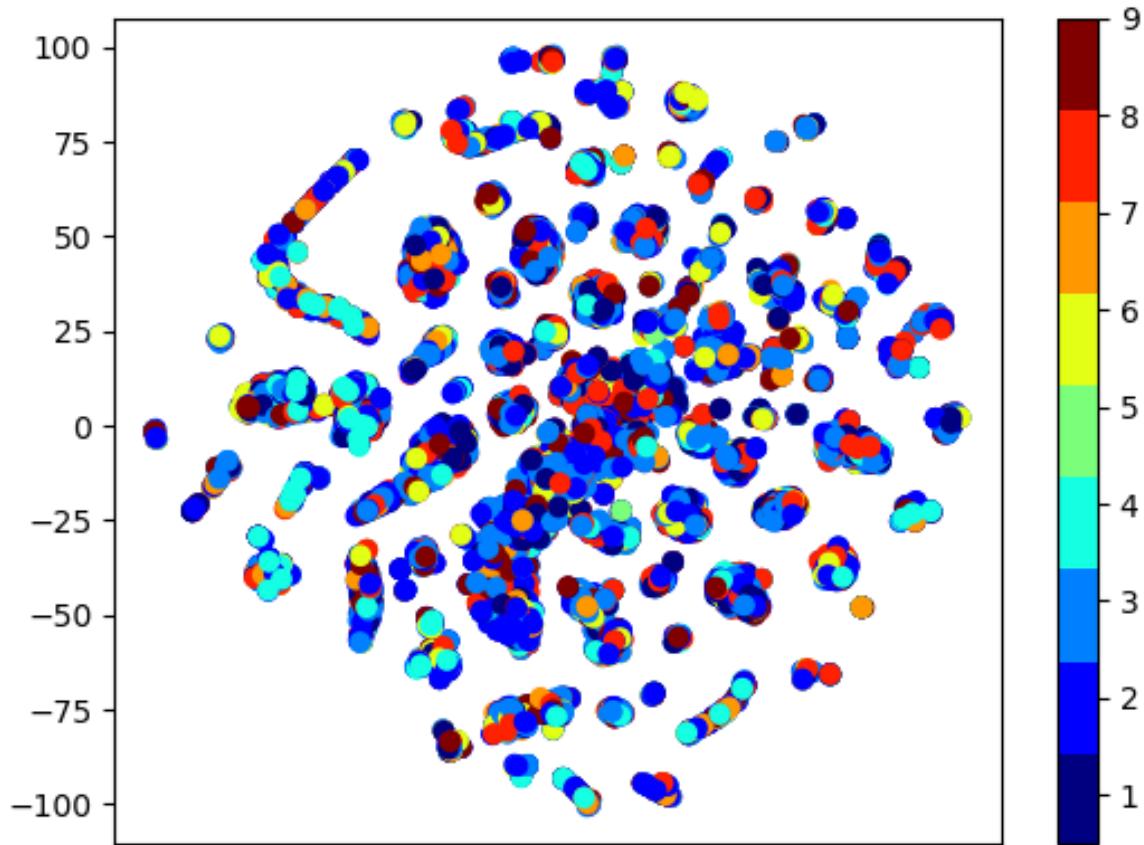
```
plt.colorbar(ticks=range(10))
plt.clim(0.5, 9)
plt.show()
```



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

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





TSNE for asm data with perplexity 50

4.2.3 Conclusion on EDA

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

4.3 Train and test split

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

X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x,asm_y ,stratify=asm
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_asm, stra

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



```
HEADER:    False
.text:    False
.Pav:    False
...
```

4.4. Machine Learning models on features of .asm files

```
.rsrc:    False
```

4.4.1 K-Nearest Neighbors

```
...    False
```

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

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

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

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

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

best_alpha = np.argmin(cv_log_error_array)

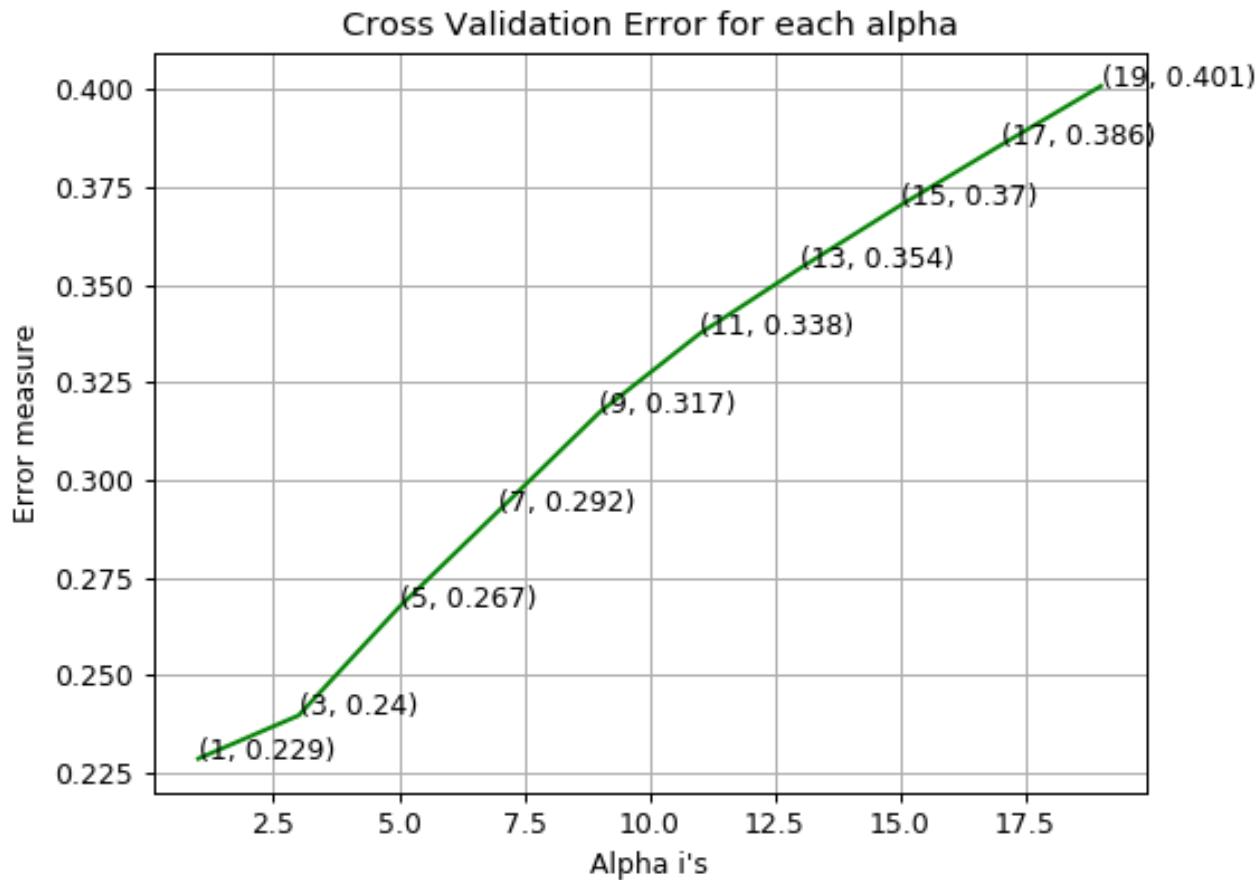
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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



```
log_loss for k = 1 is 0.2286951008264786
log_loss for k = 3 is 0.23974604909767921
log_loss for k = 5 is 0.26743767182569295
log_loss for k = 7 is 0.2922812711849662
log_loss for k = 9 is 0.3173517943176062
log_loss for k = 11 is 0.3375272301343973
log_loss for k = 13 is 0.3542581717184334
log_loss for k = 15 is 0.3703567252351854
log_loss for k = 17 is 0.3857570471590401
log_loss for k = 19 is 0.40090334916939535
```



```
log loss for train data 0.07371820550676085
log loss for cv data 0.2286951008264786
log loss for test data 0.2135354369723588
Number of misclassified points 4.001839926402944
```

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

	1	289.000	3.000	3.000	1.000	0.000	1.000	1.000	7.000
	2	1.000	484.000	4.000	2.000	0.000	0.000	0.000	1.000
	3	0.000	1.000	586.000	0.000	0.000	0.000	0.000	1.000
class	4	0.000	0.000	0.000	89.000	1.000	3.000	0.000	2.000

Original Class	1	2	3	4	5	6	7	8
5	0.000	0.000	0.000	0.000	5.000	0.000	0.000	2.000
6	4.000	3.000	0.000	4.000	0.000	136.000	1.000	2.000
7	1.000	3.000	1.000	0.000	0.000	0.000	75.000	0.000
8	7.000	3.000	0.000	4.000	1.000	4.000	0.000	223.000
9	2.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000

Predicted Class

Precision matrix

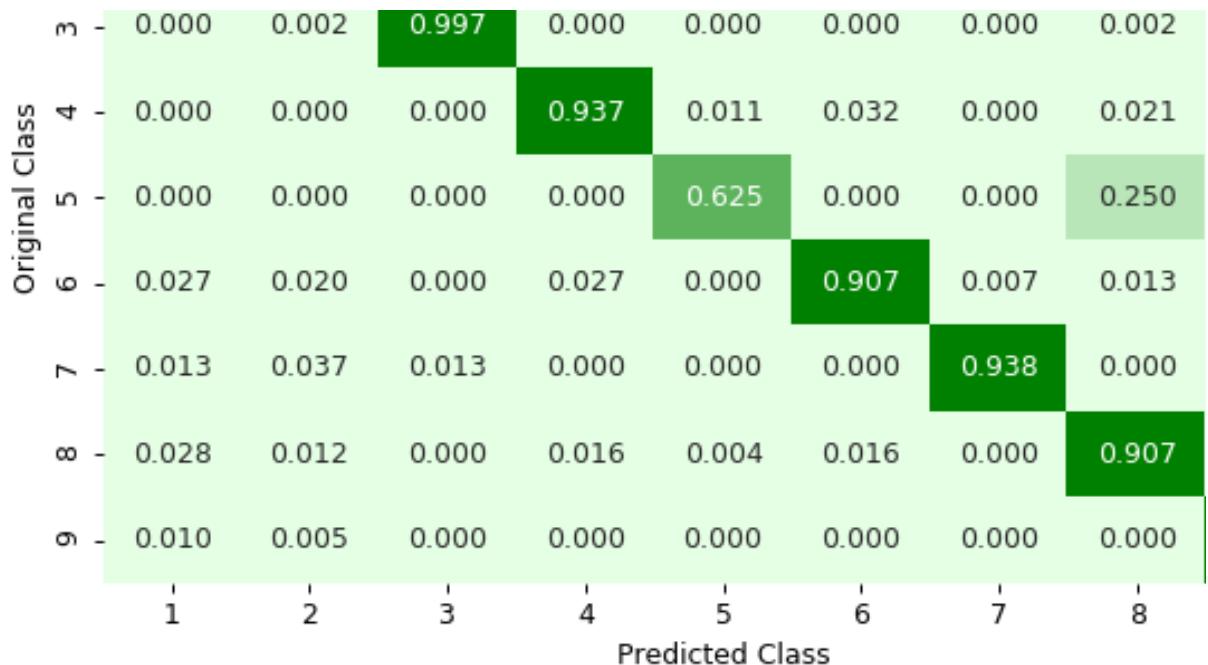
Original Class	1	2	3	4	5	6	7	8
1	0.951	0.006	0.005	0.010	0.000	0.007	0.013	0.029
2	0.003	0.972	0.007	0.020	0.000	0.000	0.000	0.004
3	0.000	0.002	0.987	0.000	0.000	0.000	0.000	0.004
4	0.000	0.000	0.000	0.890	0.143	0.021	0.000	0.008
5	0.000	0.000	0.000	0.000	0.714	0.000	0.000	0.008
6	0.013	0.006	0.000	0.040	0.000	0.944	0.013	0.008
7	0.003	0.006	0.002	0.000	0.000	0.000	0.974	0.000
8	0.023	0.006	0.000	0.040	0.143	0.028	0.000	0.937
9	0.007	0.002	0.000	0.000	0.000	0.000	0.000	0.000

Predicted Class

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

Recall matrix

Original Class	1	2	3	4	5	6	7	8
1	0.938	0.010	0.010	0.003	0.000	0.003	0.003	0.023
2	0.002	0.976	0.008	0.004	0.000	0.000	0.000	0.002



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

4.4.2 Logistic Regression

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

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

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

```
alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
```

```
--o_-----\_, ,_-----, predict_y = sig_clf.predict_proba(X_cv_asm)
cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.classes_, eps=1e

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

best_alpha = np.argmin(cv_log_error_array)

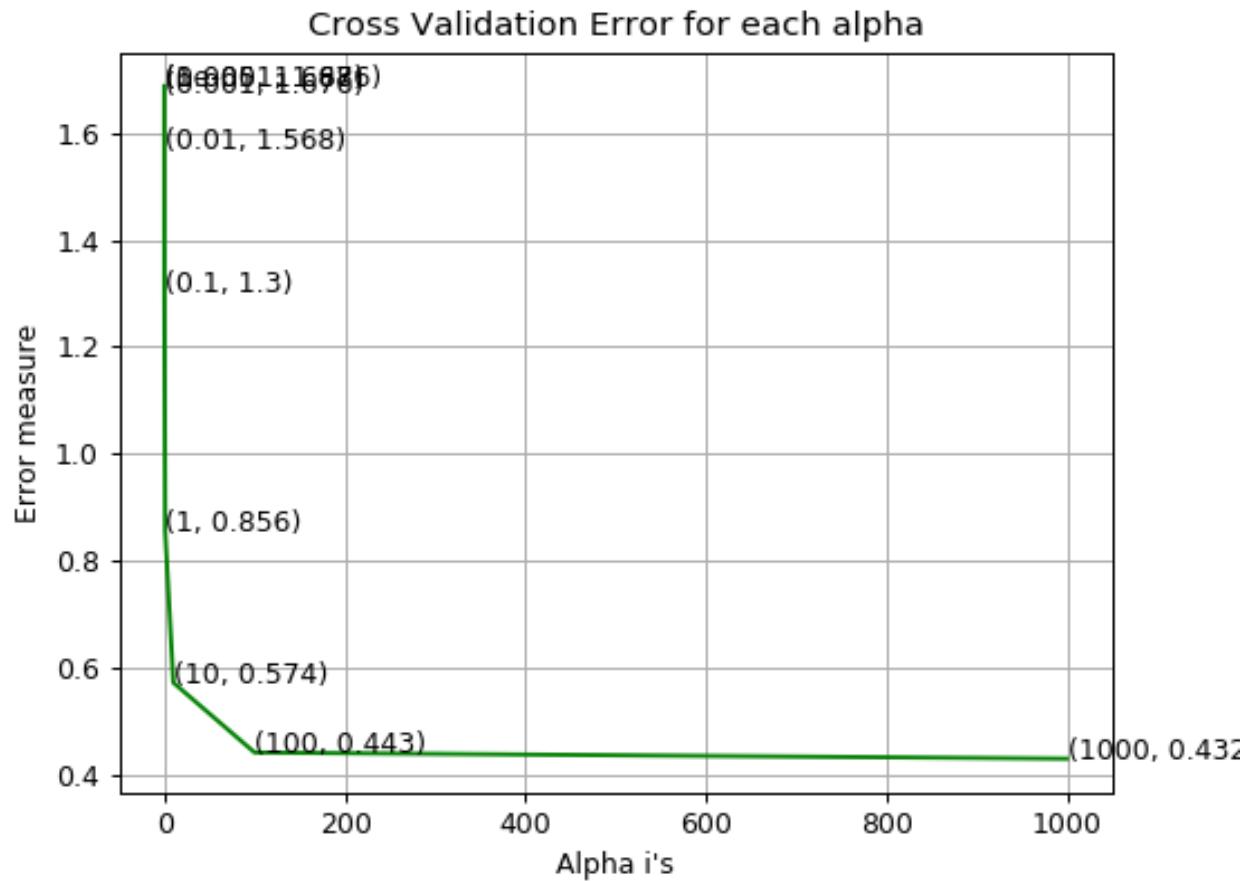
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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



```
log_loss for c = 1e-05 is 1.6869859868957804
log_loss for c = 0.0001 is 1.6855010192472757
log_loss for c = 0.001 is 1.6755781562152487
log_loss for c = 0.01 is 1.5677273322121714
log_loss for c = 0.1 is 1.3002573116338927
log_loss for c = 1 is 0.856048258533692
log_loss for c = 10 is 0.5735687649879864
log_loss for c = 100 is 0.4431214718098947
log_loss for c = 1000 is 0.43157353232283385
```



log loss for train data 0.38150548196180933

log loss for cv data 0.43157353232283385

log loss for test data 0.38308926013384326

Number of misclassified points 7.405703771849126

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

		1	2	3	4	5	6	7	8
1	-	281.000	2.000	9.000	2.000	0.000	0.000	4.000	6.000
2	-	0.000	496.000	0.000	0.000	0.000	0.000	0.000	0.000
3	-	0.000	0.000	586.000	0.000	0.000	1.000	0.000	1.000
4	-	1.000	2.000	13.000	74.000	0.000	0.000	5.000	0.000
5	-	1.000	0.000	0.000	2.000	4.000	0.000	0.000	1.000

Original Class	1	2	3	4	5	6	7	8
Predicted Class	1	2	3	4	5	6	7	8
1	1.000	0.000	0.000	2.000	4.000	0.000	0.000	1.000
2	0.000	0.000	23.000	0.000	0.000	121.000	3.000	2.000
3	2.000	4.000	3.000	0.000	0.000	0.000	69.000	0.000
4	34.000	4.000	1.000	0.000	0.000	0.000	6.000	192.000
5	0.000	0.000	0.000	0.000	0.000	8.000	1.000	4.000

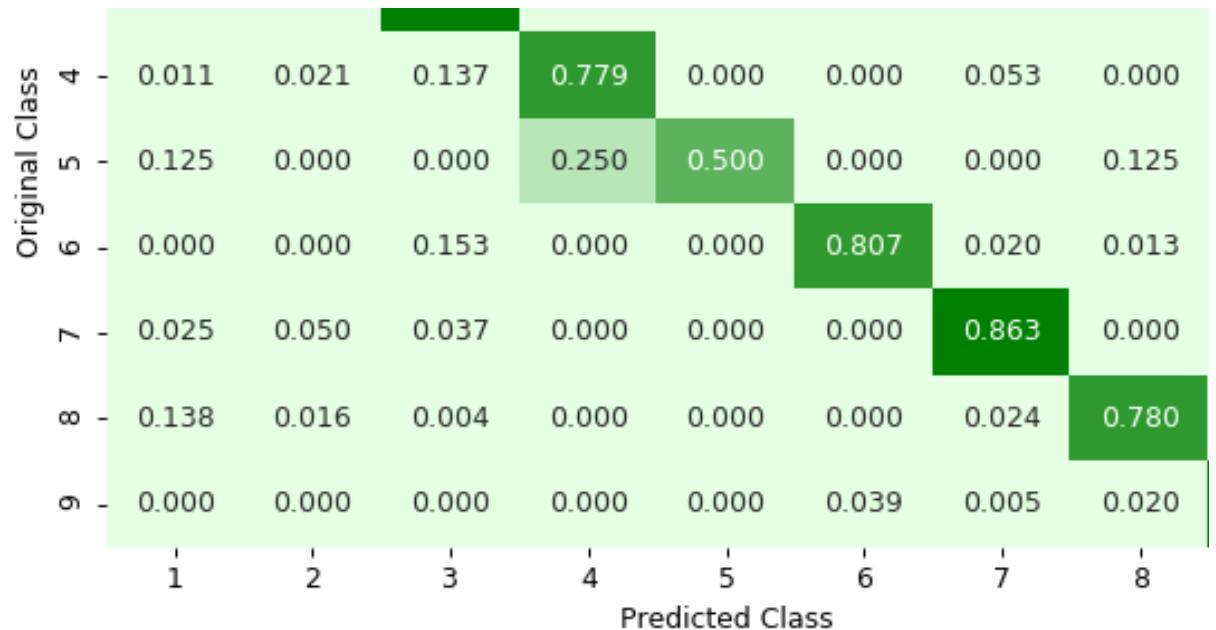
----- Precision matrix -----

Original Class	1	2	3	4	5	6	7	8
Predicted Class	1	2	3	4	5	6	7	8
1	0.881	0.004	0.014	0.026	0.000	0.000	0.045	0.029
2	0.000	0.976	0.000	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	0.923	0.000	0.000	0.008	0.000	0.005
4	0.003	0.004	0.020	0.949	0.000	0.000	0.057	0.000
5	0.003	0.000	0.000	0.026	1.000	0.000	0.000	0.005
6	0.000	0.000	0.036	0.000	0.000	0.931	0.034	0.010
7	0.006	0.008	0.005	0.000	0.000	0.000	0.784	0.000
8	0.107	0.008	0.002	0.000	0.000	0.000	0.068	0.932
9	0.000	0.000	0.000	0.000	0.000	0.062	0.011	0.019

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

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

Original Class	1	2	3	4	5	6	7	8
Predicted Class	1	2	3	4	5	6	7	8
1	0.912	0.006	0.029	0.006	0.000	0.000	0.013	0.019
2	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	0.997	0.000	0.000	0.002	0.000	0.002



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

4.4.3 Random Forest Classifier

```

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

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

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

# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-classifier
# -----
```

```

alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    r_f1=RandomForestClassifier(n_estimators=i random_state=42 n_inhcs=-1)
    https://colab.research.google.com/drive/1R4RdmKhYl7r3RUjrmP13A1dgDCHE2y2H#scrollTo=33N93b0ldWK2&printMode=true
    68/108

```

```
r_cfl=RandomForestClassifier(n_estimators=50,random_state=42,n_jobs=-1)
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)
predict_y = sig_clf.predict_proba(X_cv_asm)
cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_, eps=1e-15))

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

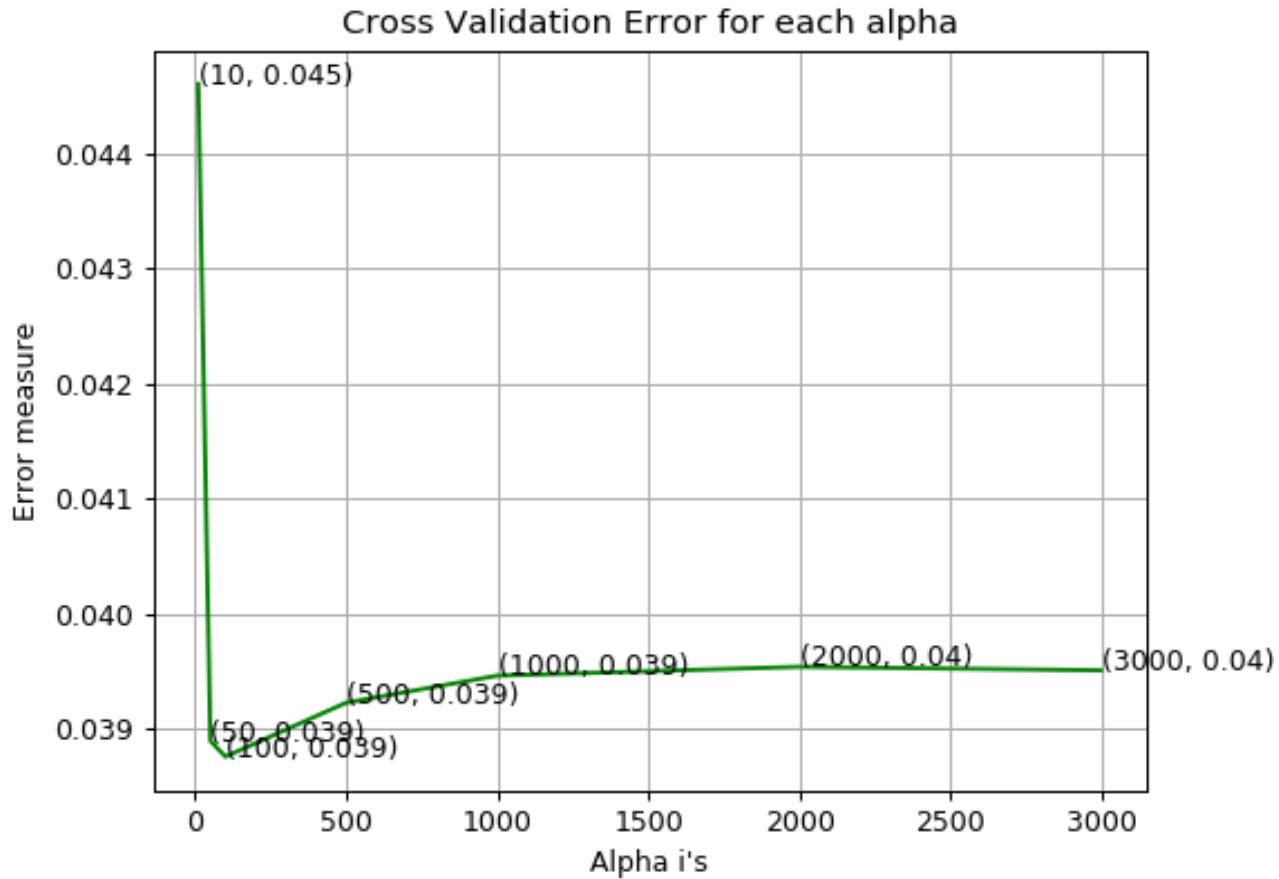
best_alpha = np.argmin(cv_log_error_array)

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

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



```
log_loss for c = 10 is 0.044604000720488056
log_loss for c = 50 is 0.038892329851189296
log_loss for c = 100 is 0.03875524544813011
log_loss for c = 500 is 0.039224440805809314
log_loss for c = 1000 is 0.03945941790783839
log_loss for c = 2000 is 0.03953659123286974
log_loss for c = 3000 is 0.03950608587732239
```



```
log loss for train data 0.012379247850927044
log loss for cv data 0.03875524544813011
log loss for test data 0.039428378936875376
Number of misclassified points 0.8279668813247469
```

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

Original Class	1	2	3	4	5	6	7	8
1	306.000	0.000	0.000	0.000	0.000	0.000	0.000	2.000
2	0.000	495.000	0.000	0.000	0.000	0.000	0.000	1.000
3	0.000	0.000	587.000	0.000	0.000	0.000	0.000	1.000
4	0.000	0.000	0.000	94.000	0.000	0.000	0.000	1.000
5	0.000	0.000	0.000	0.000	8.000	0.000	0.000	0.000

Original Class	6	3.000	0.000	0.000	0.000	0.000	146.000	0.000	1.000
7	-	1.000	0.000	0.000	0.000	0.000	0.000	79.000	0.000
8	-	3.000	0.000	0.000	1.000	0.000	1.000	1.000	240.000
9	-	1.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000
	1	2	3	4	5	6	7	8	
	Predicted Class								

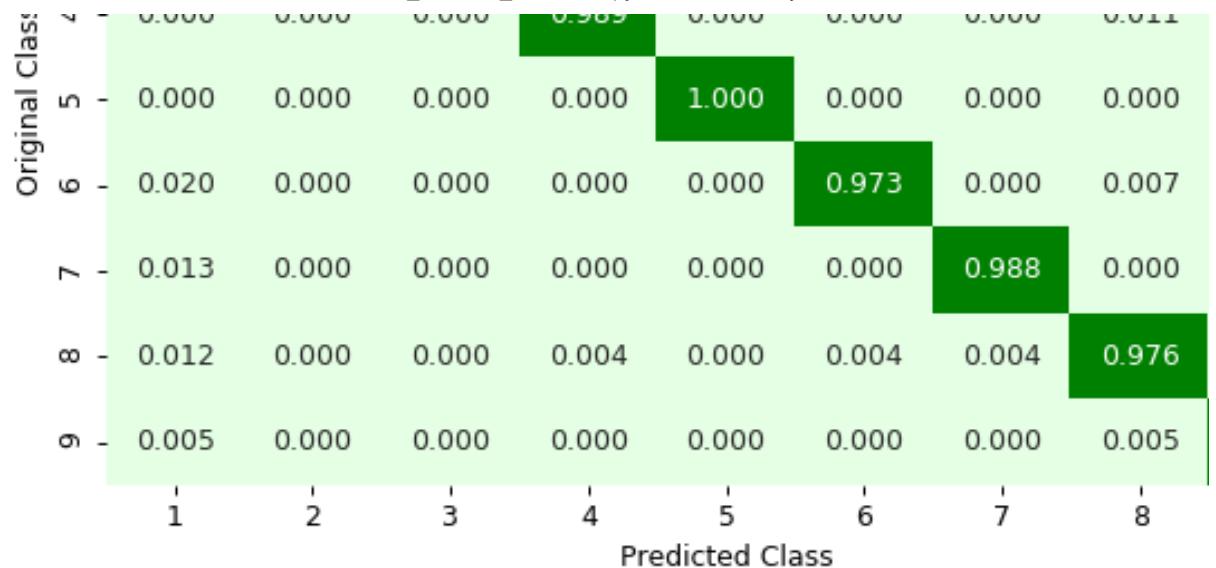
----- Precision matrix -----

Original Class	1	0.975	0.000	0.000	0.000	0.000	0.000	0.000	0.008
2	-	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.004
3	-	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.004
4	-	0.000	0.000	0.000	0.989	0.000	0.000	0.000	0.004
5	-	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000
6	-	0.010	0.000	0.000	0.000	0.000	0.993	0.000	0.004
7	-	0.003	0.000	0.000	0.000	0.000	0.000	0.988	0.000
8	-	0.010	0.000	0.000	0.011	0.000	0.007	0.013	0.972
9	-	0.003	0.000	0.000	0.000	0.000	0.000	0.000	0.004
	1	2	3	4	5	6	7	8	
	Predicted Class								

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

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

Original Class	1	0.994	0.000	0.000	0.000	0.000	0.000	0.000	0.006
2	-	0.000	0.998	0.000	0.000	0.000	0.000	0.000	0.002
3	-	0.000	0.000	0.998	0.000	0.000	0.000	0.000	0.002
4	-	0.000	0.000	0.000	0.999	0.000	0.000	0.000	0.011



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

4.4.4 XgBoost Classifier

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

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

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

```
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_cfl = CalibratedClassifierCV(x_cfl, method="sigmoid")
```

```
sig_clf.fit(X_train_asm, y_train_asm)
predict_y = sig_clf.predict_proba(X_cv_asm)
cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=x_cfl.classes_, eps=1e-15)

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

best_alpha = np.argmin(cv_log_error_array)

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

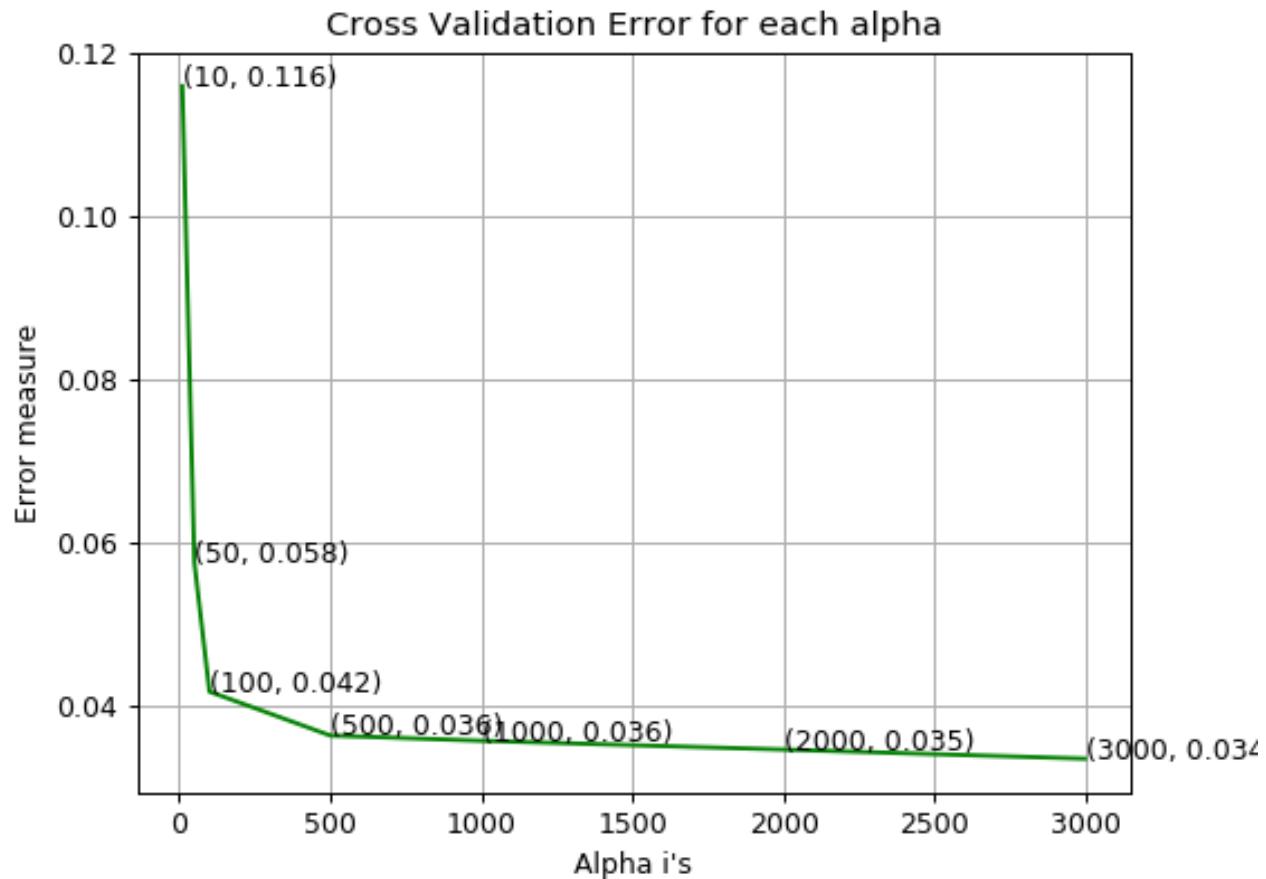
x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
x_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)

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



```
log_loss for c = 10 is 0.11588105338340265
log_loss for c = 50 is 0.057658250882591494
log_loss for c = 100 is 0.04186141305711363
log_loss for c = 500 is 0.03649854125696994
log_loss for c = 1000 is 0.035859619519393905
log_loss for c = 2000 is 0.03478236752207586
log_loss for c = 3000 is 0.033667303437409195
```



For values of best alpha = 3000 The train log loss is: 0.00982726018742022

For values of best alpha = 3000 The cross validation log loss is: 0.033667303437409195

For values of best alpha = 3000 The test log loss is: 0.042877055973511075

Number of misclassified points 0.78196872125115

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

Original Class	1	2	3	4	5	6	7	8
1	306.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000
2	0.000	494.000	1.000	0.000	0.000	0.000	1.000	0.000
3	0.000	0.000	587.000	0.000	0.000	0.000	0.000	1.000
4	1.000	0.000	0.000	94.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	8.000	0.000	0.000	0.000

Original Class	6	3.000	0.000	0.000	0.000	0.000	147.000	0.000	0.000
7	-	0.000	0.000	0.000	0.000	0.000	0.000	80.000	0.000
8	-	1.000	1.000	0.000	2.000	0.000	1.000	1.000	240.000
9	-	1.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000
	1	2	3	4	5	6	7	8	
	Predicted Class								

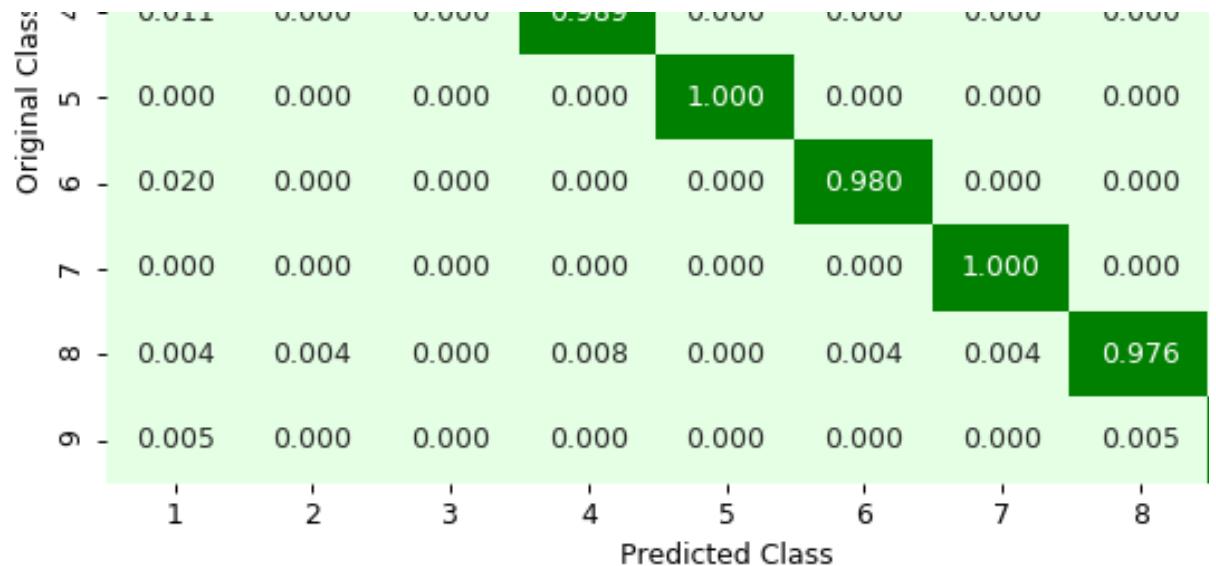
----- Precision matrix -----

Original Class	1	0.981	0.000	0.000	0.000	0.000	0.007	0.000	0.000
2	-	0.000	0.998	0.002	0.000	0.000	0.000	0.012	0.000
3	-	0.000	0.000	0.998	0.000	0.000	0.000	0.000	0.004
4	-	0.003	0.000	0.000	0.979	0.000	0.000	0.000	0.000
5	-	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000
6	-	0.010	0.000	0.000	0.000	0.000	0.987	0.000	0.000
7	-	0.000	0.000	0.000	0.000	0.000	0.000	0.976	0.000
8	-	0.003	0.002	0.000	0.021	0.000	0.007	0.012	0.992
9	-	0.003	0.000	0.000	0.000	0.000	0.000	0.000	0.004
	1	2	3	4	5	6	7	8	
	Predicted Class								

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

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

Original Class	1	0.994	0.000	0.000	0.000	0.000	0.003	0.000	0.000
2	-	0.000	0.996	0.002	0.000	0.000	0.000	0.002	0.000
3	-	0.000	0.000	0.998	0.000	0.000	0.000	0.000	0.002
4	-	0.011	0.000	0.000	0.989	0.000	0.000	0.000	0.000
	1	2	3	4	5	6	7	8	
	Predicted Class								



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

4.4.5 Xgboost Classifier with best hyperparameters

```
x_cfl=XGBClassifier()

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



```

Fitting 3 folds for each of 10 candidates, totalling 30 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done  5 tasks      | elapsed:  5.6min
[Parallel(n_jobs=-1)]: Done 10 tasks      | elapsed:  6.4min
[Parallel(n_jobs=-1)]: Done 17 tasks      | elapsed:  9.3min

print (random_cfl.best_params_)

👤 {'subsample': 0.5, 'n_estimators': 2000, 'max_depth': 10, 'learning_rate': 0.01, 'colsample_bytree': 1, 'gamma': 0, 'learning_rate': 0.1, 'max_delta_step': 0, 'subsample': 1, 'colsample_bytree': 1, 'colsample_bylevel': 1, 'reg_alpha': 0, 'reg_lambda': 1, 'scale_pos_weight': 1, 'base_score': 0.5, 'random_state': 0, 'seed': None, 'missing': None, **kwargs}

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

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

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


x_cfl=XGBClassifier(n_estimators=2000,subsample=0.5,learning_rate=0.01,colsample_bytree=0.5,max_depth=10)
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 0.010626478719576738
cv loss 0.033016089856804966
test loss 0.04082419520278345

```

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

4.5.1. Merging both asm and byte file features

```
result.head()
```

	Unnamed: 0	ID	0	1	2	3	4
0	0.000000	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048
1	0.000092	01IsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303
2	0.000184	01jsnpXSAlgw6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464
3	0.000276	01kcPWA9K2B0xQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770
4	0.000368	01SuzwMJEIXsK7A8dQbI	0.008629	0.001000	0.000168	0.000234	0.000342

5 rows × 261 columns

```
result_asm.head()
```

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:	.eda
0	01kcPWA9K2B0xQeS5Rju	19	744	0	127	57	0	323	
1	1E93CpP60RHFNiT5Qfvn	17	838	0	103	49	0	0	
2	3ekVow2ajZHbTnBcsDfX	17	427	0	50	43	0	145	
3	3X2nY7iQaPBIWDrAZqJe	17	227	0	43	19	0	0	
4	46OZzdsSKDCFV8h7XWxf	17	402	0	59	170	0	0	

5 rows × 54 columns

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

```
(10868, 261)
(10868, 54)
```

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

	Unnamed: 0	0	1	2	3	4	5	6	7
0	0.000000	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946
1	0.000092	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984
2	0.000184	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155
3	0.000276	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481
4	0.000368	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229

5 rows × 308 columns

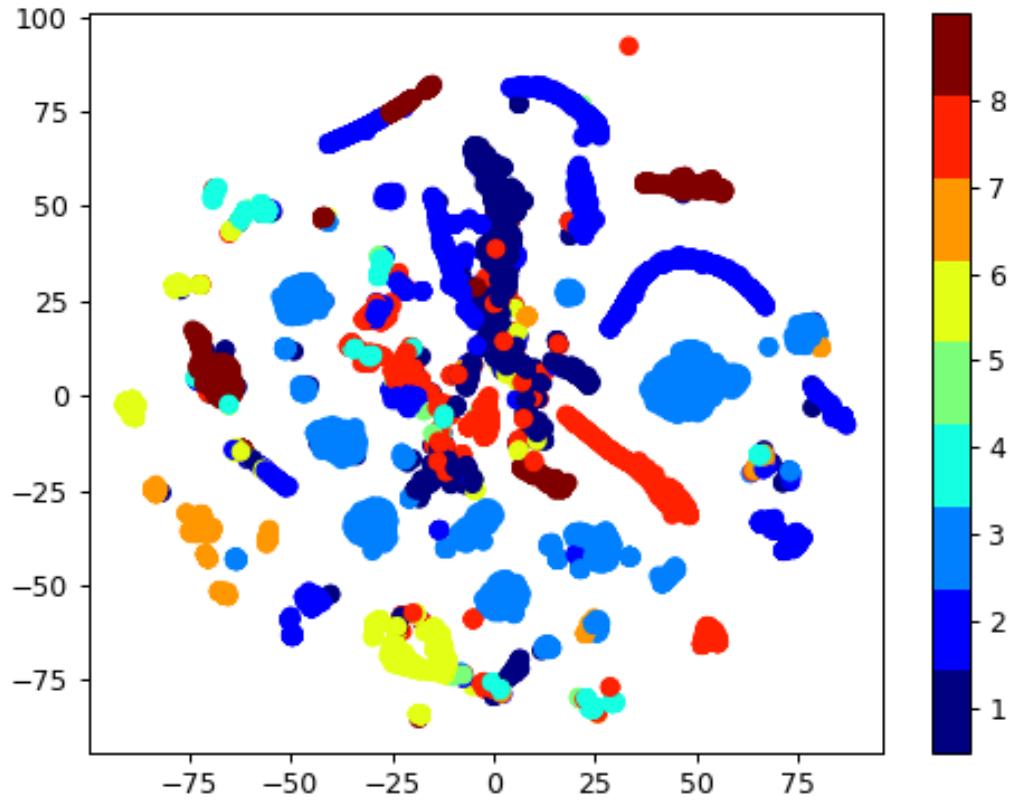
result_y.head()

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

4.5.2. Multivariate Analysis on final features

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





4.5.3. Train and Test split

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

4.5.4. Random Forest Classifier on final features

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

# Some of methods of RandomForestClassifier()
# fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
```

```
# predict(X) Perform classification on samples in X.
# predict_proba (X) Perform classification on samples in X.

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

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


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

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

best_alpha = np.argmin(cv_log_error_array)

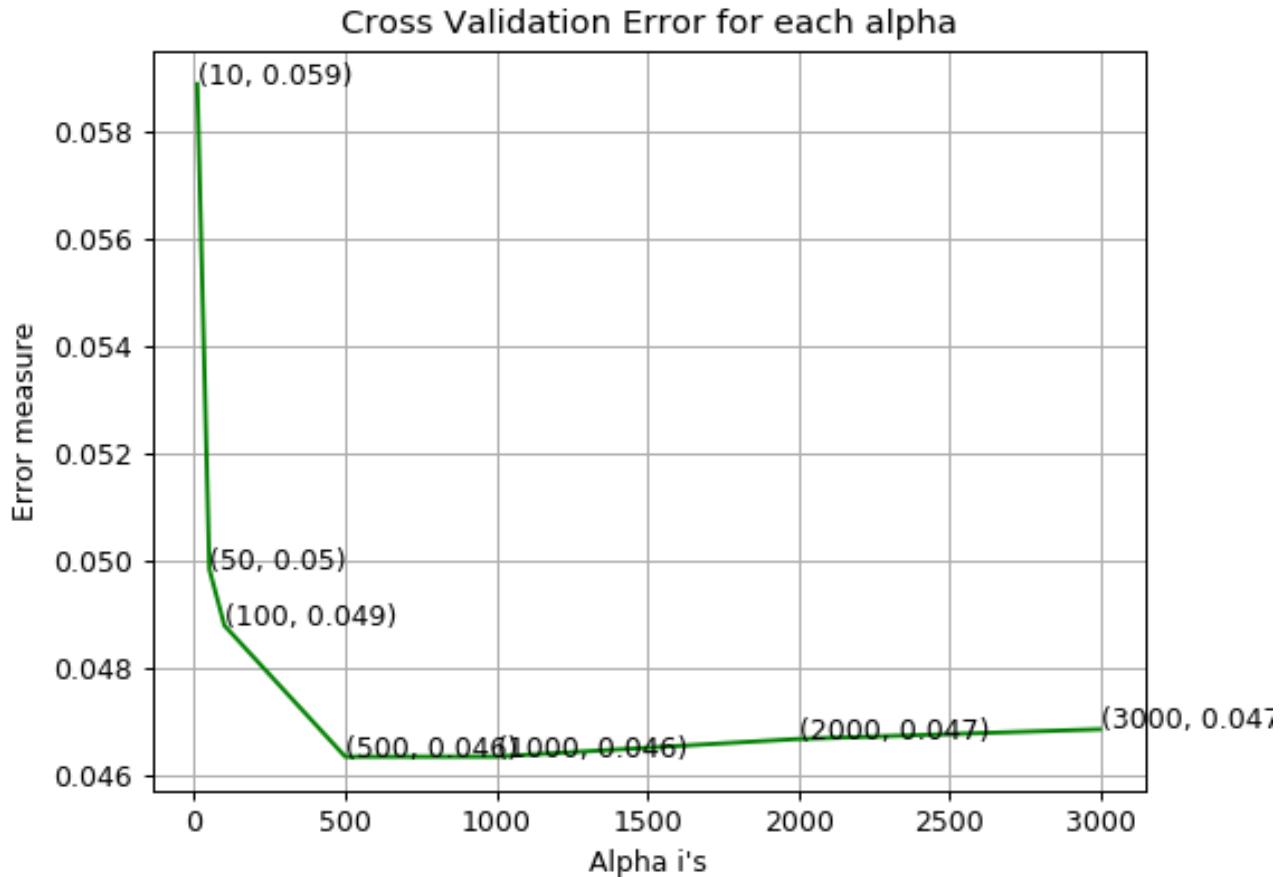
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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



```
log_loss for c = 10 is 0.058864988008900165
log_loss for c = 50 is 0.04982171352389583
log_loss for c = 100 is 0.04877439563993806
log_loss for c = 500 is 0.04633136949419593
log_loss for c = 1000 is 0.04633282669842955
log_loss for c = 2000 is 0.04666148931304081
log_loss for c = 3000 is 0.04684161733430787
```



For values of best alpha = 500 The train log loss is: 0.015045746557915482

For values of best alpha = 500 The cross validation log loss is: 0.04633136949419593

For values of best alpha = 500 The test log loss is: 0.0419437056294099

4.5.5. XgBoost Classifier on final features

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

```
# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/
# -----
# default paramters
# class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
# # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_w
```

```

# max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_la
# scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)

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

```

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, labels=x_cfl.classes_, eps=1e-1))

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

```

best_alpha = np.argmin(cv_log_error_array)

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

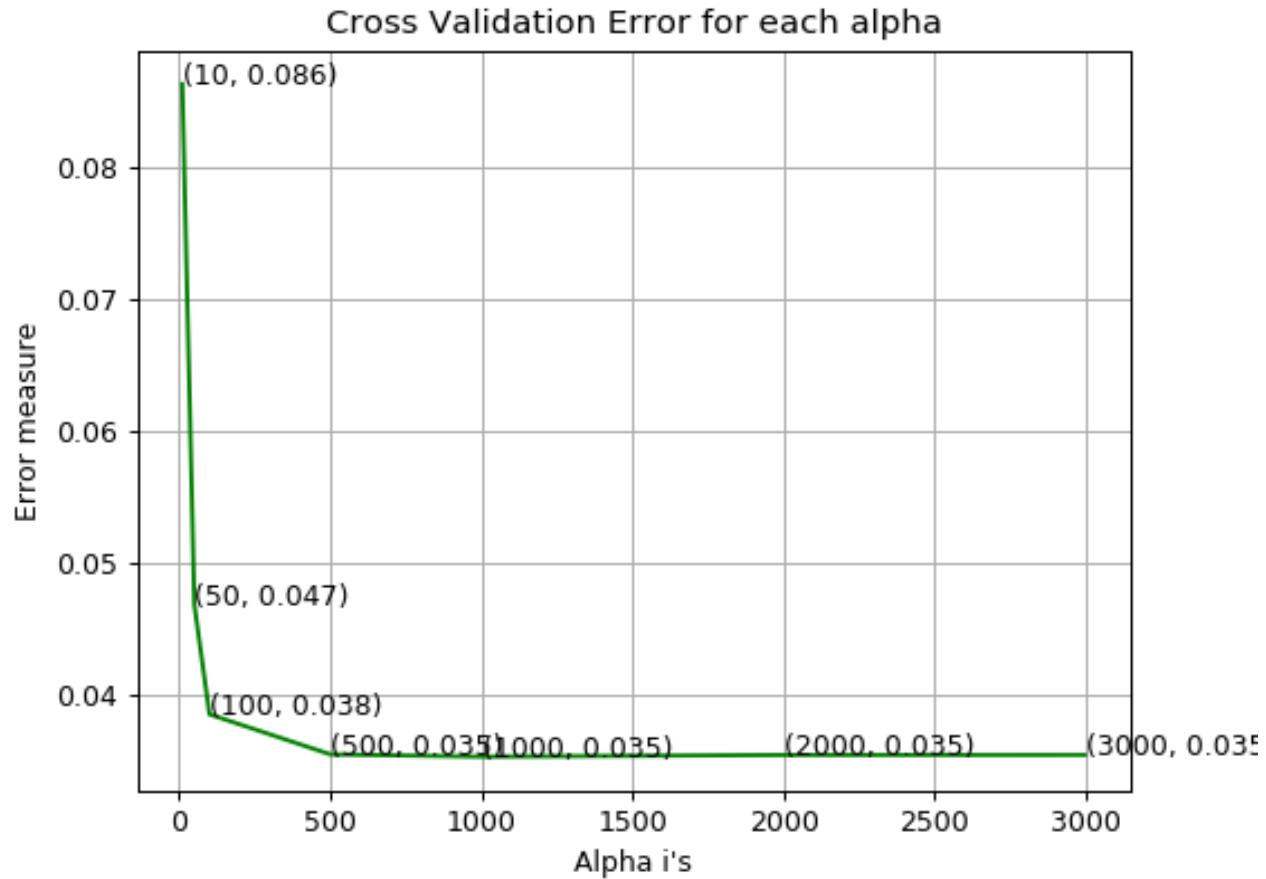
```

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

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



```
log_loss for c = 10 is 0.08634410259197668
log_loss for c = 50 is 0.0467962200270487
log_loss for c = 100 is 0.03846464669244138
log_loss for c = 500 is 0.03542509345482663
log_loss for c = 1000 is 0.03524790113745623
log_loss for c = 2000 is 0.03537820448736872
log_loss for c = 3000 is 0.035384159245550155
```



For values of best alpha = 1000 The train log loss is: 0.010771162453744454

For values of best alpha = 1000 The cross validation log loss is: 0.035384159245550155

For values of best alpha = 1000 The test log loss is: 0.024834218493213808

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

```
x_cfl=XGBClassifier()

prams={
  'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
  'n_estimators':[100,200,500,1000,2000],
  'max_depth':[3,5,10],
  'colsample_bytree':[0.1,0.3,0.5,1],
```

```

'subsample':[0.1,0.3,0.5,1]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,n_jobs=-1, )
random_cfl.fit(X_train_merge, y_train_merge)

❷ Fitting 3 folds for each of 10 candidates, totalling 30 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 5 tasks      | elapsed: 4.6min
[Parallel(n_jobs=-1)]: Done 10 tasks     | elapsed: 11.0min
[Parallel(n_jobs=-1)]: Done 17 tasks     | elapsed: 17.5min
[Parallel(n_jobs=-1)]: Done 27 out of 30 | elapsed: 28.6min remaining: 3.2min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 37.0min finished
RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
                    estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                    colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
                    max_depth=3, min_child_weight=1, missing=None, n_estimators=100,
                    n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
                    reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                    silent=True, subsample=1),
                    fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
                    param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15, 0.2], 'n_e
                    pre_dispatch='2*n_jobs', random_state=None, refit=True,
                    return_train_score='warn', scoring=None, verbose=10)

```

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

❸ {'subsample': 1, 'n_estimators': 200, 'max_depth': 5, 'learning_rate': 0.1, 'colsample_b

```

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

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

```
x_cfl=XGBClassifier(n_estimators=1000,max_depth=5,learning_rate=0.1,colsample_bytree=0.5,subs
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)
```

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



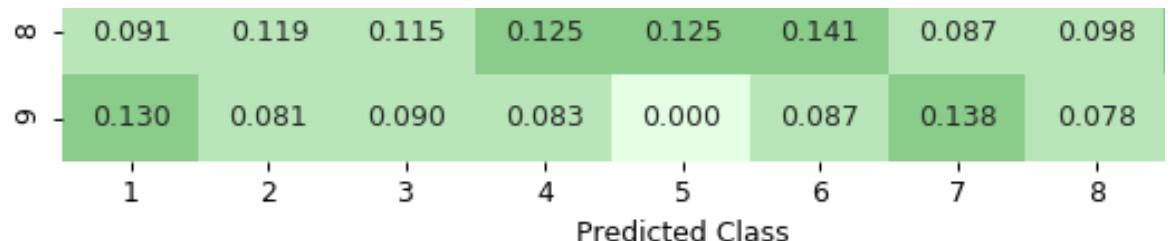
For values of best alpha = 1000 The train log loss is: 0.010849938040281054
 For values of best alpha = 1000 The cross validation log loss is: 0.03269108838914283
 For values of best alpha = 1000 The test log loss is: 0.02814277993749233
 Number of misclassified points 81.73873045078197

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

Original Class	1	2	3	4	5	6	7	8
	1	2	3	4	5	6	7	8
1	44.000	67.000	84.000	14.000	0.000	17.000	7.000	39.000
2	76.000	120.000	133.000	20.000	1.000	27.000	17.000	56.000
3	76.000	134.000	162.000	26.000	2.000	46.000	20.000	69.000
4	15.000	18.000	30.000	10.000	0.000	6.000	5.000	7.000
5	0.000	2.000	3.000	0.000	0.000	0.000	1.000	2.000
6	16.000	43.000	35.000	2.000	3.000	12.000	6.000	20.000
7	13.000	13.000	22.000	4.000	1.000	7.000	6.000	8.000
8	28.000	59.000	68.000	12.000	1.000	21.000	7.000	24.000
9	40.000	40.000	53.000	8.000	0.000	13.000	11.000	19.000

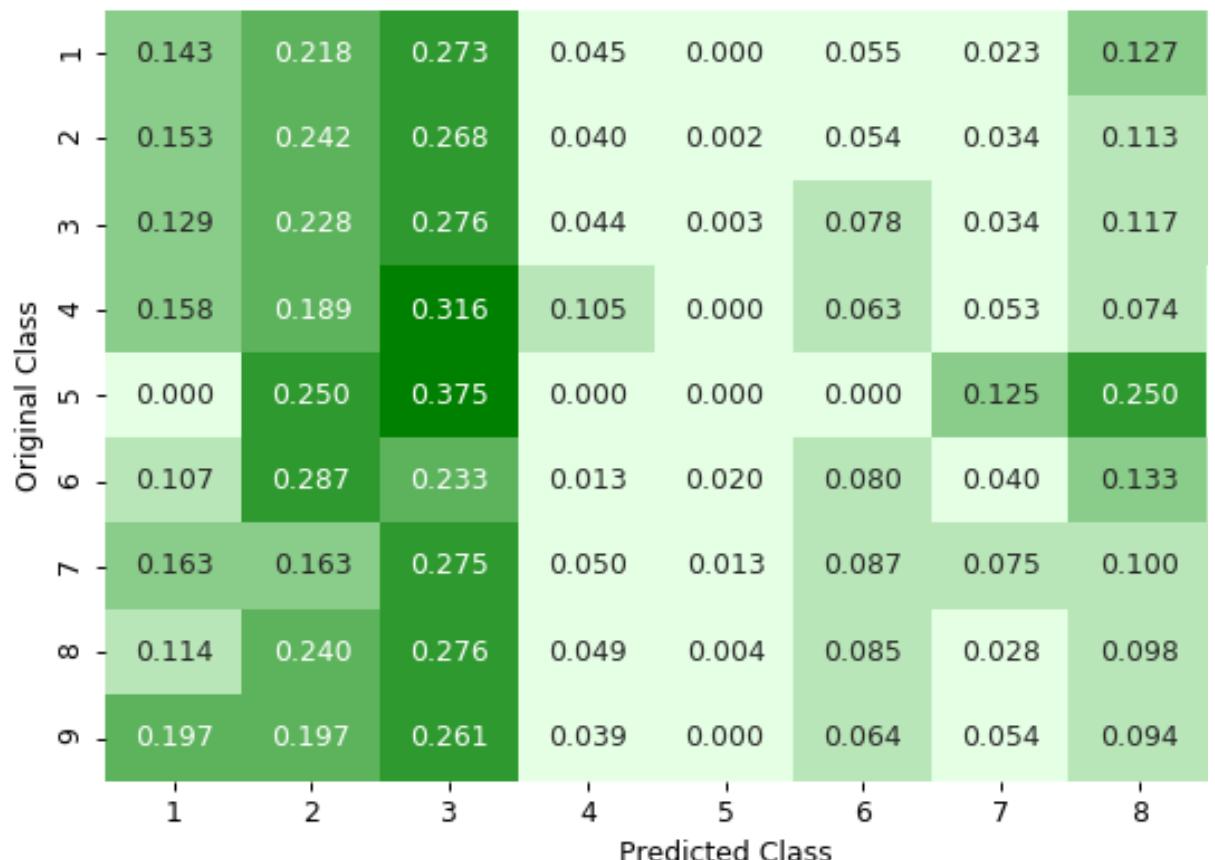
----- Precision matrix -----

Original Class	1	2	3	4	5	6	7	8
	1	2	3	4	5	6	7	8
1	0.143	0.135	0.142	0.146	0.000	0.114	0.087	0.160
2	0.247	0.242	0.225	0.208	0.125	0.181	0.212	0.230
3	0.247	0.270	0.275	0.271	0.250	0.309	0.250	0.283
4	0.049	0.036	0.051	0.104	0.000	0.040	0.062	0.029
5	0.000	0.004	0.005	0.000	0.000	0.000	0.013	0.008
6	0.052	0.087	0.059	0.021	0.375	0.081	0.075	0.082
7	0.042	0.026	0.037	0.042	0.125	0.047	0.075	0.033



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

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



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

▼ byte features

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

```
byte_vocab = "00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,17,18,19,1
```

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

```
byte_bigram()
```

 66049

```
byte_bigram_vocab[:5]
```

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

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

```
byte_trigram()
```

 16974593

```
byte_trigram_vocab[:5]
```

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

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

vector = CountVectorizer(lowercase=False,ngram_range=(2,2), vocabulary=byte_bigram_vocab)
bytebigram_vect = scipy.sparse.csr_matrix((10868, 66049))
for i, file in tqdm(enumerate(os.listdir('byteFiles'))):
    f = open('byteFiles/' + file)
    a[i:]+= scipy.sparse.csr_matrix(vect.fit_transform([f.read().replace('\n', ' ').lower()]))
    f.close()
```

 10868it [3:49:23, 2.10it/s]

```
bytebigram_vect
```

 <10868x66049 sparse matrix of type '<class 'numpy.float64'>' with 0 stored elements in Compressed Sparse Row format>

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

from sklearn.preprocessing import normalize
byte_bigram_vect = normalize(scipy.sparse.load_npz('bytebigram.npz'), axis = 0)
```

▼ N-Gram(2-Gram, 3-Gram, 4-Gram) Opcode Vectorization

```
opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc', 'dec', 'a
```

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

```
asmopcodebigram
```

 676

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

```
asmopcodetrigram
```

 17576

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

```
asmopcodetetragram
```

 456976

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

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

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

```
opcodebivect
```

```
👤 <10868x676 sparse matrix of type '<class 'numpy.float64'>'  
      with 1877309 stored elements in Compressed Sparse Row format>
```

```
scipy.sparse.save_npz('opcodebigram.npz', opcodebivect)
```

```
vect = CountVectorizer(ngram_range=(3, 3), vocabulary = asmopcodetrigram)
opcodetrivect = scipy.sparse.csr_matrix((10868, len(asmopcodetrigram)))
```

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

```
opcodetrivect
```

```
👤 <10868x17576 sparse matrix of type '<class 'numpy.float64'>'  
      with 7332672 stored elements in Compressed Sparse Row format>
```

```
scipy.sparse.save_npz('opcodetrigram.npz', opcodetrivect)
```

```
vect = CountVectorizer(ngram_range=(4, 4), vocabulary = asmopcodetetragram)
opcodetetraVect = scipy.sparse.csr_matrix((10868, len(asmopcodetetragram)))
```

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

```
opcodetetraVect
```

```
👤 <10868x456976 sparse matrix of type '<class 'numpy.float64'>'  
with 16605229 stored elements in Compressed Sparse Row format>
```

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

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

```
opcodetriVect=scipy.sparse.load_npz('opcodetrigram.npz')
```

```
opcodebiVect=scipy.sparse.load_npz('opcodebigram.npz')
```

▼ Image Feature Extraction From ASM Files

```
import array

def collect_img_asm():
    for asmfile in os.listdir("asmFiles"):
        filename = asmfile.split('.')[0]
        file = codecs.open("asmFiles/" + asmfile, 'rb')
        filelen = os.path.getsize("asmFiles/" + asmfile)
        width = int(filelen ** 0.5)
        rem = int(filelen / width)
        arr = array.array('B')
        arr.frombytes(file.read())
        file.close()
        reshaped = np.reshape(arr[:width * width], (width, width))
        reshaped = np.uint8(reshaped)
        scipy.misc.imsave('asm_image/' + filename + '.png', reshaped)

collect_img_asm()

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

▼ First 200 Image Pixels

```
import cv2
imagefeatures = np.zeros((10868, 200))

for i, asmfile in enumerate(os.listdir("asmFiles")):
    img = cv2.imread("asm_image/" + asmfile.split('.')[0] + '.png')
    img_arr = img.flatten()[:200]
```

```

imagefeatures[i, :] += img_arr

imgfeatures_name = []
for i in range(200):
    img_features_name.append('pix' + str(i))
imgdf = pd.DataFrame(normalize(imagefeatures, axis = 0), columns = imgfeatures_name)

imgdf['ID'] = result.ID

imgdf.head()

```

👤

	pix0	pix1	pix2	pix3	pix4	pix5	pix6	pix7	pix8
0	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320
1	0.006560	0.006560	0.006560	0.013504	0.013504	0.013504	0.012927	0.012927	0.012927
2	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320
3	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320
4	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320

5 rows × 201 columns

```
joblib.dump(imgdf, 'img_df')
```

👤 ['img_df']

```
img_df=joblib.load('img_df')
```

```
img_df.head()
```

👤

	pix0	pix1	pix2	pix3	pix4	pix5	pix6	pix7	pix8
0	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320
1	0.006560	0.006560	0.006560	0.013504	0.013504	0.013504	0.012927	0.012927	0.012927
2	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320
3	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320
4	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.008320	0.008320

5 rows × 201 columns

▼ Important Feature Selection Using Random Forest

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

► Important Feature Among Opcode Bi-Gram

↳ 5 cells hidden

▼ Important Feature Among Opcode 3-Gram

```
op_tri_indexes = imp_features(normalize(opcodetrivect, axis = 0), asmopcodetrigram, 200)  
  
op_tri_df = pd.SparseDataFrame(normalize(opcodetrivect, axis = 0), columns = asmopcodetrigram  
op_tri_df = op_tri_df.loc[:, np.intersect1d(op_tri_df.columns, np.take(asmopcodetrigram, op_t  
  
op_tri_df.to_dense().to_csv('op_tri.csv')  
  
op_tri_df = pd.read_csv('op_tri.csv').drop('Unnamed: 0', axis = 1).fillna(0)  
  
op_tri_df['ID'] = result.ID  
op_tri_df.head()
```



	add	cmp	jmp	add	mov	add	mov	cmp	add	mov	jmp	add	mov	mov	add	pop	call	add	pop	mov	add	pop	pop	add	pop	push
0	0.000000			0.002183		0.001340		0.001563		0.003593		0.0		0.005354		0.000342		0.000000		0.0						
1	0.000000			0.001364		0.000670		0.000625		0.002705		0.0		0.001785		0.000000		0.000000		0.000000		0.0				
2	0.000000			0.000000		0.000000		0.000000		0.000000		0.0		0.000000		0.000000		0.000000		0.000000		0.0				
3	0.000000			0.000000		0.000000		0.000000		0.000000		0.0		0.000000		0.000000		0.000000		0.000000		0.0				
4	0.001292			0.001091		0.004914		0.002814		0.014009		0.0		0.000000		0.000000		0.000441		0.0						

5 rows x 201 columns

▼ Important Feature Among Opcode 4-Gram

```
op_tetra_indexes = imp_features(normalize(opcodetetraVect, axis = 0), asmopcodetetragram, 200)

op_tetra_df = pd.SparseDataFrame(normalize(opcodetetraVect, axis = 0), columns = asmopcodetet
op_tetra_df = op_tetra_df.loc[:, np.intersect1d(op_tetra_df.columns, np.take(asmopcodetetrag
op_tetra_df = op_tetra_df.to_dense().to_csv('op_tetra.csv')

op_tetra_df = pd.read_csv('op_tetra.csv').drop('Unnamed: 0', axis = 1).fillna(0)

op_tetra_df['ID'] = result.ID
op_tetra_df.head()
```



	add	mov	add	mov	add	mov	add	mov	add	pop	add	pop	add	pop	add	retn	call	add	mov	..
	add	mov	add	pop	cmp	jnb	mov	add	mov	mov	mov	push	pop	pop	push	push	push	sub		
0	0.001593	0.007668	0.000000	0.002031	0.002517		0.0	0.0	0.0	0.0	0.00116	0.000000							.	
1	0.000000	0.007668	0.000000	0.001625	0.002760		0.0	0.0	0.0	0.0	0.000000	0.000000							.	
2	0.000000	0.000000	0.000000	0.000000	0.000000		0.0	0.0	0.0	0.0	0.000000	0.000000							.	
3	0.000000	0.000000	0.000000	0.000000	0.000000		0.0	0.0	0.0	0.0	0.000000	0.000000							.	
4	0.002125	0.000000	0.023352	0.023558	0.006657		0.0	0.0	0.0	0.0	0.000000	0.009682							.	

5 rows × 201 columns

▼ Important Feature Among Byte Bi-Gram

```
byte_bi_idxes = imp_features(normalize(bytebigram_vect, axis = 0), byte_bigram_vocab, 300)

np.save('byte_bi_idx', byte_bi_idxes)

byte_bi_idxes = np.load('byte_bi_idx.npy')

top_byte_bi = np.zeros((10868, 0))
for i in byte_bi_idxes:
    sliced = bytebigram_vect[:, i].todense()
    top_byte_bi = np.hstack([top_byte_bi, sliced])

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

byte_bi_df.to_dense().to_csv('byte_bi.csv')

byte_bi_df = pd.read_csv('byte_bi.csv').drop('Unnamed: 0', axis = 1).fillna(0)

byte_bi_df['ID'] = result.ID

byte_bi_df.head()
```



	??	55 ??	55 95	55 b3	55 b2	55 b1	55 b0	55 af	55 ae	55 ad	55 ac	...	54 b3	54 b4
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0

5 rows × 301 columns

▼ Advanced features

▼ Adding 300 bytebigram,200 opcode bigram,200 opcode trigram,200 or image pixels

```
final_data = pd.concat([result_x, op_bi_df, op_tri_df, op_tetra_df, byte_bi_df, img_df], axis = 1)

final_data = final_data.drop('ID', axis = 1)

final_data.head()
```



	Unnamed: 0	0	1	2	3	4	5	6	7
0	0.000000	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946
1	0.000092	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984
2	0.000184	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155
3	0.000276	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481
4	0.000368	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229

5 rows × 1408 columns

```
final_data.to_csv('final_data.csv')
```

```
final_data = pd.read_csv('final_data.csv')

x_train_final, x_test_final, y_train_final, y_test_final = train_test_split(final_data, resul
x_trn_final, x_cv_final, y_trn_final, y_cv_final = train_test_split(x_train_final, y_train_fi
```

▼ Machine Learning Models on ASM Features + Byte Features +

```
alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(x_trn_final,y_trn_final)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(x_trn_final,y_trn_final)
    predict_y = sig_clf.predict_proba(x_cv_final)
    cv_log_error_array.append(log_loss(y_cv_final, predict_y, labels=logisticR.classes_, eps=1e-9))

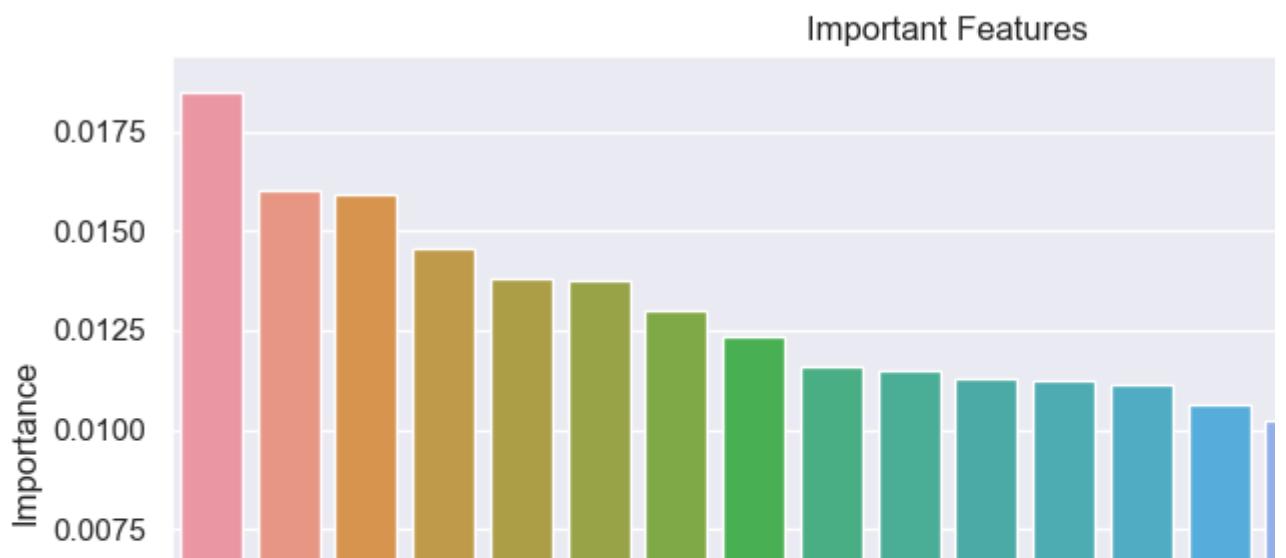
for i in range(len(cv_log_error_array)):
    print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])

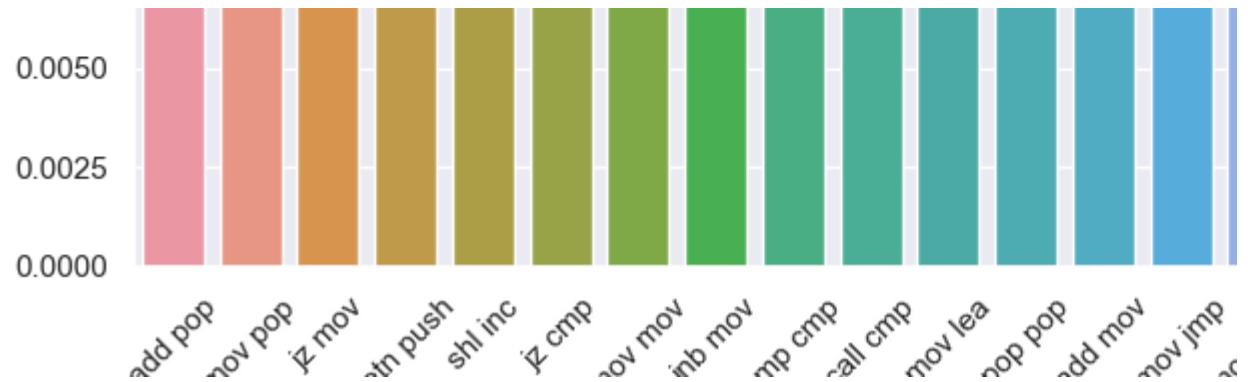
best_alpha = np.argmin(cv_log_error_array)

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

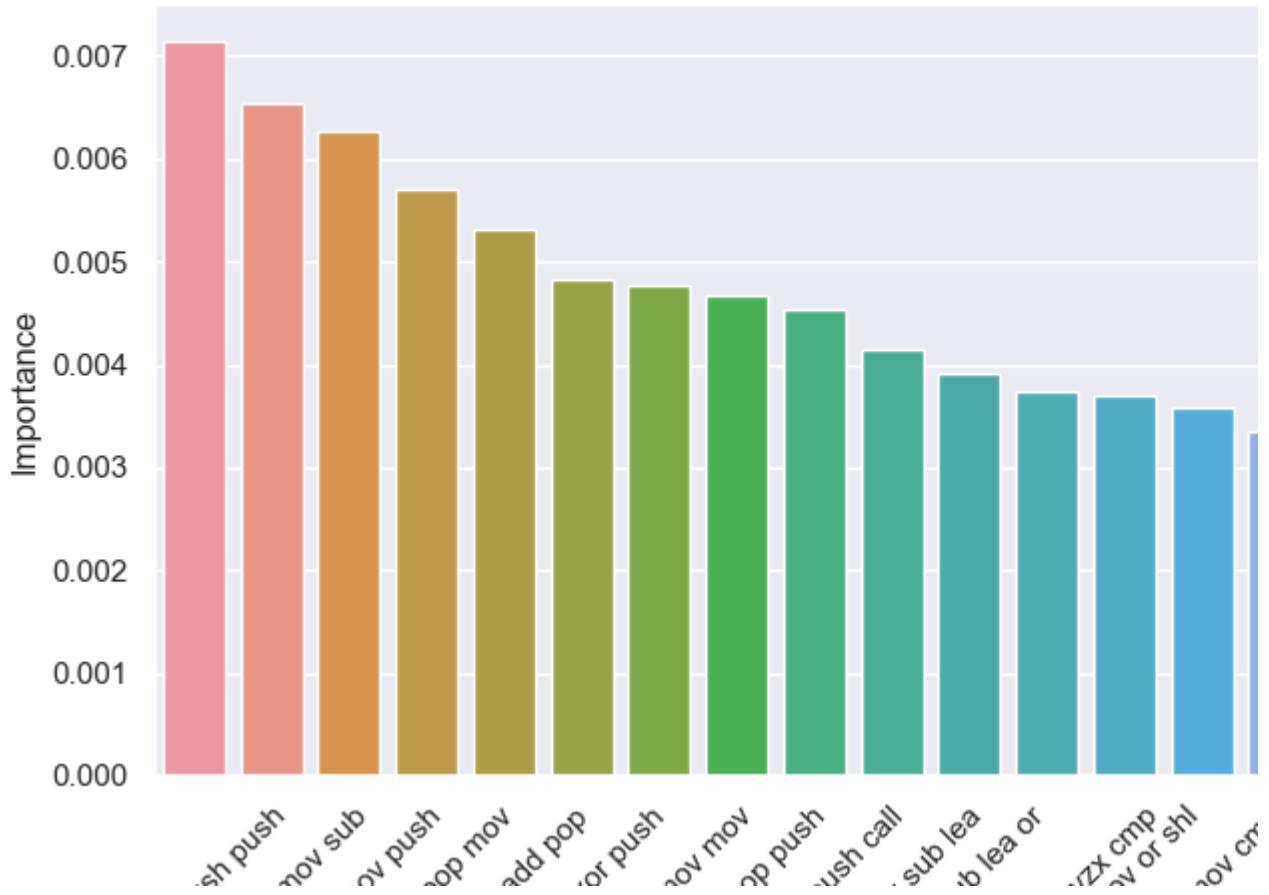


```
log_loss for c = 1e-05 is 1.1840039867727614
log_loss for c = 0.0001 is 1.1217881098745714
log_loss for c = 0.001 is 1.174936322460997
log_loss for c = 0.01 is 1.0741224260174453
log_loss for c = 0.1 is 1.1761396828975654
log_loss for c = 1 is 1.2362570810343723
log_loss for c = 10 is 1.1804717850739066
log_loss for c = 100 is 1.1684083137157295
log_loss for c = 1000 is 1.1061521197568476
```

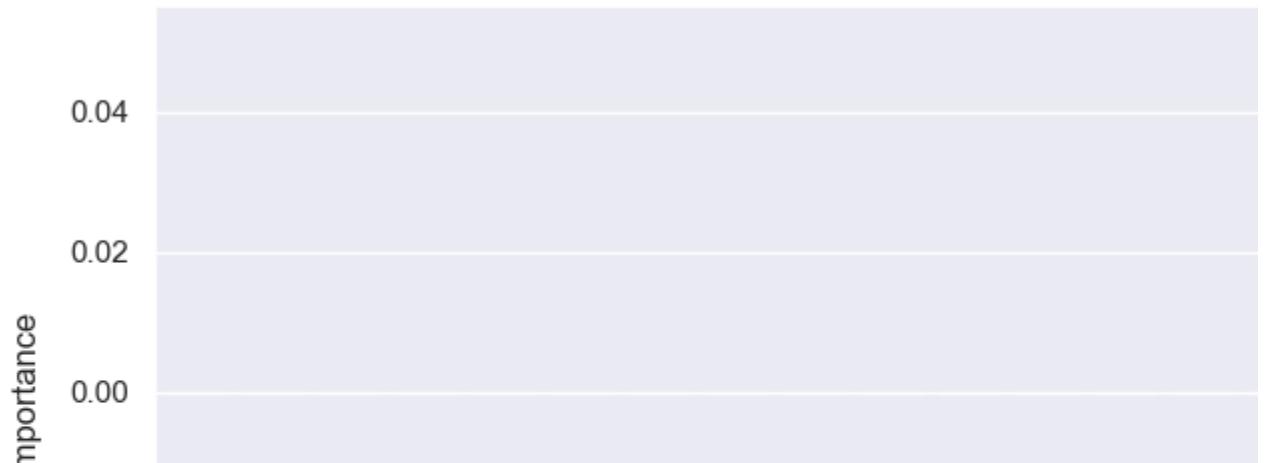


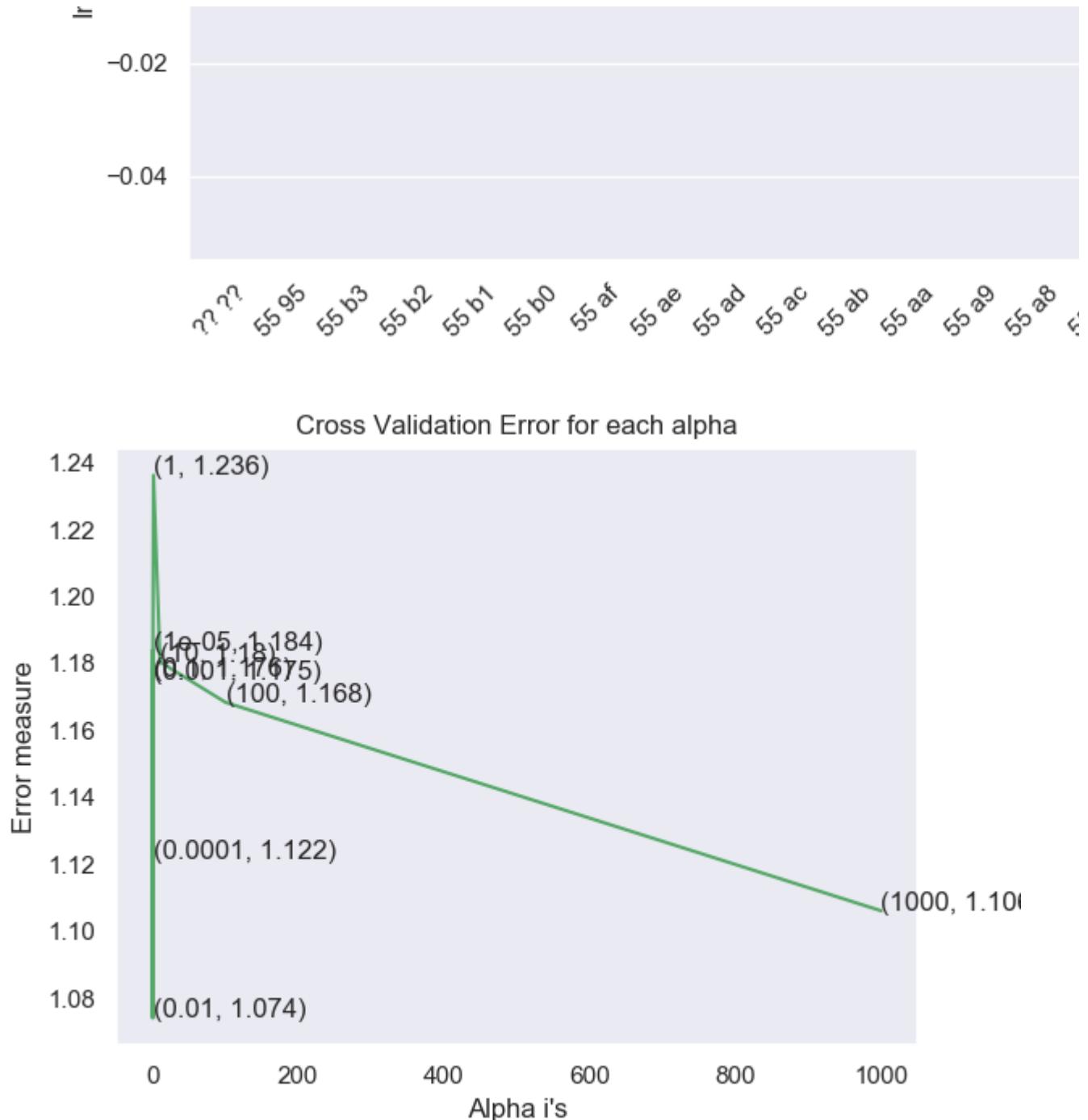


Important Features



Important Features





```

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

predict_y = sig_clf.predict_proba(x_trn_final)
print ('log loss for train data',(log_loss(y_trn_final, predict_y, labels=logisticR.classes_))
predict_y = sig_clf.predict_proba(x_cv_final)
  
```

```
print ('log loss for cv data',(log_loss(y_cv_final, predict_y, labels=logisticR.classes_, eps
predict_y = sig_clf.predict_proba(x_test_final)
print ('log loss for test data',(log_loss(y_test_final, predict_y, labels=logisticR.classes_,
```

👤 C:\Users\Sai charan\Anaconda3\lib\site-packages\sklearn\svm\base.py:922: ConvergenceWarning
"the number of iterations.", ConvergenceWarning)
log loss for train data 1.193174530266704
log loss for cv data 1.1785070578048291
log loss for test data 1.2060464393477006

```
plot_confusion_matrix(y_test_final,sig_clf.predict(x_test_final))
```



Number of misclassified points 31.324747010119598

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

Original Class	1	2	3	4	5	6	7	8
	1	2	3	4	5	6	7	8
1	252.000	3.000	50.000	0.000	0.000	0.000	0.000	3.000
2	7.000	471.000	17.000	0.000	0.000	0.000	0.000	1.000
3	0.000	6.000	581.000	0.000	0.000	0.000	0.000	1.000
4	0.000	7.000	87.000	0.000	0.000	0.000	0.000	1.000
5	0.000	0.000	7.000	0.000	0.000	0.000	0.000	1.000
6	4.000	0.000	145.000	0.000	0.000	0.000	0.000	1.000
7	2.000	5.000	69.000	0.000	0.000	0.000	4.000	0.000
8	15.000	5.000	41.000	0.000	0.000	0.000	0.000	185.000
9	3.000	14.000	185.000	0.000	0.000	0.000	0.000	1.000

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

Original Class	1	2	3	4	5	6	7	8
	1	2	3	4	5	6	7	8
1	0.890	0.006	0.042					0.000 0.015
2	0.025	0.922	0.014					0.000 0.005
3	0.000	0.012	0.492					0.000 0.005
4	0.000	0.014	0.074					0.000 0.005
5	0.000	0.000	0.006					0.000 0.005
6	0.014	0.000	0.123					0.000 0.005
7	0.007	0.010	0.058				1.000 0.000	
8	0.053	0.010	0.035				0.000 0.954	

9 -	0.011	0.027	0.157			0.000	0.005	
	1	2	3	4	5	6	7	8

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

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

Original Class	1 -	0.818	0.010	0.162	0.000	0.000	0.000	0.000	0.010
	2 -	0.014	0.950	0.034	0.000	0.000	0.000	0.000	0.002
	3 -	0.000	0.010	0.988	0.000	0.000	0.000	0.000	0.002
	4 -	0.000	0.074	0.916	0.000	0.000	0.000	0.000	0.011
	5 -	0.000	0.000	0.875	0.000	0.000	0.000	0.000	0.125
	6 -	0.027	0.000	0.967	0.000	0.000	0.000	0.000	0.007
	7 -	0.025	0.062	0.863	0.000	0.000	0.000	0.050	0.000
	8 -	0.061	0.020	0.167	0.000	0.000	0.000	0.000	0.752
	9 -	0.015	0.069	0.911	0.000	0.000	0.000	0.000	0.005
	1	2	3	4	5	6	7	8	

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

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

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

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

```
alpha=[10,100,1000,2000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=i)
    x_cfl.fit(x_trn_final,y_trn_final)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(x_trn_final, y_trn_final)
    predict_y = sig_clf.predict_proba(x_cv_final)
    cv_log_error_array.append(log_loss(y_cv_final, predict_y, labels=x_cfl.classes_, eps=1e-1))

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

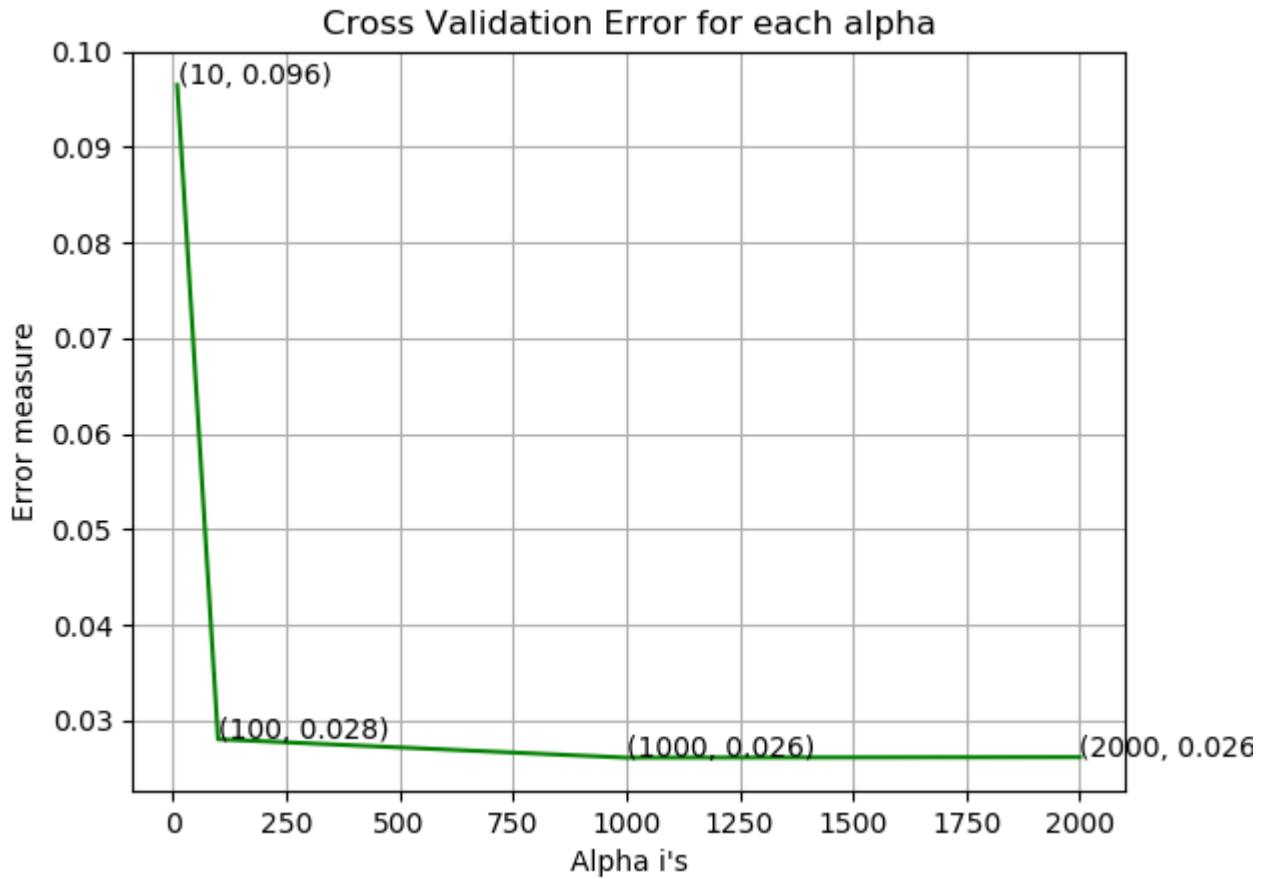


```
best_alpha = np.argmin(cv_log_error_array)

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



```
log_loss for c = 10 is 0.09649648467635132
log_loss for c = 100 is 0.028026994875892948
log_loss for c = 1000 is 0.02610102301724636
log_loss for c = 2000 is 0.026155764643162237
```



```
x_cfl=XGBClassifier(n_estimators=2000,nthread=-1)
x_cfl.fit(x_trn_final,y_trn_final,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(x_trn_final, y_trn_final)

predict_y = sig_clf.predict_proba(x_trn_final)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_
predict_y = sig_clf.predict_proba(x_cv_final)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",l
predict_y = sig_clf.predict_proba(x_test_final)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_te
```

- 👤 For values of best alpha = 0.01 The train log loss is: 0.010187974436441512
- 👤 For values of best alpha = 0.01 The cross validation log loss is: 0.02395762856614576
- 👤 For values of best alpha = 0.01 The test log loss is: 0.018309505637434106

▼ Procedure:

1. First I took the byte file and made Exploratory Data Analysis.
2. Used uni-gram count features and applied machine learning models.
3. Preprocessed the asm file and extracted various segment count as features.
4. Applied machine learning models on asm segment count.
5. Combined byte features and asm segment features.
6. Applied machine learning models on combined features.
7. Extracted features like byte bigram,opcode bi gram,opcode trigram , opcode tetra gram and 200 pixel features.
8. Applied machine learning models on combined features.

▼ Results

```
from prettytable import PrettyTable
phtable = PrettyTable()
phtable.title = " Model Comparision "
phtable.field_names = ["Model", 'Features', 'log loss']
phtable.add_row(["random", "Byte files", "2.45"])
phtable.add_row(["knn", "Byte files", "0.48"])
phtable.add_row(["Logistic Regression", "Byte files", "0.52"])
phtable.add_row(["Random Forest Classifier ", "Byte files", "0.06"])
phtable.add_row(["XgBoost Classification", "Byte files", "0.07"])
phtable.add_row(["\n", "\n", "\n"])
phtable.add_row(["knn", "asmfiles", "0.21"])
phtable.add_row(["Logistic Regression", "asmfiles", "0.38"])
phtable.add_row(["Random Forest Classifier ", "asmfiles", "0.03"])
phtable.add_row(["XgBoost Classification", "asmfiles", "0.04"])
phtable.add_row(["\n", "\n", "\n"])
phtable.add_row(["Random Forest Classifier ", "Byte files+asmfiles", "0.04"])
phtable.add_row(["XgBoost Classification", "Byte files+asmfiles", "0.02"])
phtable.add_row(["\n", "\n", "\n"])
phtable.add_row(["Logistic Regression", "Byte files+asmfiles+advanced features", "1.12"])
phtable.add_row(["XgBoost Classification", "Byte files+asmfiles+advanced features", "0.01"])
print(phtable)
```



Model	Features	log loss
random	Byte files	2.45
knn	Byte files	0.48
Logistic Regression	Byte files	0.52
Random Forest Classifier	Byte files	0.06
XgBoost Classification	Byte files	0.07
<hr/>		
knn	asmfiles	0.21
Logistic Regression	asmfiles	0.38
Random Forest Classifier	asmfiles	0.03
XgBoost Classification	asmfiles	0.04
<hr/>		
Random Forest Classifier	Byte files+asmfiles	0.04
XgBoost Classification	Byte files+asmfiles	0.02
<hr/>		
Logistic Regression	Byte files+asmfiles+advanced features	1.12
XgBoost Classification	Byte files+asmfiles+advanced features	0.01