

Microsoft Malware detection

In [1]:

```
!pip install xgboost
```

```
Collecting xgboost
  Downloading xgboost-1.2.0-py3-none-manylinux2010_x86_64.whl (148.9 MB)
    |██████████| 148.9 MB 28 kB/s eta 0:00:01
Requirement already satisfied: numpy in /opt/conda/lib/python3.7/site-packages (from xgboost) (1.19.1)
Requirement already satisfied: scipy in /opt/conda/lib/python3.7/site-packages (from xgboost) (1.5.1)
Installing collected packages: xgboost
Successfully installed xgboost-1.2.0
```

1.Business/Real-world Problem

1.1. What is Malware?

In [1]:

```
!pip install py7zr
```

```
Requirement already satisfied: py7zr in /opt/conda/lib/python3.7/site-packages (0.9.7)
Requirement already satisfied: importlib-metadata; python_version < "3.8" in /opt/conda/lib/python3.7/site-packages (from py7zr) (1.7.0)
Requirement already satisfied: texttable in /opt/conda/lib/python3.7/site-packages (from py7zr) (1.6.3)
Requirement already satisfied: pycryptodome in /opt/conda/lib/python3.7/site-packages (from py7zr) (3.9.8)
Requirement already satisfied: zipp>=0.5 in /opt/conda/lib/python3.7/site-packages (from importlib-metadata; python_version < "3.8"->py7zr) (3.1.0)
```

The term malware is a contraction of malicious software. Put simply, malware is any piece of software that was written with the intent of doing harm to data, devices or to people.

Source: <https://www.avg.com/en/signal/what-is-malware> (<https://www.avg.com/en/signal/what-is-malware>)

In [2]:

```
!python -m py7zr x train.7z
!py7zr x train.7z
```

In [4]:

```
print('su')
```

```
su
```

1.2. Problem Statement

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

1.3 Source/Useful Links

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

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

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

1.4. Real-world/Business objectives and constraints.

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

2. Machine Learning Problem

2.1. Data

2.1.1. Data Overview

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

2.1.2. Example Data Point

.asm file

```

.text:00401000          assume es:nothing, ss:nothing
g, ds:_data, fs:nothing, gs:nothing
.text:00401000 56        push    esi
                        lea     eax, [esp+8]
.text:00401001 8D 44 24 08      push    eax
.text:00401005 50        mov     esi, ecx
.text:00401006 8B F1        call    ??0exception@_
                           std@@QAE@ABQBD@Z ; std::exception::exception(char const * const &)
.text:0040100D C7 06 08 BB 42 00      mov     dword ptr [es
i], offset off_42BB08
.text:00401013 8B C6        mov     eax, esi
.text:00401015 5E        pop    esi
.text:00401016 C2 04 00        retn    4
.text:00401016          ; -----
-----
.text:00401019 CC CC CC CC CC CC CC CC align 10h
.text:00401020 C7 01 08 BB 42 00      mov     dword ptr [ec
x], offset off_42BB08
.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 [es
i], offset off_42BB08
.text:00401039 E8 13 1C 00 00        call    sub_402C51
.text:0040103E F6 44 24 08 01      test   byte ptr [esp
+8], 1
.text:00401043 74 09        jz     short loc_40104E
.text:00401045 56        push    esi
.text:00401046 E8 6C 1E 00 00        call    ??3@YAXPAX@Z
; operator delete(void *)
.text:0040104B 83 C4 04        add    esp, 4
.text:0040104E          loc_40104E:           ; CODE XR
EF: .text:00401043;j
.text:0040104E 8B C6        mov     eax, esi
.text:00401050 5E        pop    esi
.text:00401051 C2 04 00        retn    4
.text:00401051          ; -----
-----
```

.bytes file

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

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

2.2.1. Type of Machine Learning Problem

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

2.2.2. Performance Metric

Source: [\(https://www.kaggle.com/c/malware-classification#evaluation\)](https://www.kaggle.com/c/malware-classification#evaluation)

Metric(s):

- Multi class log-loss
- Confusion matrix

2.2.3. Machine Learing Objectives and Constraints

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

Constraints:

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

2.3. Train and Test Dataset

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

2.4. Useful blogs, videos and reference papers

<http://blog.kaggle.com/2015/05/26/microsoft-malware-winners-interview-1st-place-no-to-overfitting/>
<https://arxiv.org/pdf/1511.04317.pdf>
First place solution in Kaggle competition: <https://www.youtube.com/watch?v=VLQTRILGz5Y>
<https://github.com/dchad/malware-detection>
<http://vizsec.org/files/2011/Nataraj.pdf>
https://www.dropbox.com/sh/gfqzv0ckgs4l1bf/AAB6EeInEjvvuQg2nu_pIB6ua?dl=0
" Cross validation is more trustworthy than domain knowledge."

3. Exploratory Data Analysis

In [1]:

```
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
%matplotlib inline
```

In [2]:

```
#separating byte files and asm files

source = 'train'
destination_1 = 'byteFiles'
destination_2 = 'asmFiles'

# we will check if the folder 'byteFiles' exists if it not there we will create a folder with the same name
if not os.path.isdir(destination_1):
    os.makedirs(destination_1)
if not os.path.isdir(destination_2):
    os.makedirs(destination_2)

# if we have folder called 'train' (train folder contains both .asm files and .bytes files) we will rename it 'asmFiles'
# for every file that we have in our 'asmFiles' directory we check if it is ending with .bytes, if yes we will move it to
# 'byteFiles' folder

# so by the end of this snippet we will separate all the .byte files and .asm files
if os.path.isdir(source):
    data_files = os.listdir(source)
    for file in data_files:
        print(file)
        if (file.endswith("bytes")):
            shutil.move(source+'/'+file,destination_1)
        if (file.endswith("asm")):
            shutil.move(source+'/'+file,destination_2)
```

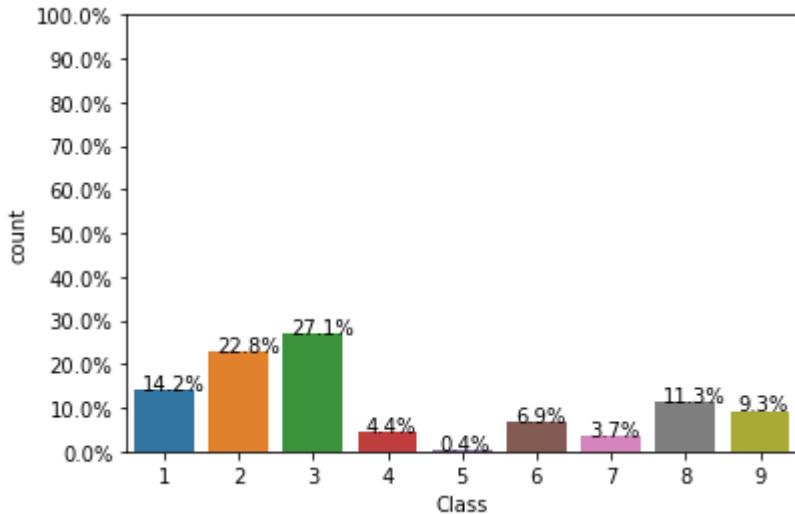
3.1. Distribution of malware classes in whole data set

In [10]:

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

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

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



3.2. Feature extraction

3.2.1 File size of byte files as a feature

In [11]:

```
#file sizes of byte files

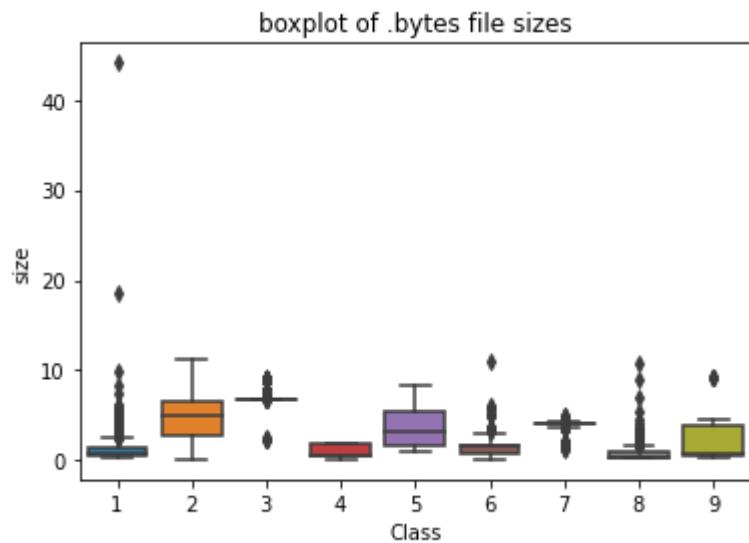
files=os.listdir('byteFiles')
filenames=Y['Id'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, st_nlink=1, st_uid=0,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

	ID	size	Class
0	9MW5Nuf0ogCEcR1YJeKG	1.113281	8
1	GFblksovaPYLI5XeC8RU	3.808594	9
2	kjQFMnf7vADCqtX4ai2u	1.628906	6
3	EKmgShQsf6a9vY0znNlU	0.691406	4
4	1WCXg8qEtJFITMilaf6k	0.539062	6

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

In [12]:

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



3.2.3 feature extraction from byte files

In [6]:

```
#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+')
        with open('byteFiles/'+file+'.bytes',"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+'.bytes')
        text_file.close()

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

In [0]:

```
#program to convert into bag of words of bytefiles
#this is custom-built bag of words this is unigram bag of words
byte_feature_file=open('result.csv','w+')
byte_feature_file.write("ID,0,1,2,3,4,5,6,7,8,9,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,
17,18,19,1a,1b,1c,1d,1e,1f,20,21,22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,
34,35,36,37,38,39,3a,3b,3c,3d,3e,3f,40,41,42,43,44,45,46,47,48,49,4a,4b,4c,4d,4e,4f,50,
51,52,53,54,55,56,57,58,59,5a,5b,5c,5d,5e,5f,60,61,62,63,64,65,66,67,68,69,6a,6b,6c,6d,
6e,6f,70,71,72,73,74,75,76,77,78,79,7a,7b,7c,7d,7e,7f,80,81,82,83,84,85,86,87,88,89,8a,
8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,9a,9b,9c,9d,9e,9f,a0,a1,a2,a3,a4,a5,a6,a7,
a8,a9,aa,ab,ac,ad,ae,af,b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,ba,bb,bc,bd,be,bf,c0,c1,c2,c3,c4,
c5,c6,c7,c8,c9,ca,cb,cc,cd,ce,cf,d0,d1,d2,d3,d4,d5,d6,d7,d8,d9,da,db,dc,dd,de,df,e0,e1,
e2,e3,e4,e5,e6,e7,e8,e9,ea,eb,ec,ed,ee,ef,f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,
ff,??")
byte_feature_file.write("\n")
for file in files:
    filenames2.append(file)
    byte_feature_file.write(file+",")
    if(file.endswith("txt")):
        with open('byteFiles/'+file,"r") as byte_flie:
            for lines in byte_flie:
                line=lines.rstrip().split(" ")
                for hex_code in line:
                    if hex_code=='??':
                        feature_matrix[k][256]+=1
                    else:
                        feature_matrix[k][int(hex_code,16)]+=1
            byte_flie.close()
    for i, row in enumerate(feature_matrix[k]):
        if i!=len(feature_matrix[k])-1:
            byte_feature_file.write(str(row)+",")
        else:
            byte_feature_file.write(str(row))
    byte_feature_file.write("\n")
    k += 1
byte_feature_file.close()
```

In [13]:

```
byte_features=pd.read_csv("result.csv")
byte_features['ID'] = byte_features['ID'].str.split('.').str[0]
byte_features.head(2)
```

Out[13]:

	ID	0	1	2	3	4	5	6	7	8	...
0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...
1	01lsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	4

2 rows × 258 columns



In [14]:

data_size_byte.head(2)

Out[14]:

	ID	size	Class
0	9MW5Nuf0ogCEcR1YJeKG	1.113281	8
1	GFblksovaPYLI5XeC8RU	3.808594	9

In [15]:

```
byte_features_with_size = byte_features.merge(data_size_byte, on='ID')
byte_features_with_size.to_csv("result_with_size.csv")
byte_features_with_size.head(2)
```

Out[15]:

	ID	0	1	2	3	4	5	6	7	8	...
0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...
1	01lsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...

2 rows × 260 columns

In [16]:

```
# 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(byte_features_with_size)
```

In [17]:

result.head(2)

Out[17]:

	ID	0	1	2	3	4	5
0	01azqd4lnC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873

2 rows × 260 columns

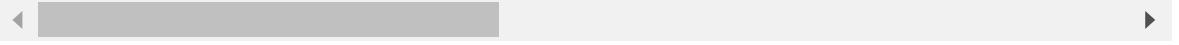
In [18]:

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

Out[18]:

	ID	0	1	2	3	4	5	
0	01azqd4InC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.0
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.0
2	01jsnpXSAlgw6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.0
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.0
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.0

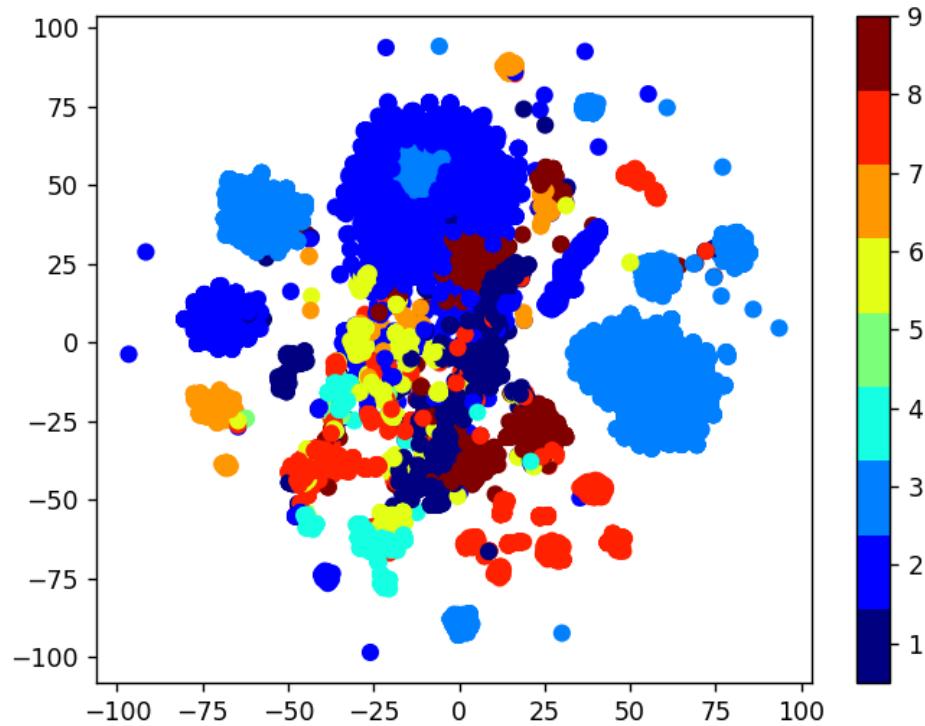
5 rows × 260 columns



3.2.4 Multivariate Analysis

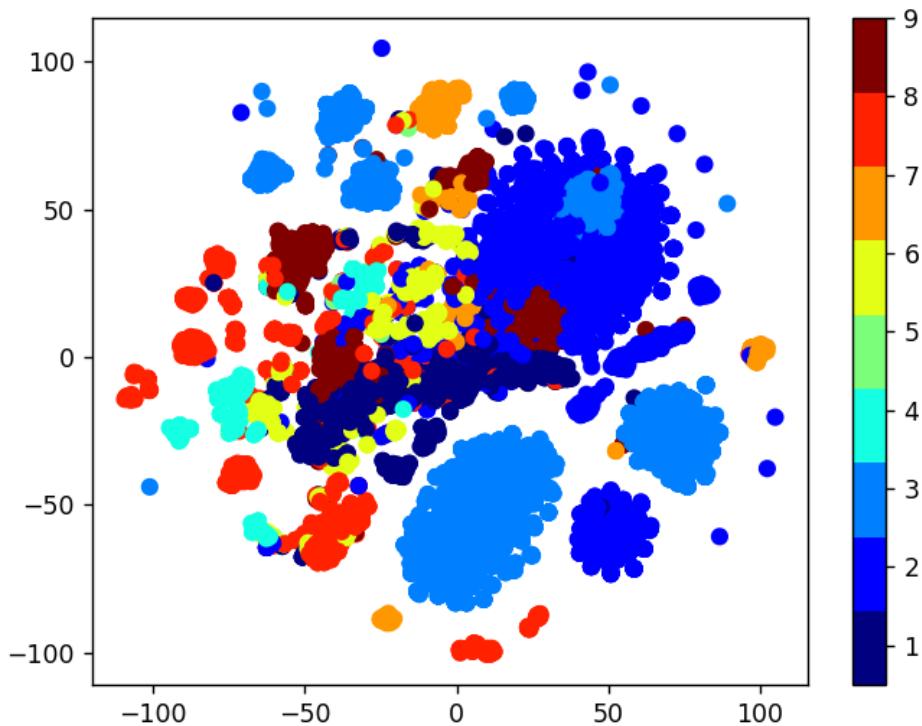
In [0]:

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



In [0]:

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



Train Test split

In [0]:

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

In [0]:

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

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

In [0]:

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

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

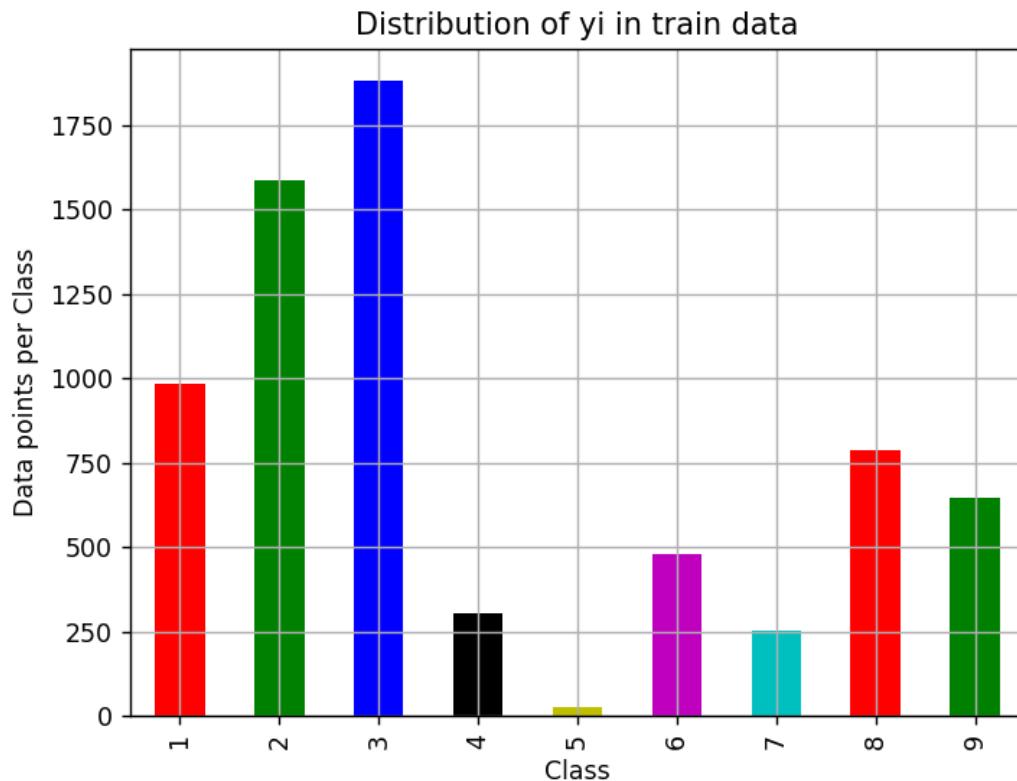
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

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

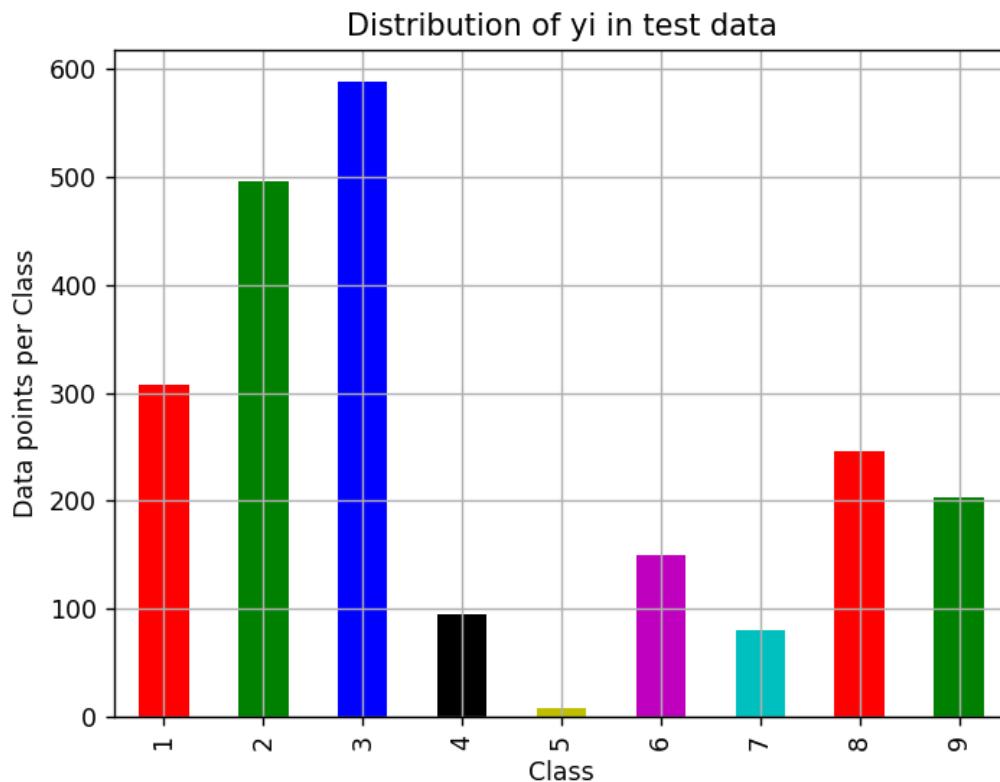
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-test_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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

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

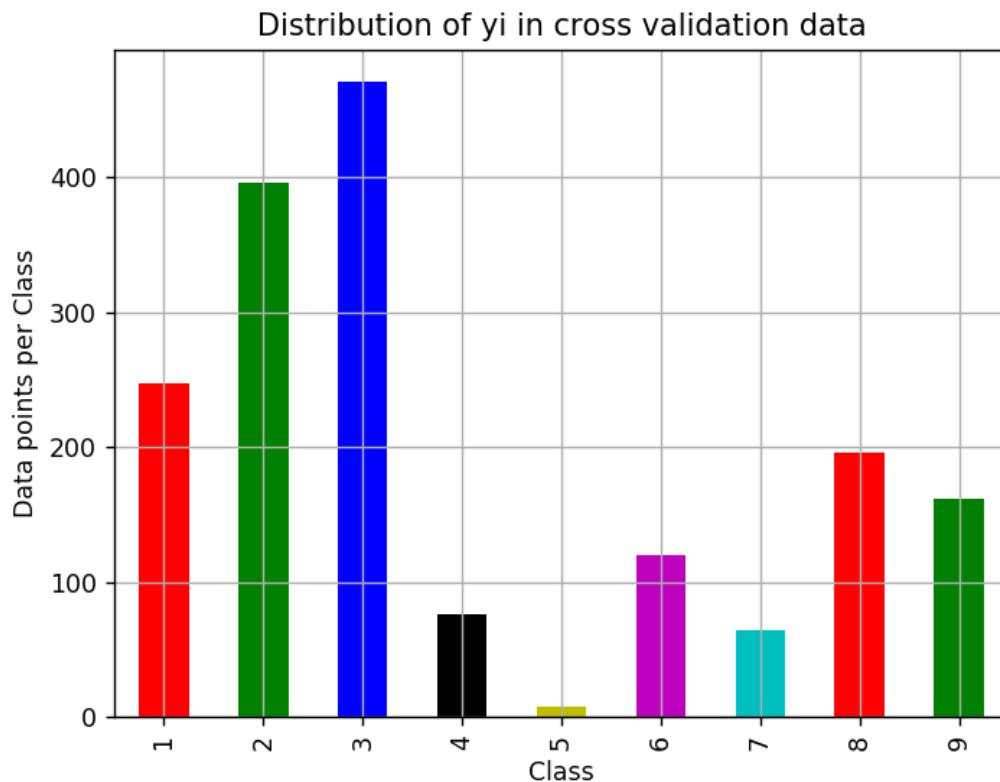
# ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
# -(train_class_distribution.values): the minus sign will give us in decreasing order
sorted_yi = np.argsort(-train_class_distribution.values)
for i in sorted_yi:
    print('Number of data points in class', i+1, ':',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



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



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



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

In [19]:

```

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

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

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

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

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

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

```

```
plt.ylabel('Original Class')
plt.show()
print("Sum of rows in precision matrix",A.sum(axis=1))
```

4. Machine Learning Models

4.1. Machine Learning Models on bytes files

4.1.1. Random Model

In [0]:

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

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

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

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

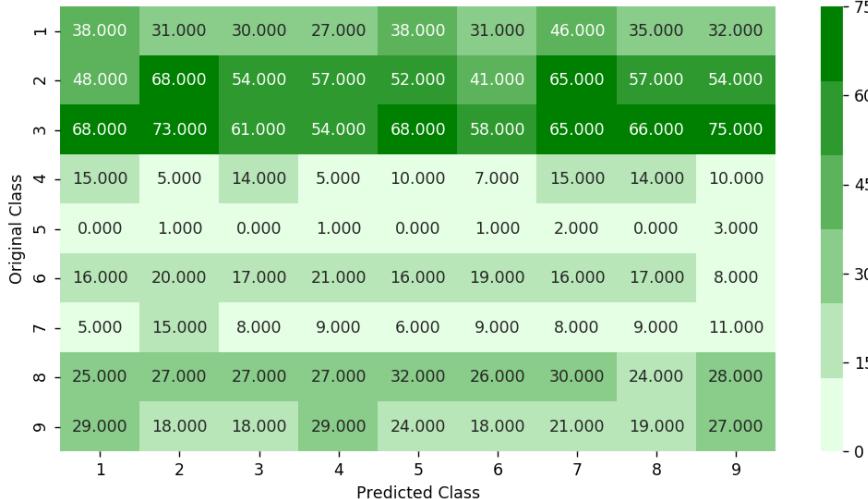
predicted_y =np.argmax(test_predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y+1)
```

Log loss on Cross Validation Data using Random Model 2.45615644965

Log loss on Test Data using Random Model 2.48503905509

Number of misclassified points 88.5004599816

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

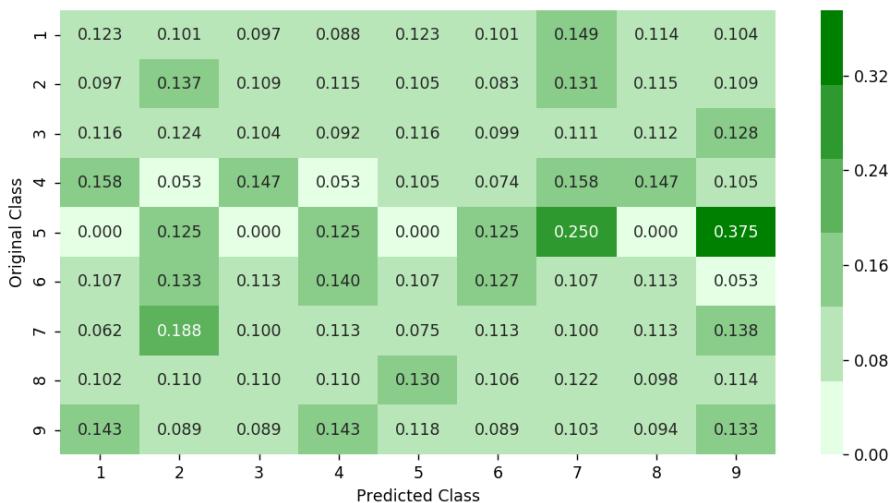


----- Precision matrix -----



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

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



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

4.1.2. K Nearest Neighbour Classification

In [0]:

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

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

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default 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])

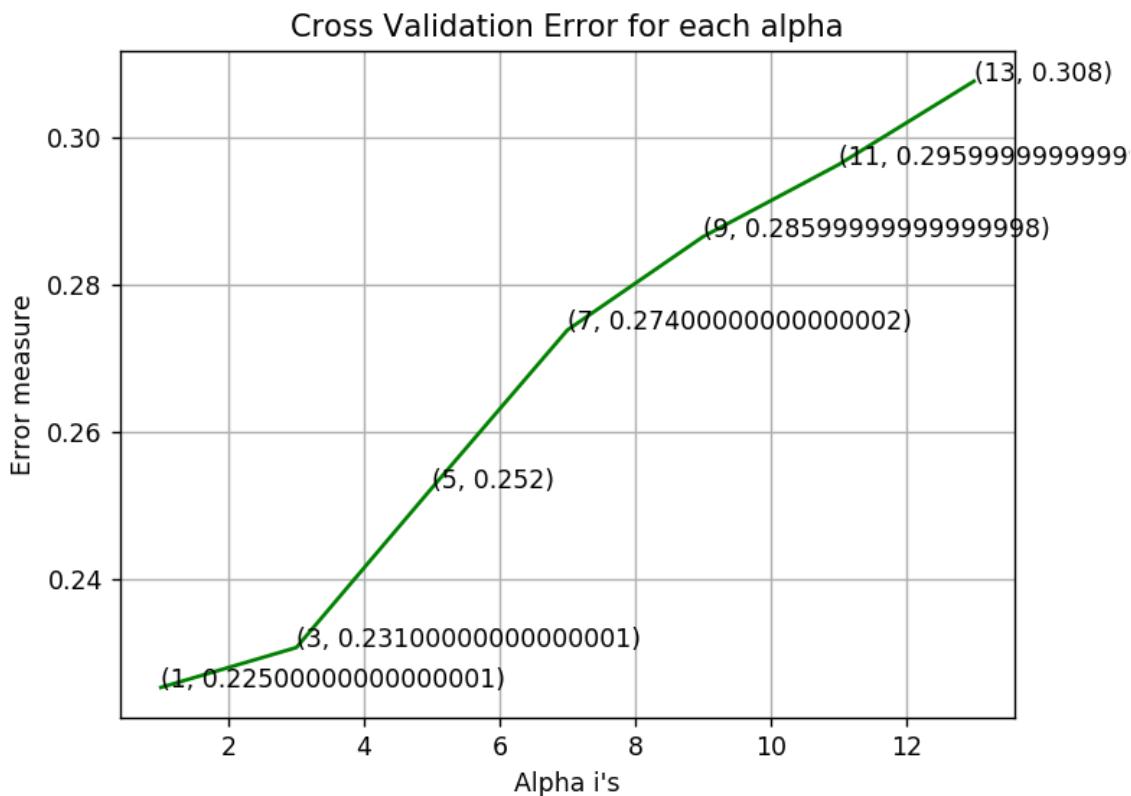
best_alpha = np.argmin(cv_log_error_array)

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()
```

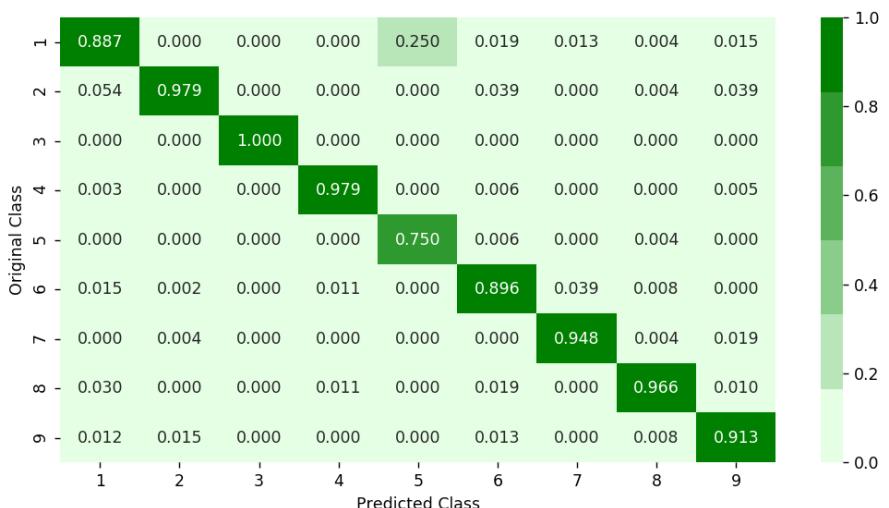
```
k_cfl=KNeighborsClassifier(n_neighbors=alpha[best_alpha])
k_cfl.fit(X_train,y_train)
sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
sig_clf.fit(X_train, y_train)

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

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

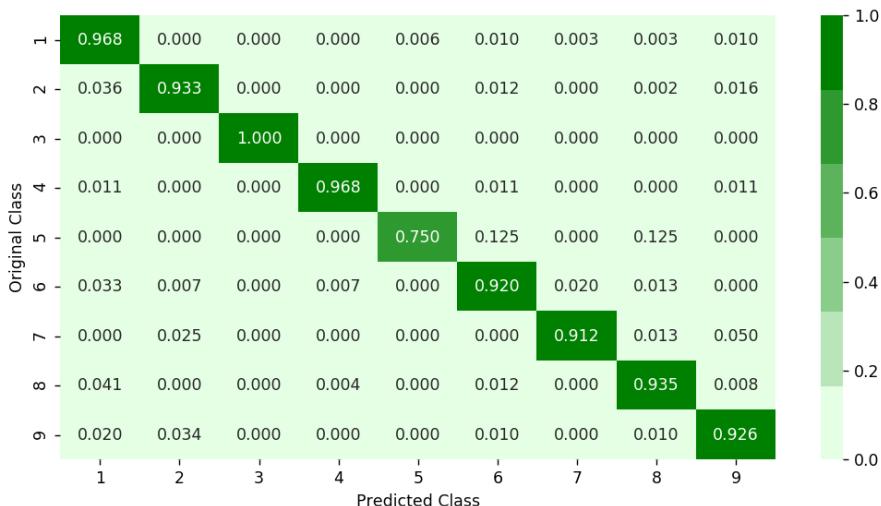


```
For values of best alpha = 1 The train log loss is: 0.0782947669247
For values of best alpha = 1 The cross validation log loss is: 0.22538623
7304
For values of best alpha = 1 The test log loss is: 0.241508604195
Number of misclassified points 4.50781968721
----- Confusion matrix -----
```

**Precision matrix**

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

Recall matrix



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

4.1.3. Logistic Regression

In [0]:

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

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

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

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

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

best_alpha = np.argmin(cv_log_error_array)

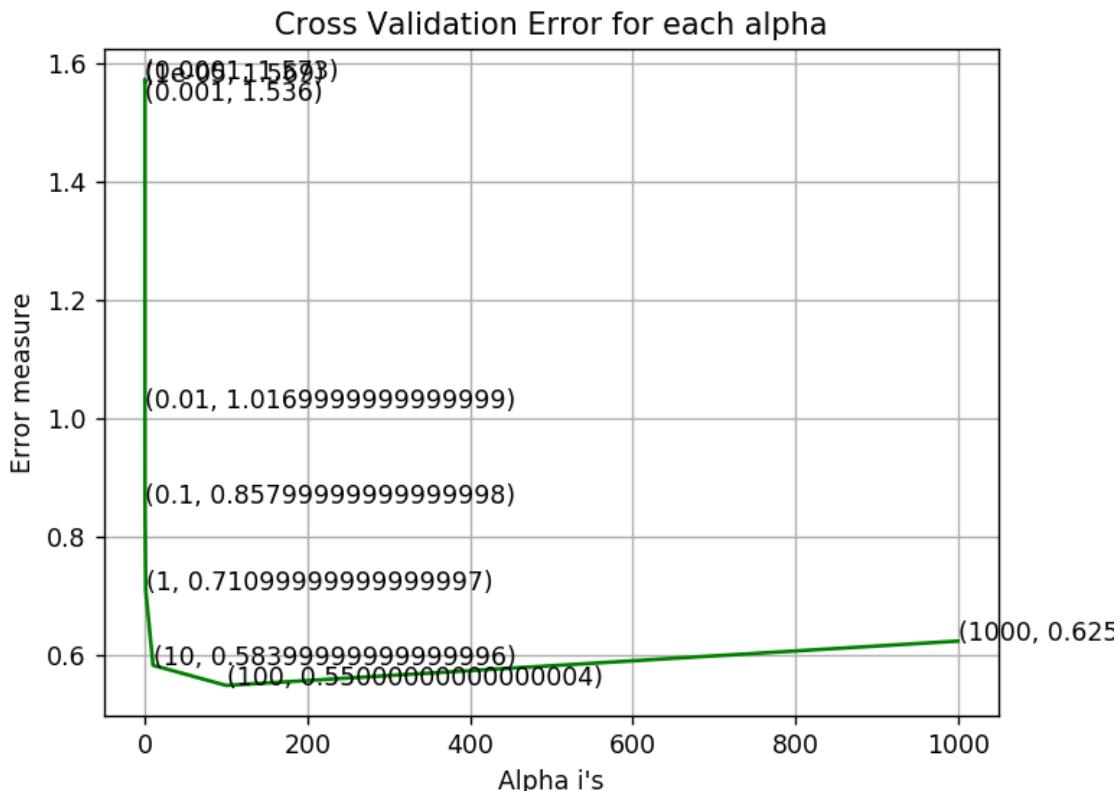
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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

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

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

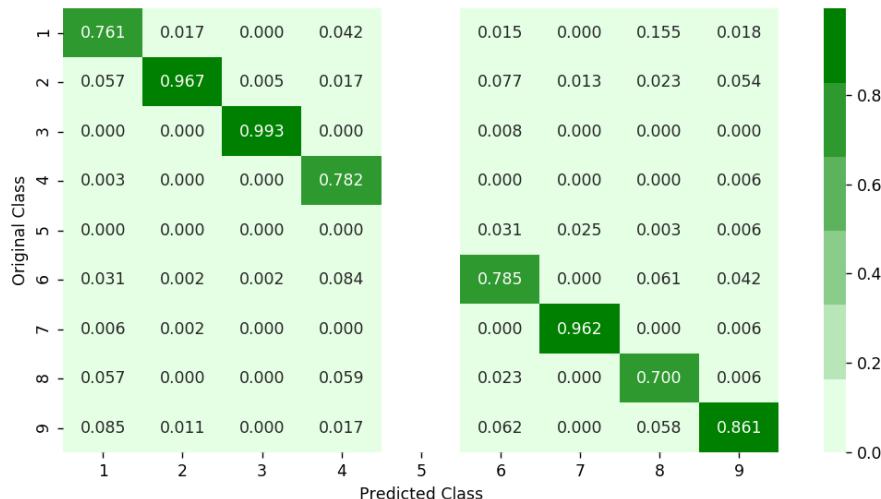


log loss for train data 0.498923428696
log loss for cv data 0.549929846589
log loss for test data 0.528347316704
Number of misclassified points 12.3275068997

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



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



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

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



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

4.1.4. Random Forest Classifier

In [0]:

```

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

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

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

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


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

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

best_alpha = np.argmin(cv_log_error_array)

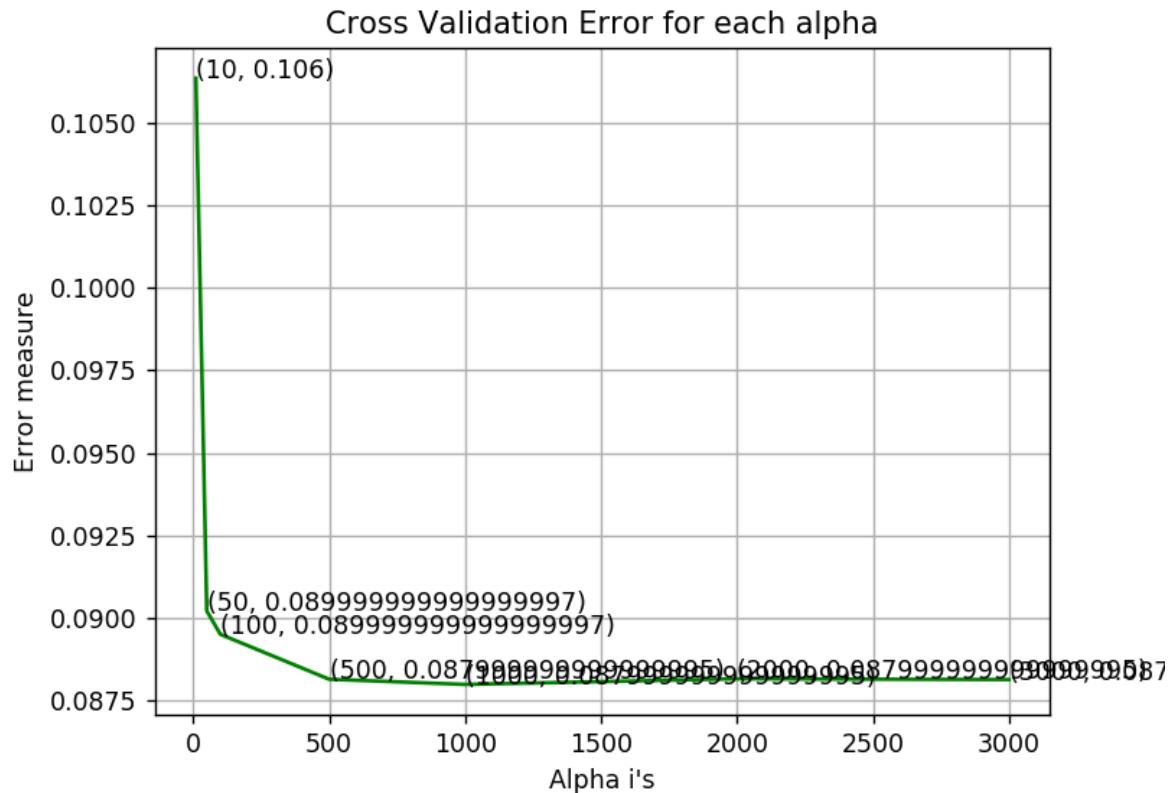
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

```

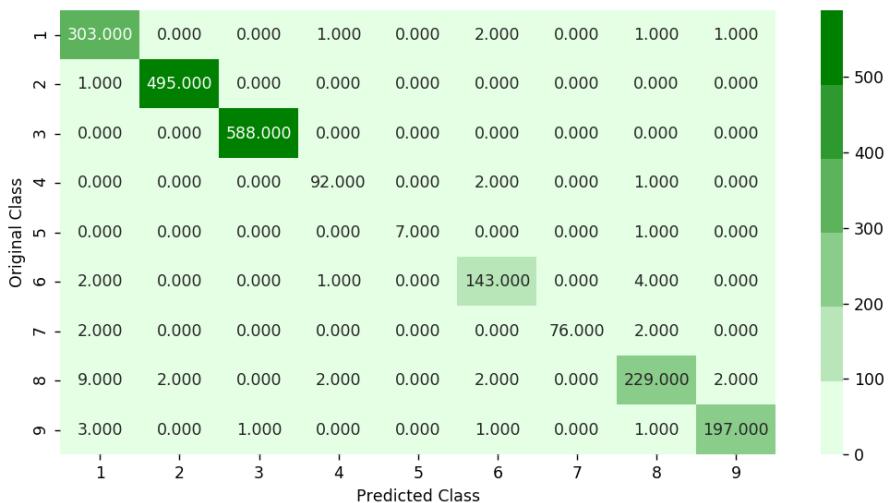
```
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv,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

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

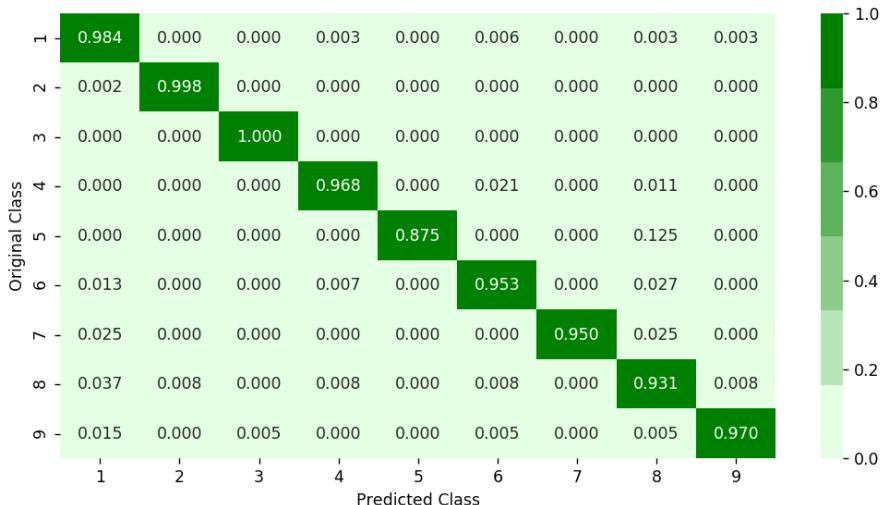


For values of best alpha = 1000 The train log loss is: 0.0266476291801
 For values of best alpha = 1000 The cross validation log loss is: 0.0879849524621
 For values of best alpha = 1000 The test log loss is: 0.0858346961407
 Number of misclassified points 2.02391904324

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

**Precision matrix**

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

Recall matrix

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

4.1.5. XgBoost Classification

In [0]:

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

# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# default 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_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
# s)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep])      Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video Link1: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/regression-using-decision-trees-2/
# video Link2: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/what-are-ensembles/
# -----


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

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

best_alpha = np.argmin(cv_log_error_array)

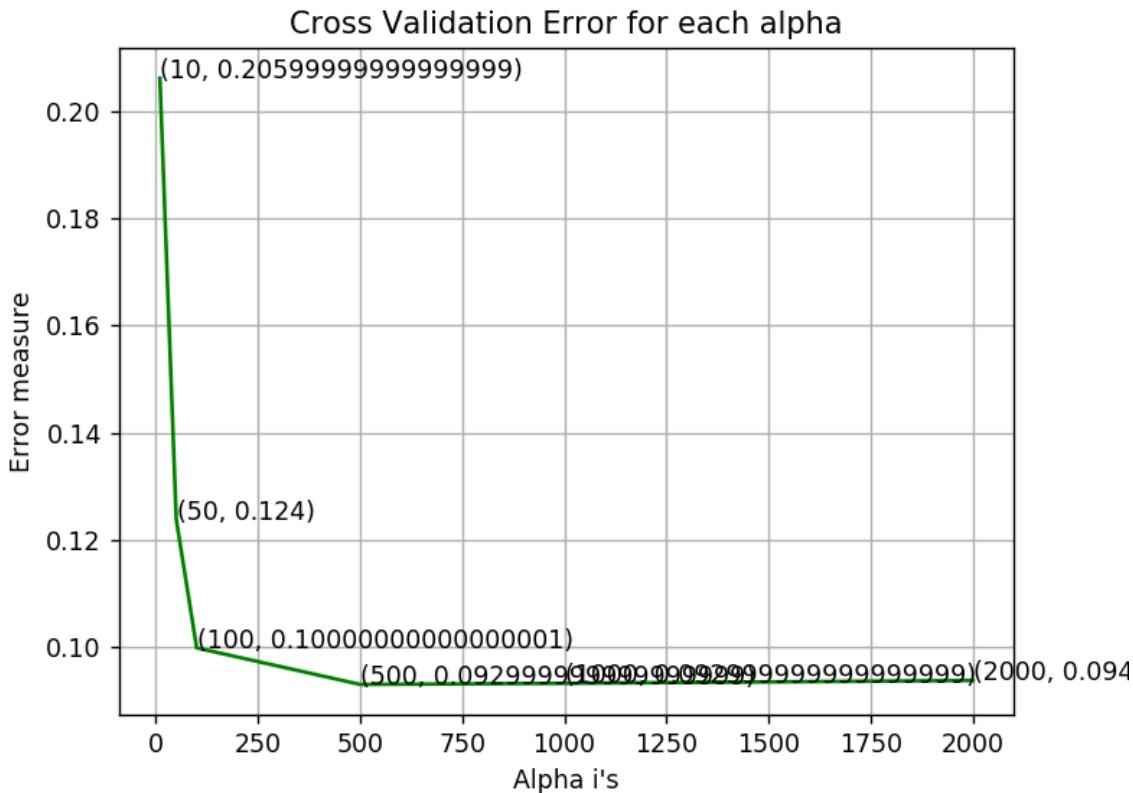
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

```
sig_clf.fit(X_train, y_train)

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

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



For values of best alpha = 500 The train log loss is: 0.0225231805824

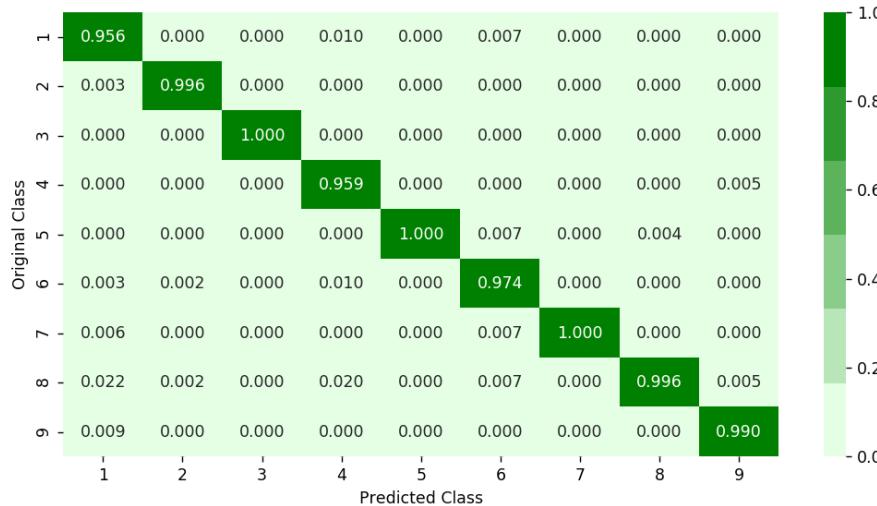
For values of best alpha = 500 The cross validation log loss is: 0.0931035681289

For values of best alpha = 500 The test log loss is: 0.0792067651731

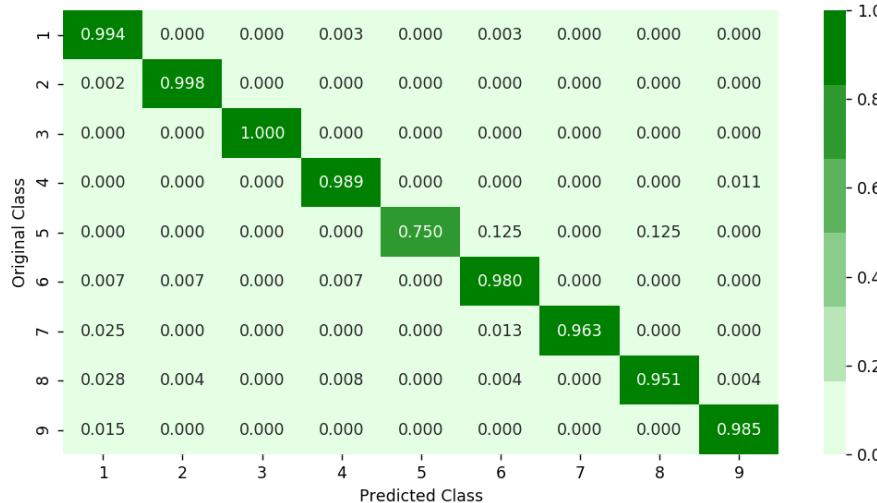
Number of misclassified points 1.24195032199

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

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

Precision matrix

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

Recall matrix

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

4.1.5. XgBoost Classification with best hyper parameters using RandomSearch

In [0]:

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

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

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

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

Out[0]:

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

In [0]:

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

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

In [0]:

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

# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# default 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_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
# s)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep])      Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video Link2: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/what-are-ensembles/
# -----


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

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

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

4.2 Modeling with .asm files

There are 10868 files of asm
All the files make up about 150 GB

The asm files contains :

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

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

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

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

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

4.2.1 Feature extraction from asm files

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

In [0]:

```
#initially create five folders
#first
#second
#third
#fourth
#fifth
#this code tells us about random split of files into five folders
folder_1 ='first'
folder_2 ='second'
folder_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]],'third')
    elif i%5 ==3:
        shutil.move(source+files[data[i]],'fourth')
    elif i%5==4:
        shutil.move(source+files[data[i]],'fifth')
```

In [0]:

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

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

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

    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc',
    'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'jz',
    'rtn', 'lea', 'movzx']
    #best keywords that are taken from different blogs
    keywords = ['.dll','std::',':dword']
    #Below taken registers are general purpose registers and special registers
    #All the registers which are taken are best
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\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',
    '.rsrc','.tls','.reloc','.BSS','.CODE']
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc',
    'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'jz',
    'rtn', 'lea', 'movzx']
    keywords = ['.dll','std::':':dword']
    registers=['edx','esi','eax','ebx','ecx','edi','ebp','esp','eip']
    file1=open("output\mediumasmfile.txt","w+")
    files = os.listdir('second')
    for f in files:
        prefixescount=np.zeros(len(prefixes),dtype=int)
        opcodescount=np.zeros(len(opcodes),dtype=int)
        keywordcount=np.zeros(len(keywords),dtype=int)
        registerscount=np.zeros(len(registers),dtype=int)
        features=[]
        f2=f.split('.')[0]
        file1.write(f2+",")
        opcodefile.write(f2+" ")
        with codecs.open('second/'+f,encoding='cp1252',errors ='replace') as fli:
            for lines in fli:
                line=lines.rstrip().split()
                l=line[0]
                for i in range(len(prefixes)):
                    if prefixes[i] in line[0]:
                        prefixescount[i]+=1
                line=line[1:]
                for i in range(len(opcodes)):
                    if any(opcodes[i]==li for li in line):
                        features.append(opcodes[i])
                        opcodescount[i]+=1
                for i in range(len(registers)):
                    for li in line:
                        if registers[i] in li and ('text' in l or 'CODE' in l):
                            registerscount[i]+=1
                for i in range(len(keywords)):
                    for li in line:
                        if keywords[i] in li:
                            keywordcount[i]+=1
        for prefix in prefixescount:
            file1.write(str(prefix)+",")
        for opcode in opcodescount:
            file1.write(str(opcode)+",")
        for register in registerscount:
            file1.write(str(register)+",")

```

```

        for key in keywordcount:
            file1.write(str(key)+",")
        file1.write("\n")
        file1.close()

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

def fourthprocess():
    prefixes = [ 'HEADER:', '.text:', '.Pav:', '.idata:', '.data:', '.bss:', '.rdata:', '.edata:', '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE' ]
    opcodes = [ 'jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc',

```

```

'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'j
z', 'rtn', 'lea', 'movzx']
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',
    '.rsrc:', '.tls:', '.reloc:', '.BSS:', '.CODE']
    opcodes = ['jmp', 'mov', 'retf', 'push', 'pop', 'xor', 'retn', 'nop', 'sub', 'inc',
    'dec', 'add', 'imul', 'xchg', 'or', 'shr', 'cmp', 'call', 'shl', 'ror', 'rol', 'jnb', 'j
z', 'rtn', 'lea', 'movzx']
    keywords = ['.dll', 'std::', ':dword']
    registers=['edx', 'esi', 'eax', 'ebx', 'ecx', 'edi', 'ebp', 'esp', 'eip']
    file1=open("output\trainasmfile.txt", "w+")
    files = os.listdir('fifth/')
    for f in files:
        prefixescount=np.zeros(len(prefixes), dtype=int)
        opcodescount=np.zeros(len(opcodes), dtype=int)
        keywordcount=np.zeros(len(keywords), dtype=int)

```

```

registerscount=np.zeros(len(registers),dtype=int)
features=[]
f2=f.split('.')[0]
file1.write(f2+",")
opcodefile.write(f2+" ")
with codecs.open('fifth/'+f,encoding='cp1252',errors ='replace') as fli:
    for lines in fli:
        line=lines.rstrip().split()
        l=line[0]
        for i in range(len(prefixes)):
            if prefixes[i] in line[0]:
                prefixescount[i]+=1
        line=line[1:]
        for i in range(len(opcodes)):
            if any(opcodes[i]==li for li in line):
                features.append(opcodes[i])
                opcodescount[i]+=1
        for i in range(len(registers)):
            for li in line:
                if registers[i] in li and ('text' in 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()
```

In [20]:

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

Out[20]:

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

5 rows × 53 columns



4.2.1.1 Files sizes of each .asm file

In [22]:

```
#file sizes of byte files

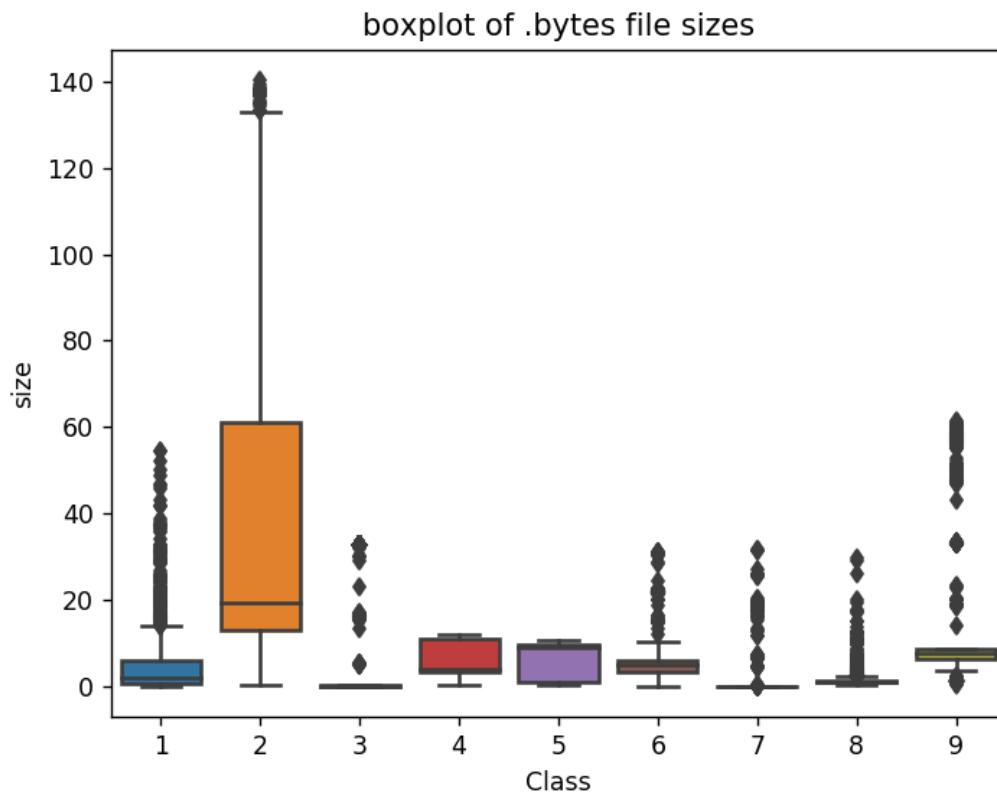
files=os.listdir('asmFiles')
filenames=Y['ID'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, st_nlink=1, st_uid=0,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

	ID	size	Class
0	0gWUIudhwovMYb3NSnZA	7.579200	9
1	aTfDJBX2PNZRIjUG1SKz	0.376585	1
2	7edjQGTXowvNgS9uVi5J	6.742073	7
3	HaWuy5IXN7gOmJG4ZCKd	60.645524	2
4	igK38Wl1FjDJSnhYZves	0.163010	3

4.2.1.2 Distribution of .asm file sizes

In [0]:

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



In [23]:

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

(10868, 53)
(10868, 3)

Out[23]:

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

5 rows × 11 columns

In [24]:

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

Out[24]:

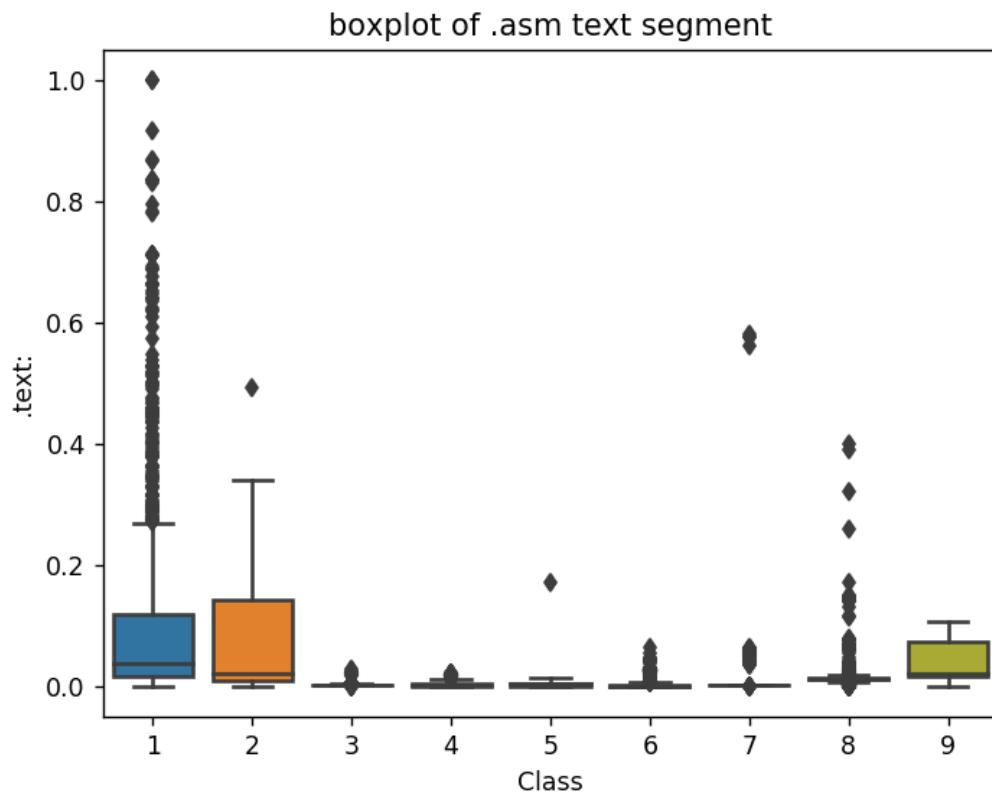
	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:
0	01kcPWA9K2BOxQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084
1	1E93CpP60RHFNiT5Qfvn	0.096045	0.001230	0.0	0.000617	0.000019	0.0	0.000000
2	3ekVow2ajZHbTnBcsDfX	0.096045	0.000627	0.0	0.000300	0.000017	0.0	0.000038
3	3X2nY7iQaPBIWDrAZqJe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

5 rows × 9 columns

4.2.2 Univariate analysis on asm file features

In [0]:

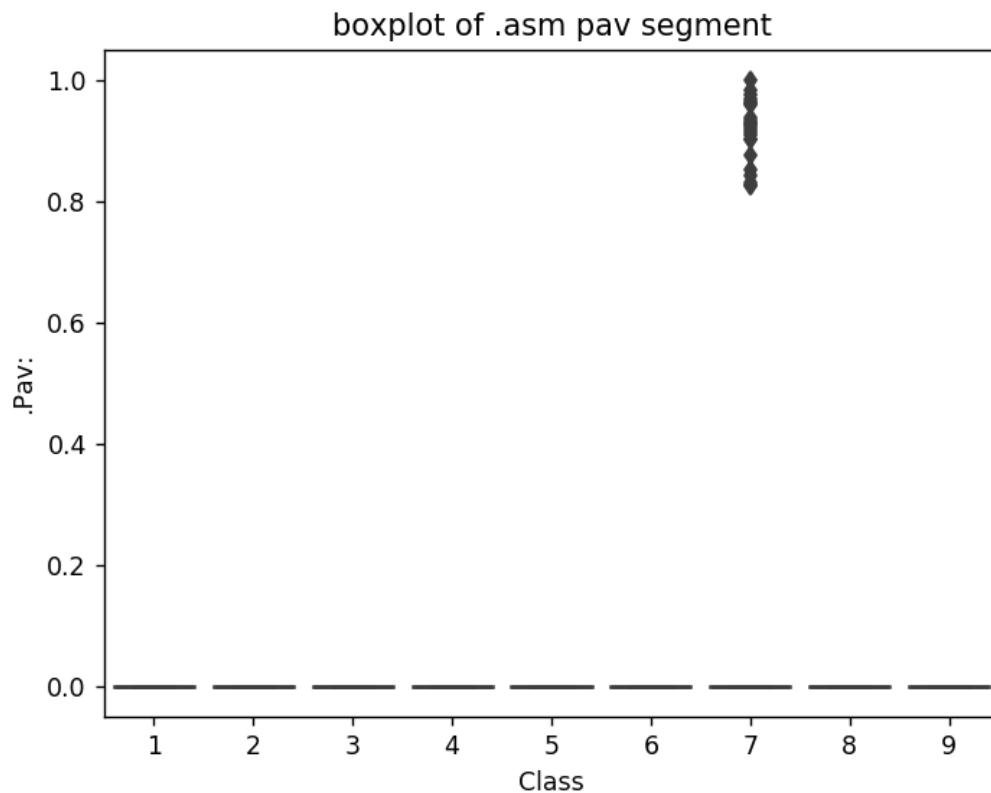
```
ax = sns.boxplot(x="Class", y=".text:", data=result_asm)
plt.title("boxplot of .asm text segment")
plt.show()
```



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

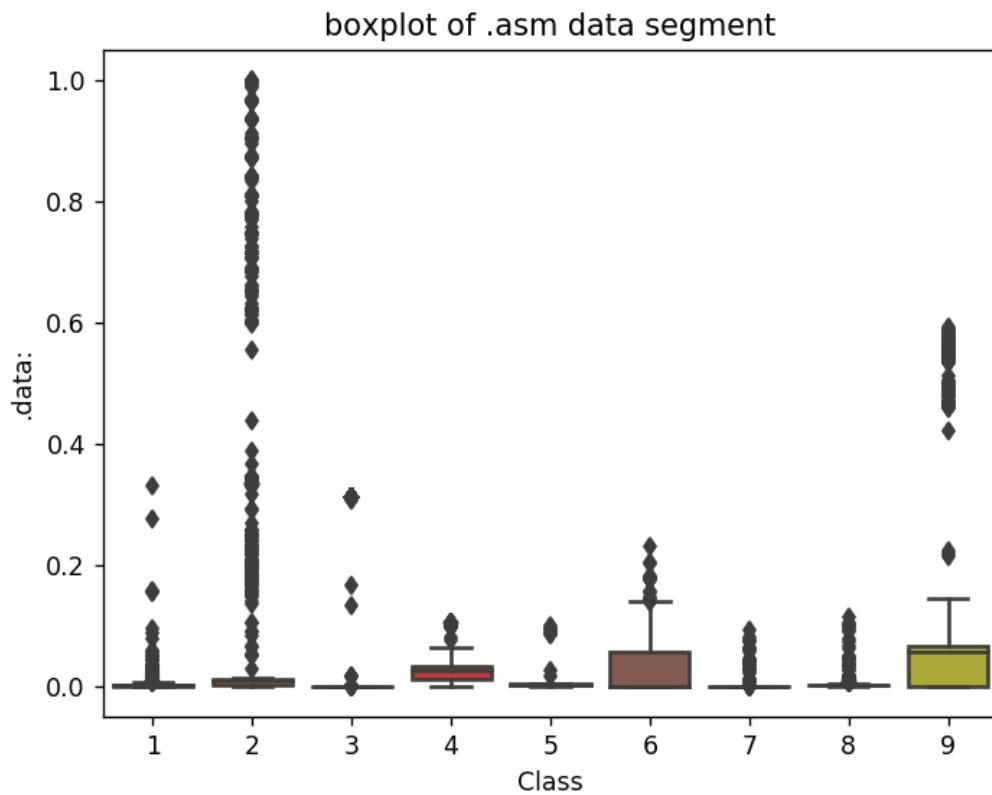
In [0]:

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



In [0]:

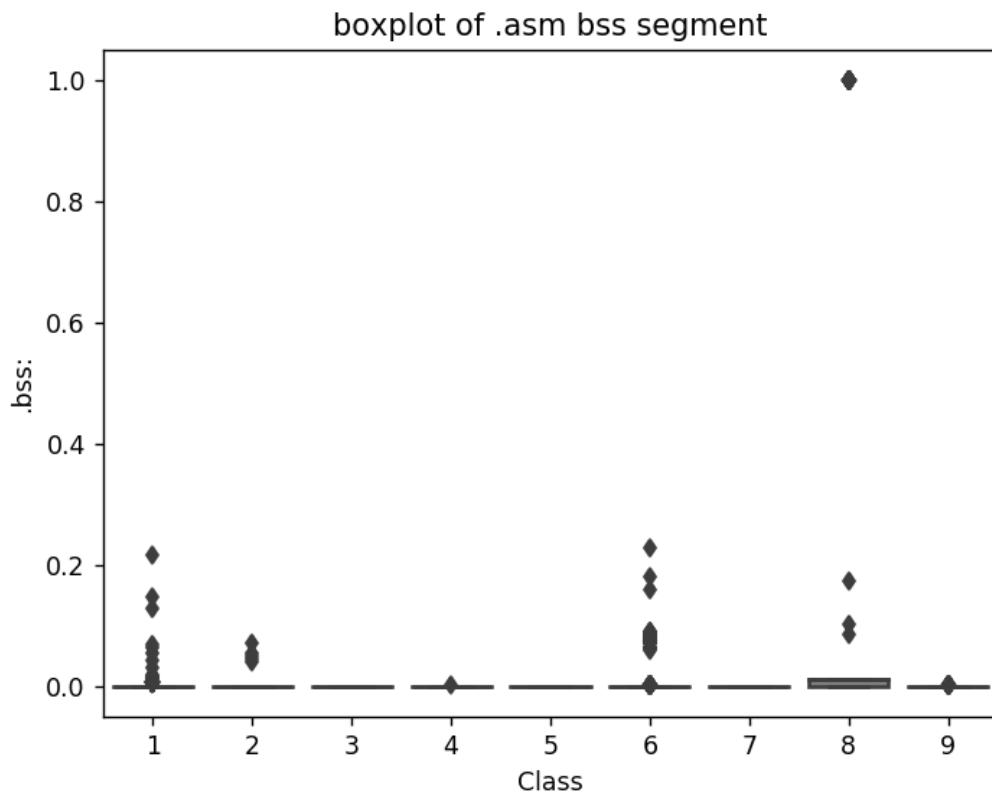
```
ax = sns.boxplot(x="Class", y=".data:", data=result_asm)
plt.title("boxplot of .asm data segment")
plt.show()
```



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

In [0]:

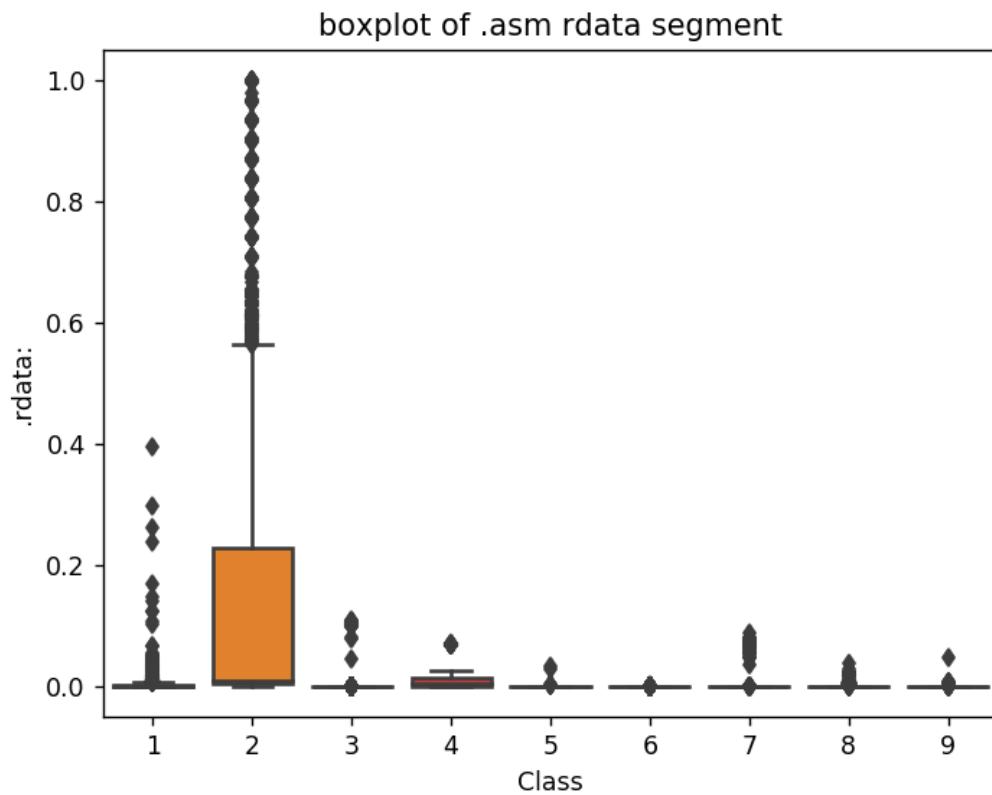
```
ax = sns.boxplot(x="Class", y=".bss:", data=result_asm)
plt.title("boxplot of .asm bss segment")
plt.show()
```



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

In [0]:

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

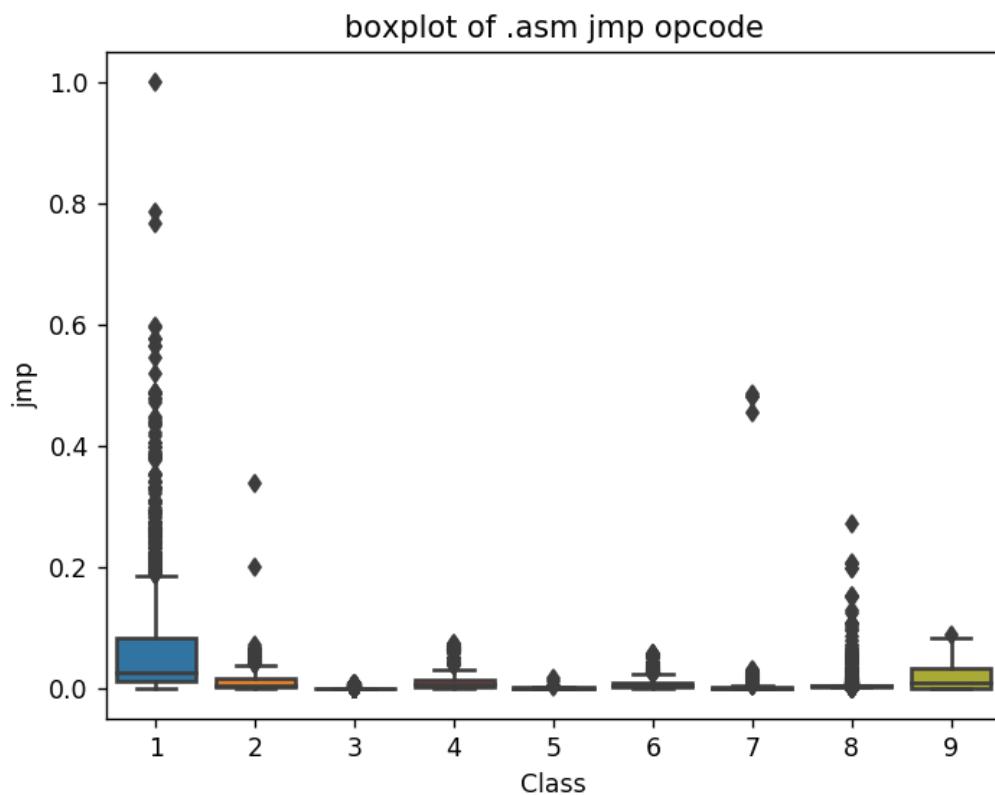


Plot between rdata segment and Class segment

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

In [0]:

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

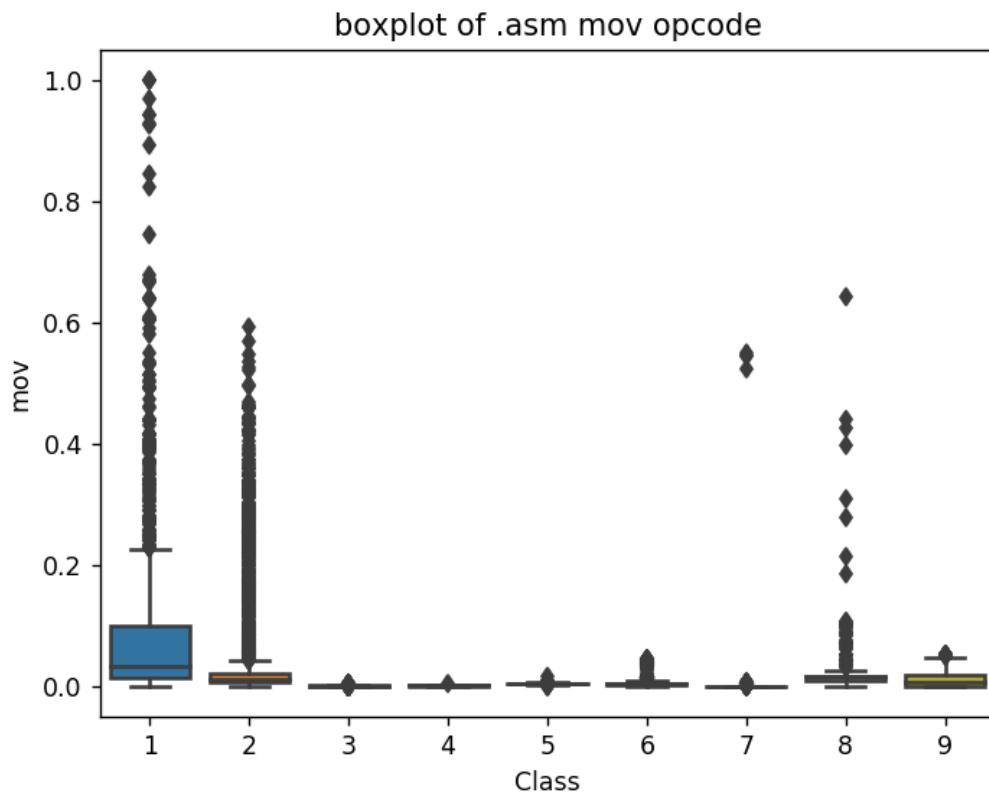


plot between jmp and Class label

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

In [0]:

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

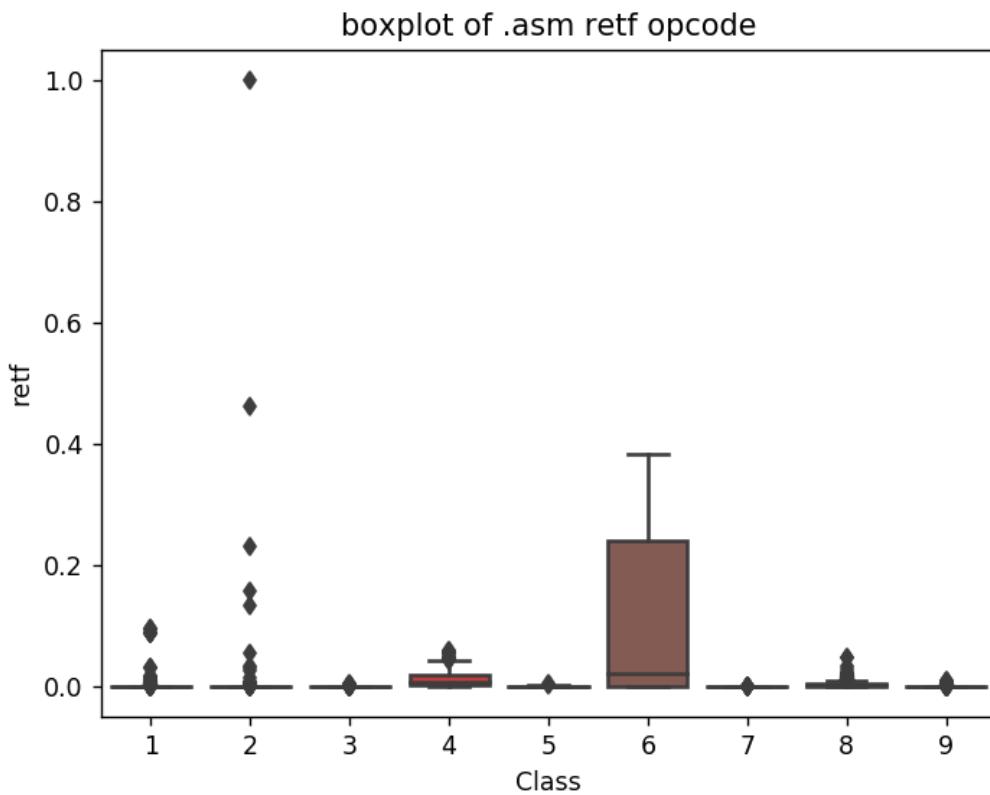


plot between Class label and mov opcode

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

In [0]:

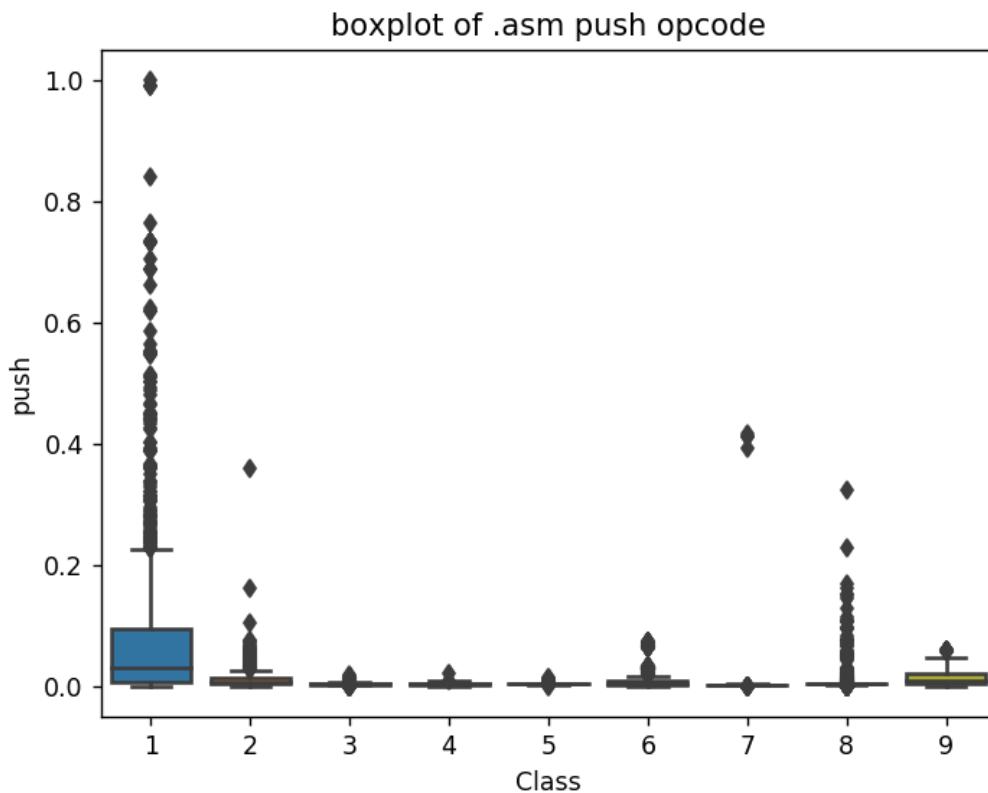
```
ax = sns.boxplot(x="Class", y="retf", data=result_asm)
plt.title("boxplot of .asm retf opcode")
plt.show()
```



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

In [0]:

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



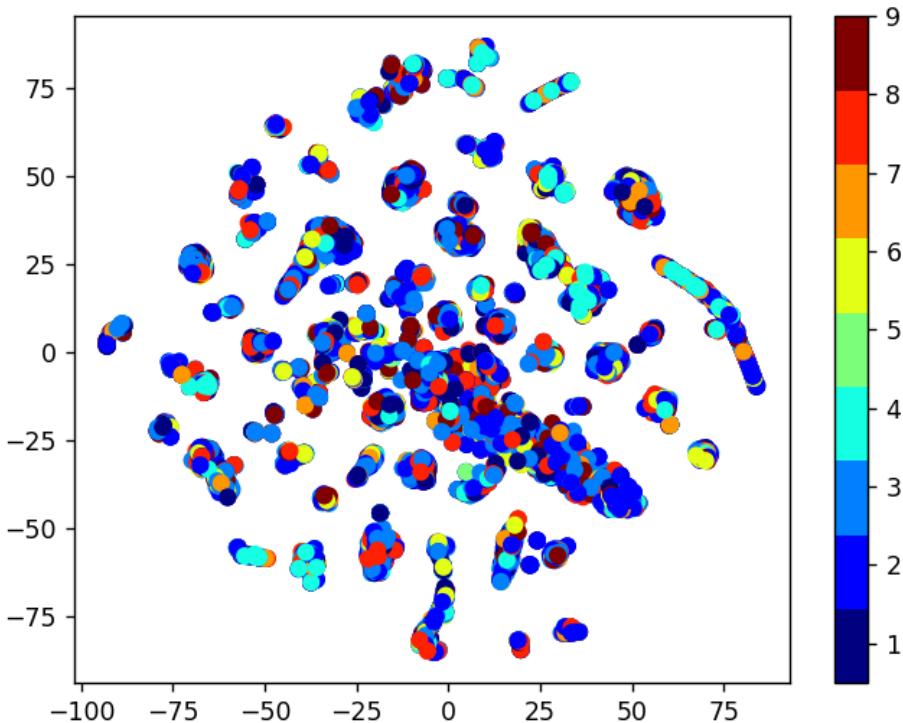
plot between push opcode and Class label
Class 1 is having 75 percentile files with push opcodes of frequency 1000

4.2.2 Multivariate Analysis on .asm file features

In [0]:

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

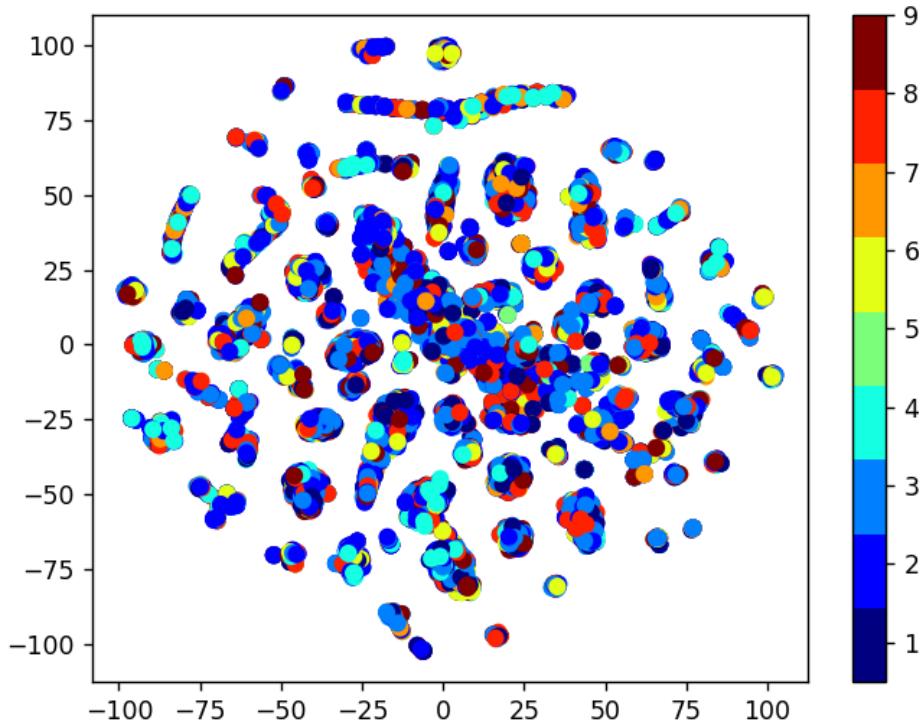
#multivariate analysis on byte files
#this is with perplexity 50
xtsne=TSNE(perplexity=50)
results=xtsne.fit_transform(result_asm.drop(['ID','Class'], axis=1).fillna(0))
vis_x = results[:, 0]
vis_y = results[:, 1]
plt.scatter(vis_x, vis_y, c=data_y, cmap=plt.cm.get_cmap("jet", 9))
plt.colorbar(ticks=range(10))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



In [0]:

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

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



TSNE for asm data with perplexity 50

4.2.3 Conclusion on EDA

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

4.3 Train and test split

In [25]:

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

In [0]:

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

In [0]:

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

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

4.4. Machine Learning models on features of .asm files

4.4.1 K-Nearest Neighbors

In [0]:

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

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

# find more about CalibratedClassifierCV here at http://scikit-learn.org/stable/modules/generated/sklearn.calibration.CalibratedClassifierCV.html
# -----
# default 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])

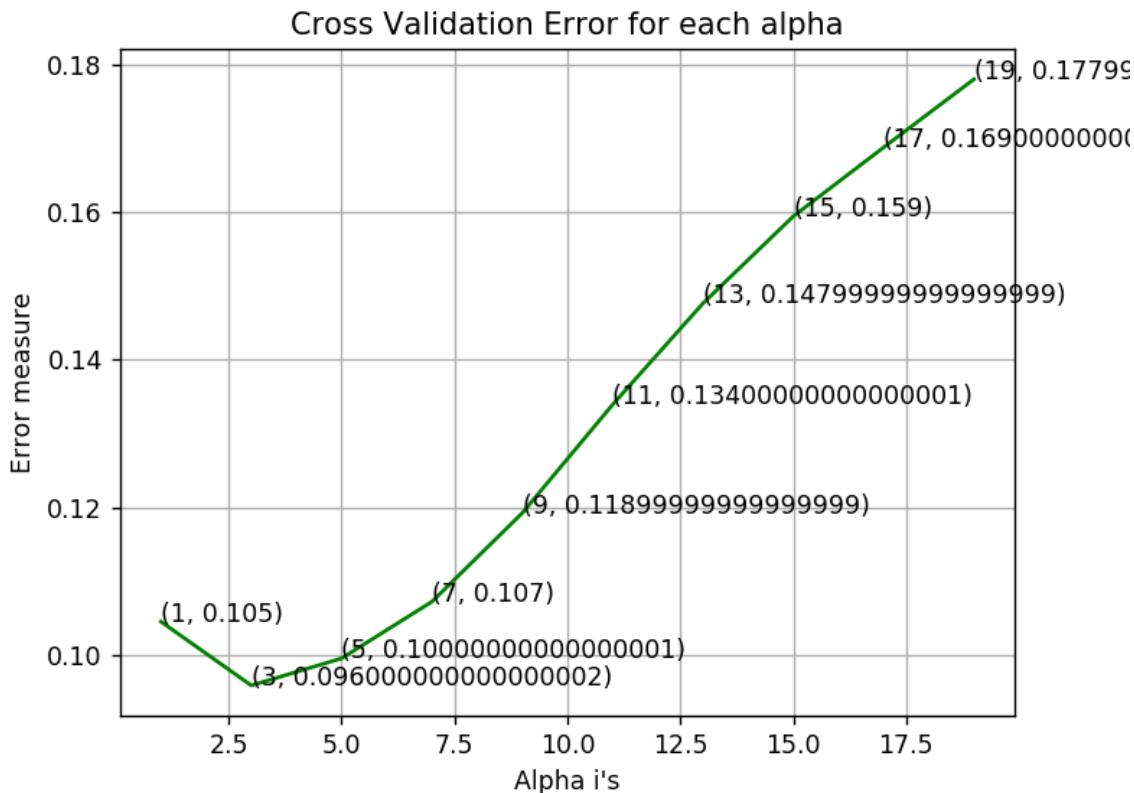
best_alpha = np.argmin(cv_log_error_array)

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

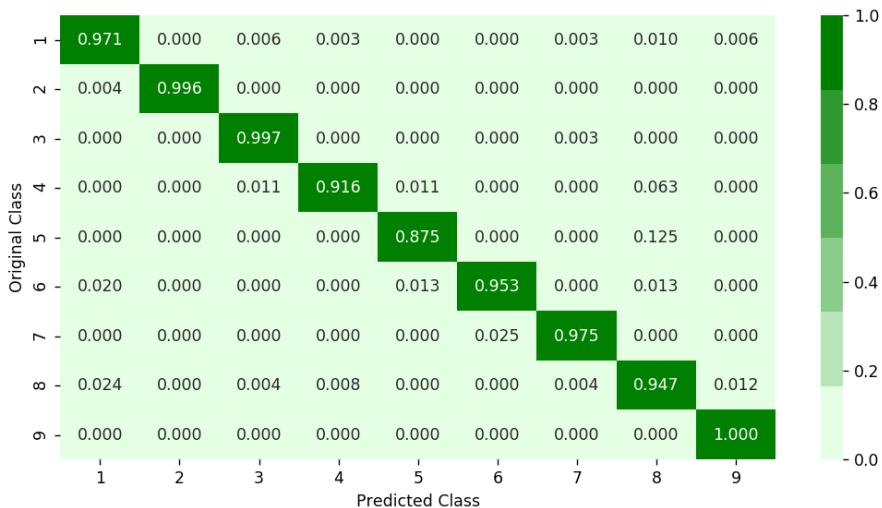


```
log loss for train data 0.0476773462198
log loss for cv data 0.0958800580948
log loss for test data 0.0894810720832
Number of misclassified points 2.02391904324
----- Confusion matrix -----
```

**Precision matrix**

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

Recall matrix



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

4.4.2 Logistic Regression

In [0]:

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

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

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

alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
 logisticR=LogisticRegression(penalty='12',C=i,class_weight='balanced')
 logisticR.fit(X_train_asm,y_train_asm)
 sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
 sig_clf.fit(X_train_asm, y_train_asm)
 predict_y = sig_clf.predict_proba(X_cv_asm)
 cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.classes_,
eps=1e-15))

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

best_alpha = np.argmin(cv_log_error_array)

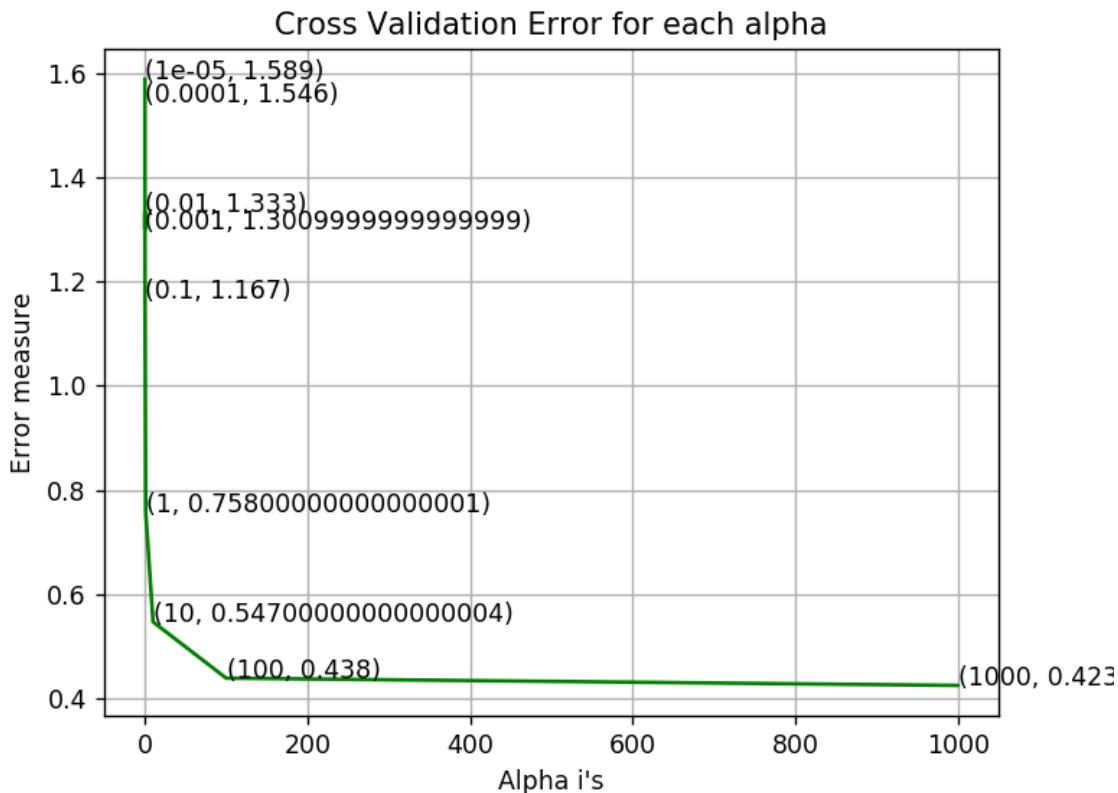
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
 ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

logisticR=LogisticRegression(penalty='12',C=alpha[best_alpha],class_weight='balanced')
logisticR.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

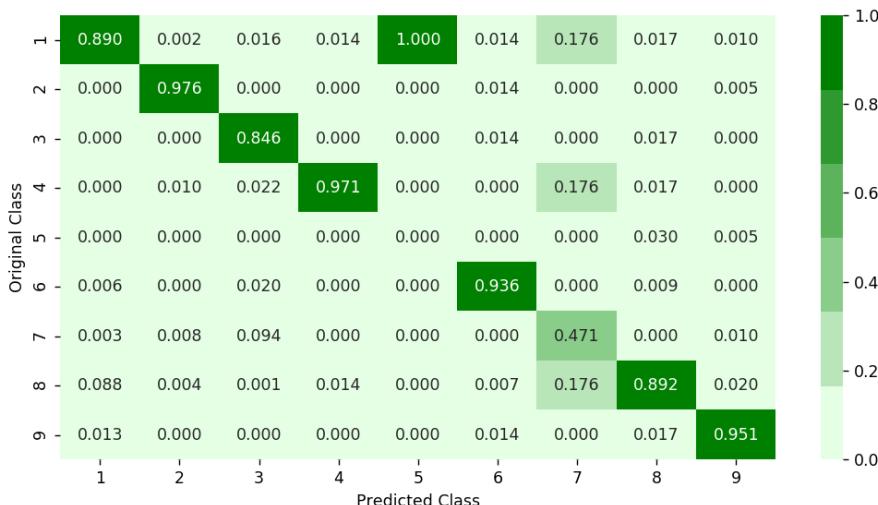
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=logisticR.classes_,
eps=1e-15)))
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=logisticR.classes_,
eps=1e-15)))
predict_y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

```
print ('log loss for test data',(log_loss(y_test_asm, predict_y, labels=logisticR.classes_, eps=1e-15)))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```

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

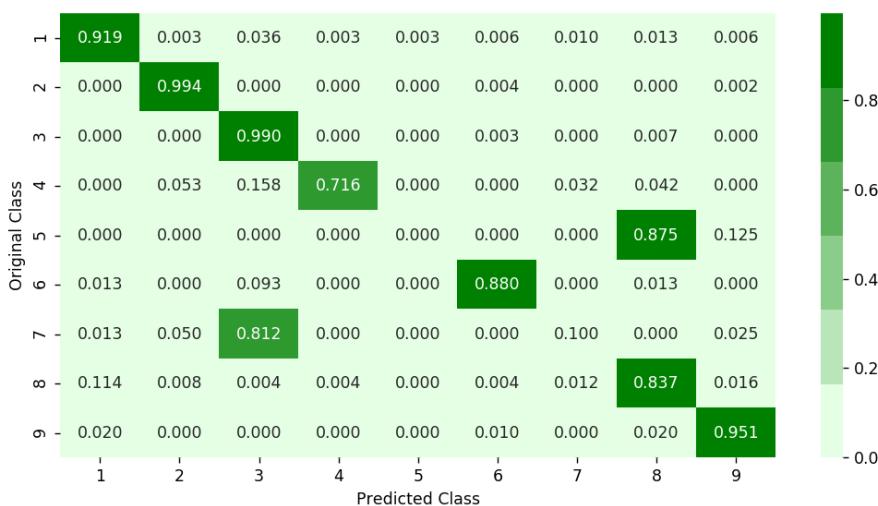


```
log loss for train data 0.396219394701
log loss for cv data 0.424423536526
log loss for test data 0.415685592517
Number of misclassified points 9.61361545538
----- Confusion matrix -----
```

**Precision matrix**

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

Recall matrix



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

4.4.3 Random Forest Classifier

In [0]:

```

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

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

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

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


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

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

best_alpha = np.argmin(cv_log_error_array)

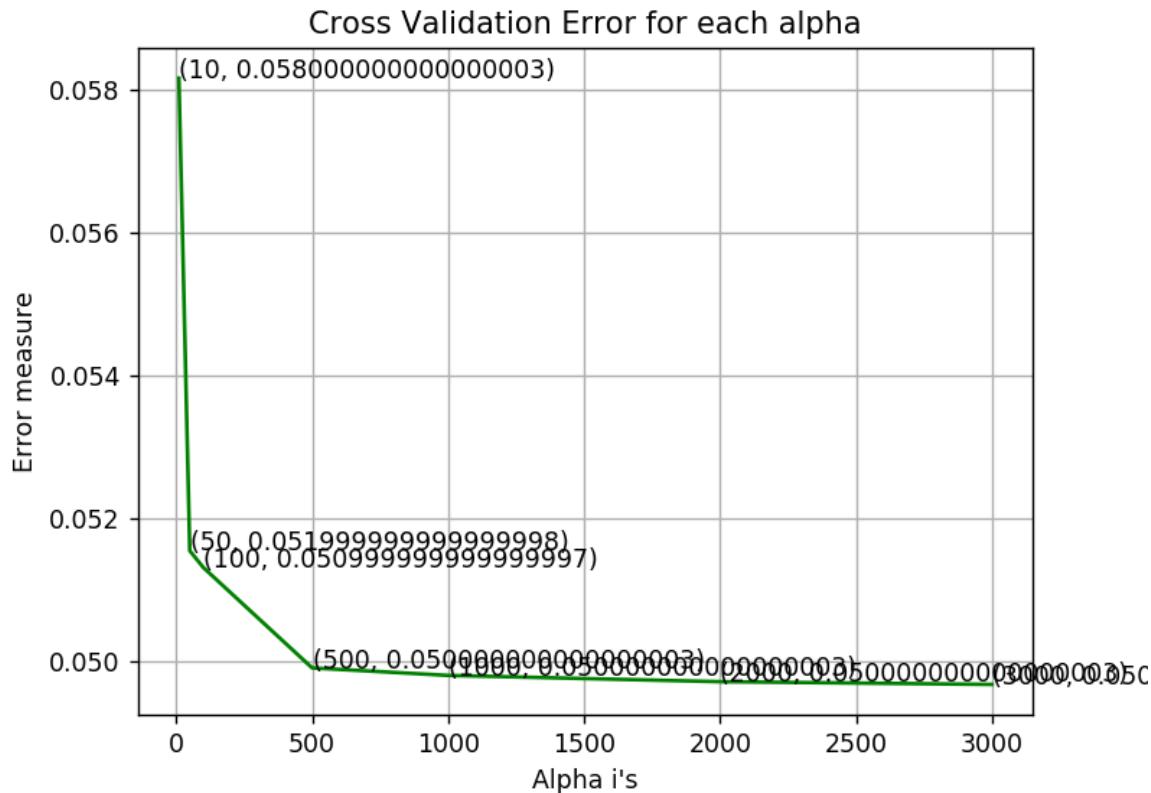
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

```

```
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=sig_clf.classes_, e
ps=1e-15)))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',(log_loss(y_test_asm, predict_y, labels=sig_clf.classes
_, eps=1e-15)))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```

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

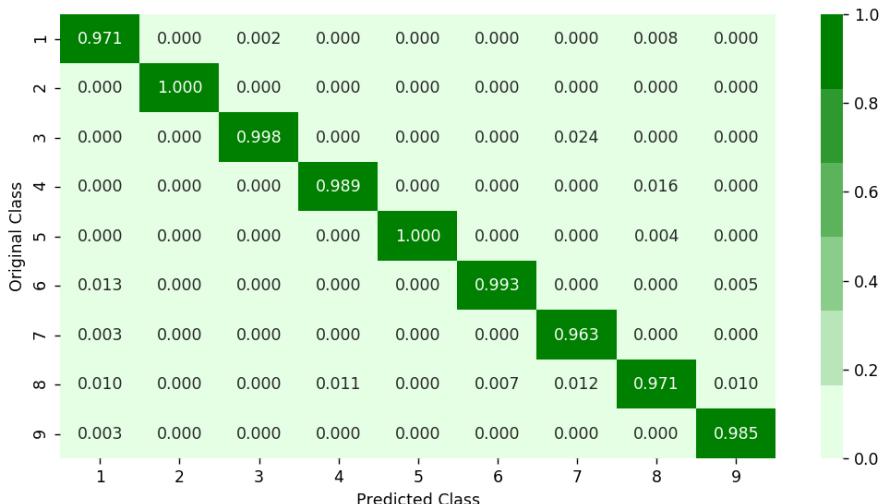


```
log loss for train data 0.0116517052676
log loss for cv data 0.0496706817633
log loss for test data 0.0571239496453
Number of misclassified points 1.14995400184
```

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

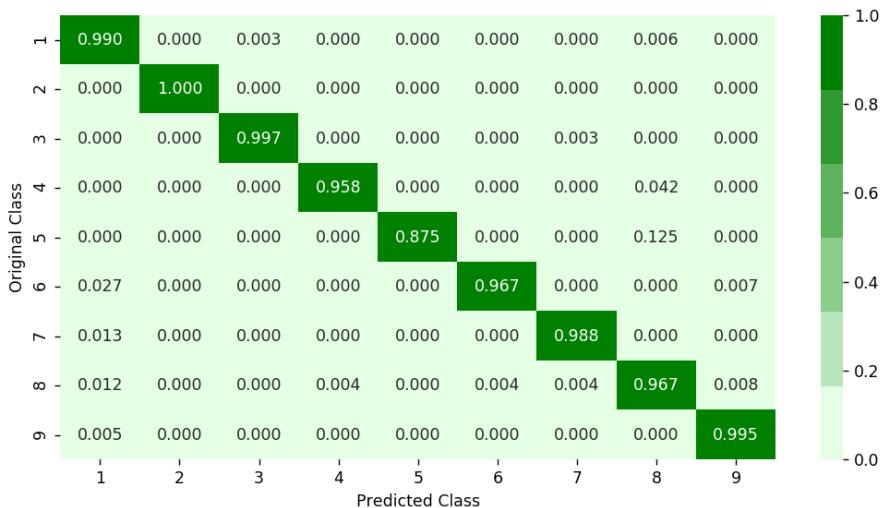


----- Precision matrix -----



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

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



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

4.4.4 XgBoost Classifier

In [0]:

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

# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# 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_weight=1,
# # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
# # reg_lambda=1,
# # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
# )

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep])      Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video Link2: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/what-are-ensembles/
# -----
```

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

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

best_alpha = np.argmin(cv_log_error_array)

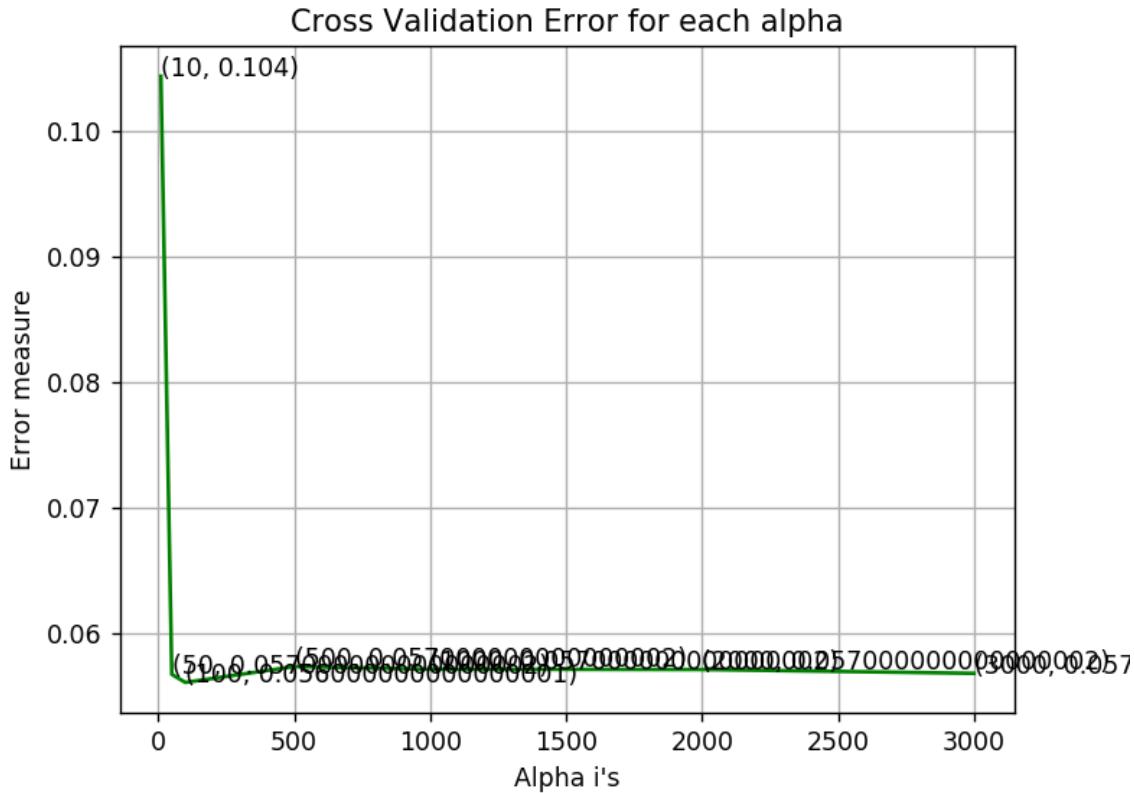
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
 ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

```
predict_y = sig_clf.predict_proba(X_train_asm)

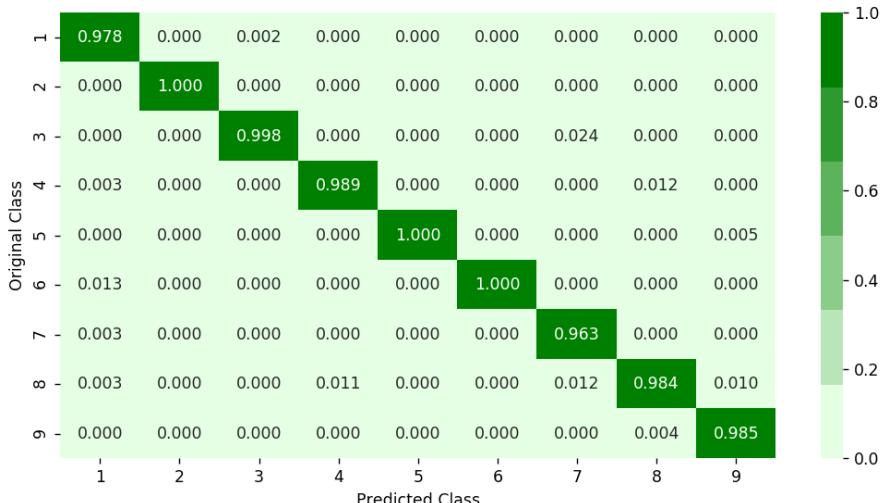
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_1
oss(y_train_asm, predict_y))
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
is:",log_loss(y_cv_asm, predict_y))
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_los
s(y_test_asm, predict_y))
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))
```

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

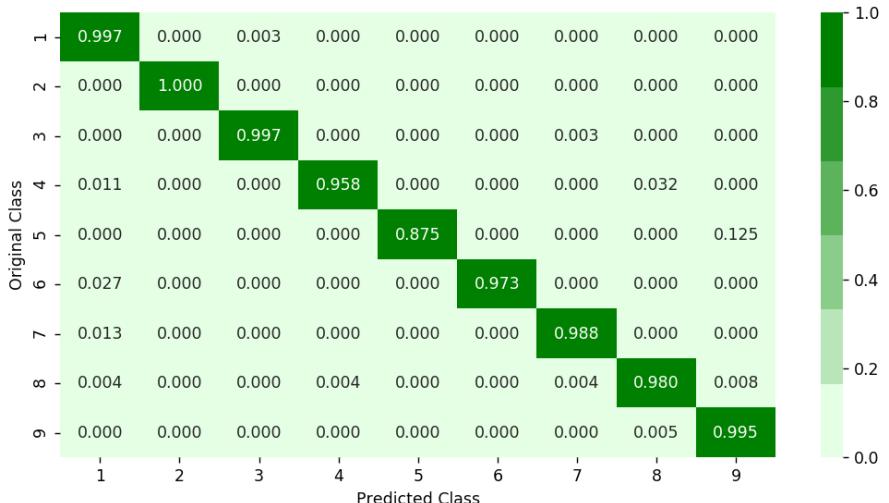


For values of best alpha = 100 The train log loss is: 0.0117883742574
For values of best alpha = 100 The cross validation log loss is: 0.056075
038646
For values of best alpha = 100 The test log loss is: 0.0491647763845
Number of misclassified points 0.873965041398

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

**Precision matrix**

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

Recall matrix

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

4.4.5 Xgboost Classifier with best hyperparameters

In [0]:

```
x_cfl=XGBClassifier()

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

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

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

Out[0]:

```
RandomizedSearchCV(cv=None, error_score='raise',
                    estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, colsample_bytree=1,
                                            gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
                                            min_child_weight=1, missing=None, n_estimators=100, nthread=-1,
                                            objective='binary:logistic', reg_alpha=0, reg_lambda=1,
                                            scale_pos_weight=1, seed=0, silent=True, subsample=1),
                    fit_params=None, iid=True, n_iter=10, n_jobs=-1,
                    param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.
15, 0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'colsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
                    pre_dispatch='2*n_jobs', random_state=None, refit=True,
                    return_train_score=True, scoring=None, verbose=10)
```

In [0]:

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

```
{'subsample': 1, 'n_estimators': 200, 'max_depth': 5, 'learning_rate': 0.1
5, 'colsample_bytree': 0.5}
```

In [0]:

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

# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# default 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_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
# s)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep])      Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video Link2: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/what-are-ensembles/
# -----
```

x_cfl=XGBClassifier(n_estimators=200,subsample=0.5,learning_rate=0.15,colsample_bytree=0.5,max_depth=3)
x_cfl.fit(X_train_asm,y_train_asm)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train_asm,y_train_asm)

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

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

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

4.5.1. Merging both asm and byte file features

In [0]:

```
result.head()
```

Out[0]:

	ID	0	1	2	3	4	5
0	01azqd4InC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835 0.0
1	01lsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873 0.0
2	01jsnpXSAlg6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280 0.0
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354 0.0
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232 0.0

5 rows × 260 columns

```
◀ ▶
```

In [0]:

```
result_asm.head()
```

Out[0]:

	ID	HEADER:	.text:	.Pav:	.idata:	.data:	.bss:	.rdata:
0	01kcPWA9K2BOxQeS5Rju	0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084
1	1E93CpP60RHFNiT5Qfvn	0.096045	0.001230	0.0	0.000617	0.000019	0.0	0.000000
2	3ekVow2ajZHbTnBcsDfX	0.096045	0.000627	0.0	0.000300	0.000017	0.0	0.000038
3	3X2nY7iQaPBIWDrAZqJe	0.096045	0.000333	0.0	0.000258	0.000008	0.0	0.000000
4	46OZZdsSKDCFV8h7XWxf	0.096045	0.000590	0.0	0.000353	0.000068	0.0	0.000000

5 rows × 54 columns

```
◀ ▶
```

In [0]:

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

(10868, 260)
(10868, 54)

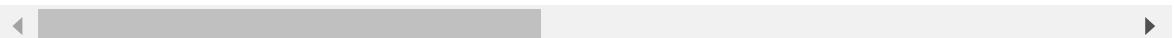
In [26]:

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

Out[26]:

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

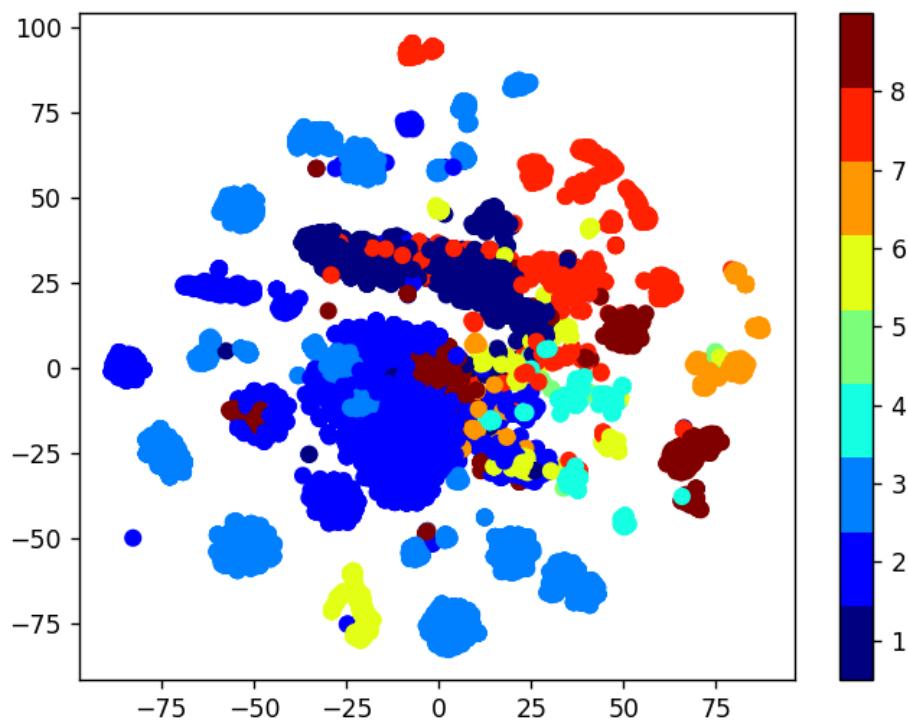
5 rows × 307 columns



4.5.2. Multivariate Analysis on final features

In [0]:

```
xtsne=TSNE(perplexity=50)
results=xtsne.fit_transform(result_x, axis=1)
vis_x = results[:, 0]
vis_y = results[:, 1]
plt.scatter(vis_x, vis_y, c=result_y, cmap=plt.cm.get_cmap("jet", 9))
plt.colorbar(ticks=range(9))
plt.ylim(0, 100)
plt.show()
```



4.5.3. Train and Test split

In [0]:

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

4.5.4. Random Forest Classifier on final features

In [0]:

```

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

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

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

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


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

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

best_alpha = np.argmin(cv_log_error_array)

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

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

predict_y = sig_clf.predict_proba(X_train_merge)

```

```

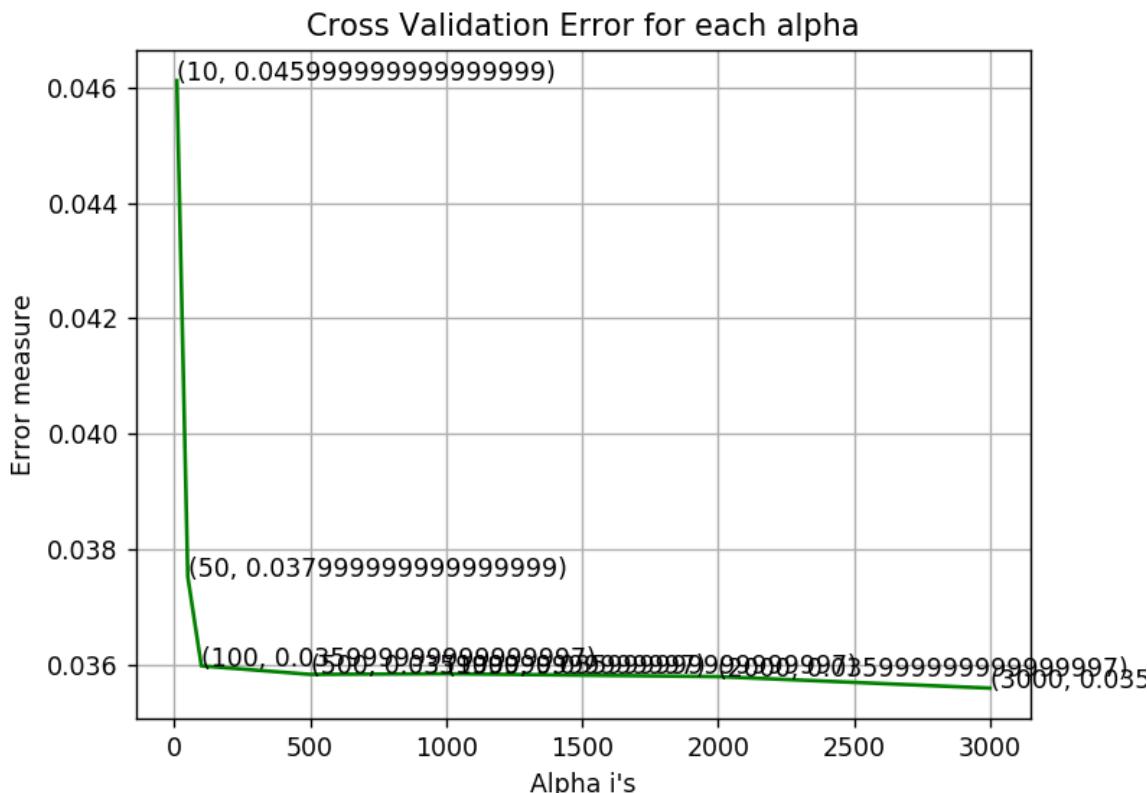
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_1
loss(y_train_merge, predict_y))
predict_y = sig_clf.predict_proba(X_cv_merge)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss
is:",log_loss(y_cv_merge, predict_y))
predict_y = sig_clf.predict_proba(X_test_merge)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_los
s(y_test_merge, predict_y))

```

```

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

```



```

For values of best alpha = 3000 The train log loss is: 0.0166267614753
For values of best alpha = 3000 The cross validation log loss is: 0.03559
09487962
For values of best alpha = 3000 The test log loss is: 0.0401141303589

```

4.5.5. XgBoost Classifier on final features

In [0]:

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

# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# 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_weight=1,
# # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0,
# # reg_lambda=1,
# # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
# )

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
# get_params([deep])      Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video Link2: https://www.appliedaicourse.com/course/applied-ai-course-online/Lessons/what-are-ensembles/
# -----
```

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

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

best_alpha = np.argmin(cv_log_error_array)

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

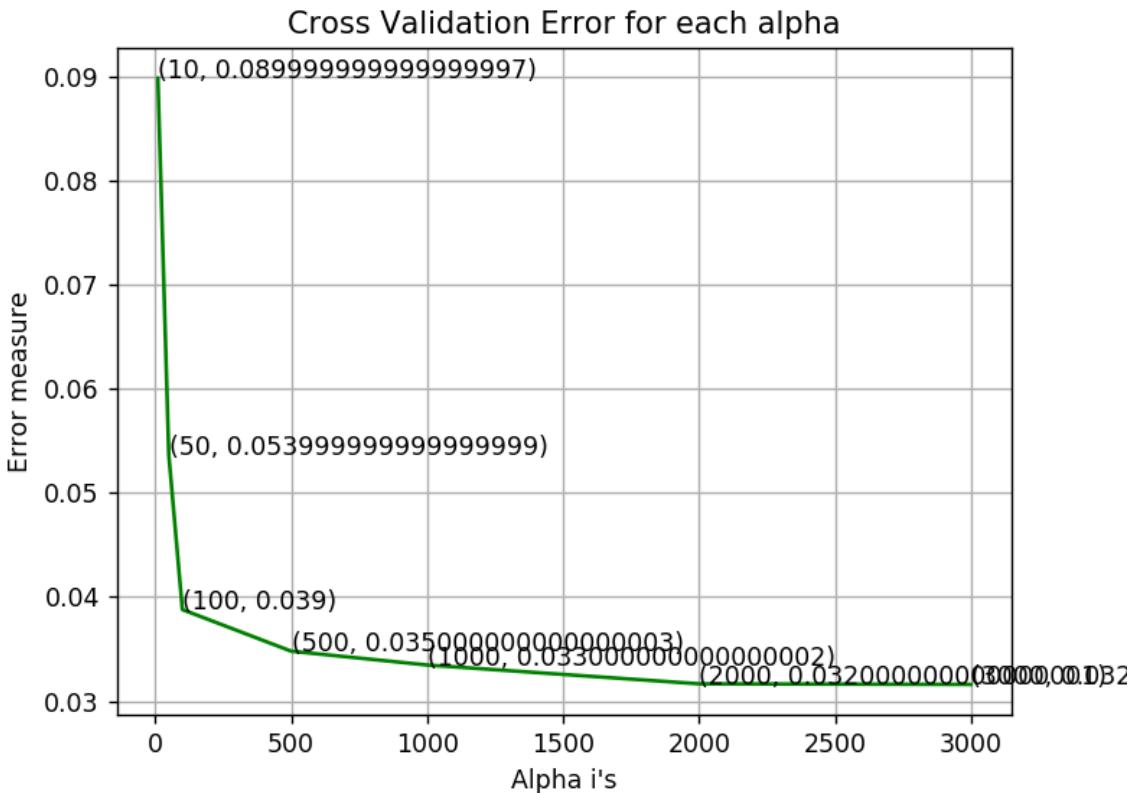
x_cfl=XGBClassifier(n_estimators=3000,nthread=-1)
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

```

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

```

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



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

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

In [0]:

```
x_cfl=XGBClassifier()

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

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

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

Out[0]:

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

In [0]:

```
print (random_cfl.best_params_)

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

In [0]:

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

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

x_cfl=XGBClassifier(n_estimators=1000,max_depth=10,learning_rate=0.15,colsample_bytree=0.3,subsample=1,nthread=-1)
x_cfl.fit(X_train_merge,y_train_merge,verbose=True)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_merge, y_train_merge)

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

For values of best alpha = 3000 The train log loss is: 0.0121922832297
For values of best alpha = 3000 The cross validation log loss is: 0.03449
55487471
For values of best alpha = 3000 The test log loss is: 0.0317041132442

5. Assignments

1. Add bi-grams on byte files and improve the log-loss
2. Watch the video ([video \(https://www.youtube.com/watch?v=VLQTRILGz5Y#t=13m11s\)](https://www.youtube.com/watch?v=VLQTRILGz5Y#t=13m11s)) and include pixel intensity features to improve the logloss

1. you need to download the train from kaggle, which is of size ~17GB, after extracting it will occupy ~128GB data your dirve
2. if you are having computation power limitations, you can try using google colab, with GPU option enabled (you can search for how to enable GPU in colab) or you can work with the Google Cloud, check this tutorials by one of our student: https://www.youtube.com/channel/UCRH_z-oM0LR0vHPe_KYR4Wg (we suggest you to use GCP over Colab)
3. To Extract the .7z file in google cloud, once after you upload the file into server, in your ipython notebook create a new cell and write these commands
 - a. !sudo apt-get install p7zip
 - b. !7z x file_name.7z -o path/where/you/want/to/extract

<https://askubuntu.com/a/341637>

Got the code from

[\(https://nbviewer.jupyter.org/github/bazinho/Microsoft/blob/master/byte2gram_CountVectorizer-SVM_Linear-0.8980acc.ipynb\)](https://nbviewer.jupyter.org/github/bazinho/Microsoft/blob/master/byte2gram_CountVectorizer-SVM_Linear-0.8980acc.ipynb)

In [38]:

```
import glob
import csv
import os
import pandas as pd

import matplotlib.pyplot as plt

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import StratifiedKFold
from sklearn.metrics import confusion_matrix,accuracy_score
from sklearn import svm

import numpy as np
from numpy.random import RandomState
np.random.seed(1)
```

In [6]:

```
flist = glob.glob("byteFiles/*.txt")
```

In [7]:

```
flist
```

Out[7]:

```
[ 'byteFiles/9MW5Nuf0ogCEcR1YJeKG.txt',
  'byteFiles/GFblksovaPYLI5XeC8RU.txt',
  'byteFiles/kjQFMnf7vADCqtX4ai2u.txt',
  'byteFiles/EKmgShQsf6a9vY0znNlU.txt',
  'byteFiles/1WCXg8qEtJFITMilaf6k.txt',
  'byteFiles/40dBpWDhOUzz5yVAKJF3.txt',
  'byteFiles/bLNRehaM1SUd19V68ugD.txt',
  'byteFiles/FfYe8u5S2rh1B0sDXIaZ.txt',
  'byteFiles/0u3DIJTqnppkg187VSb0L.txt',
  'byteFiles/jHfu9AMl1GkNWch0Yv53.txt',
  'byteFiles/4flnV7Xk0PyZ5ojuIiQc.txt',
  'byteFiles/HooEWQVSZyThrFj1tf2b.txt',
  'byteFiles/b5qZNQRzLp6MulfHisof.txt',
  'byteFiles/bPx08C4QDMk9nFzyvHTB.txt',
  'byteFiles/1D6IwN4EdfHta1YJuA95.txt',
  'byteFiles/cPo7FESsJQvVNqmaOiiY.txt',
  'byteFiles/iJpgS9hMDuFrCTItyok1.txt',
  'byteFiles/ByjYPtU1Q0T76FekSJH5.txt',
  'byteFiles/6WauKLy80DTdQf5hjZx7.txt',
  'byteFiles/G7DQnb0fcE62usXeAmIN.txt',
  'byteFiles/6QWh10fPM4EYicpGBv1L.txt',
  'byteFiles/Esg6ijhbT8W9tQmcaH41.txt',
  'byteFiles/kMDrx5LS2P864CIwNiVs.txt',
  'byteFiles/KawMHneL7VYj8oQ4m1cZ.txt',
  'byteFiles/gaYrbUuHQn9w5q10dMcm.txt',
  'byteFiles/6tIbeDRa0jfnnZYcFqVdL.txt',
  'byteFiles/3WfGisHe71kxVFN8cjXz.txt',
  'byteFiles/bG1zA9Iwqf6gCVSoJ0ZN.txt',
  'byteFiles/2XiGfo95mCALRbu6KzVY.txt',
  'byteFiles/gpQPsBn4Hc6uoxmFbYU3.txt',
  'byteFiles/1JiaRhjycfZb8BpmX4TM.txt',
  'byteFiles/EI4Jgyofj8CHR0ts59Ur.txt',
  'byteFiles/CDFlbvHczS19akf0EBwo.txt',
  'byteFiles/kKEvpQMViqzLfeTWDtF2.txt',
  'byteFiles/FNMk3wvliVuQLCe9OTDg.txt',
  'byteFiles/bF5KEHN7mXfvktq0QxzD.txt',
  'byteFiles/JVPnwpfeExcNDt2ZTAQm.txt',
  'byteFiles/i0JebqWZLA71XSkgc2N.txt',
  'byteFiles/fs9oYJKeZ4BzXS1Iwy3E.txt',
  'byteFiles/Aghn90rSzad2DvLYPtqV.txt',
  'byteFiles/dI89k371NKgURwFqs0By.txt',
  'byteFiles/FNKR36t4qTZoy5veUDRm.txt',
  'byteFiles/2r4503Hczp6togVfAJBI.txt',
  'byteFiles/2zrDU8tlpNjoWaisTeK9.txt',
  'byteFiles/E5FzXZIDnTmau6VWepvC.txt',
  'byteFiles/d9o6Is371a5eh0cQXtvY.txt',
  'byteFiles/JBwr0MNAVT4z13kvifYb.txt',
  'byteFiles/g6C1WcYBrAHskpjGXeFE.txt',
  'byteFiles/9gpKWS6osX5Fq8E4UCLr.txt',
  'byteFiles/fX0BiTsSxvUeIMnVha3y.txt',
  'byteFiles/31KjQiNWATfPaozHtY01.txt',
  'byteFiles/0P6tEopzr2mIhkuOKiCD.txt',
  'byteFiles/7twBm0QdZXenW3vcUFJ2.txt',
  'byteFiles/2eSoCjQzDri3E0cuxUYb.txt',
  'byteFiles/GmPLUX1Y2DQhan563Etf.txt',
  'byteFiles/bazOp8SLuGgisveCP0QI.txt',
  'byteFiles/HRhfKEv12kTVqm95NQ8o.txt',
  'byteFiles/bfcBD4va0tUJG6K5imqx.txt',
  'byteFiles/gXLSZjPsvM1Bd0nJGQhN.txt',
```

'byteFiles/J1YNWgvo2LCkQRV8d74u.txt',
'byteFiles/fdGV3HjRu4nBSMNibIkX.txt',
'byteFiles/jZxB2zDWto8UTu6hPENM.txt',
'byteFiles/a0tVNG6x3empWgK2Zcyd.txt',
'byteFiles/Jtp04ZqNWST65MdH7xhD.txt',
'byteFiles/euAfyh0HkFiaKndptmsJ.txt',
'byteFiles/Fa87QkVGAHec3PfuwbMm.txt',
'byteFiles/2e5P1iKzXugGRBFF0Enc.txt',
'byteFiles/2xX0y0GfquKYWIpsPa3n.txt',
'byteFiles/CG6sL3Aug4yomNnztEIw.txt',
'byteFiles/e9M8Qjno0aqN3HhLZkY4.txt',
'byteFiles/jNUAeECmnnsxGJ1QrTubv.txt',
'byteFiles/0bjN3Kgw5OATSreRmEdi.txt',
'byteFiles/gVzZU8WGxyoR07FesaJl.txt',
'byteFiles/esQpnxYJW1PVrT8ovu43.txt',
'byteFiles/Ga17B5qZe6VoM0kCxTAF.txt',
'byteFiles/g46jFCHzY1oqtws0yx05.txt',
'byteFiles/c6tLU50oSs7hqI9kmQVM.txt',
'byteFiles/9Ur08SqELYFnQAHGva7z.txt',
'byteFiles/3G6jJKRVLz1poimkvfhA.txt',
'byteFiles/FsHojWp48ESXDbcxzK0Q.txt',
'byteFiles/50iq3CaS7H1PDWetAMLR.txt',
'byteFiles/H6pGwSYfNmZVd80aeWg3.txt',
'byteFiles/0u4QcXSvNvYnMa10BCygH.txt',
'byteFiles/0sMxy5SNLu1YCkKJzaB3.txt',
'byteFiles/Iec3SYD4Gimhob1xQ7sV.txt',
'byteFiles/3EdJ5pSORwmeINfjag1Q.txt',
'byteFiles/Dsg960KfZAq5ynarhiwT.txt',
'byteFiles/6teIPpCrYKw0zFusDSix.txt',
'byteFiles/ieFMEkfnc1pmS8taKHo9.txt',
'byteFiles/juPCJRHOoLngNAvBZ7x1.txt',
'byteFiles/ICleA7H05Kaf4jMXZG2y.txt',
'byteFiles/6tmqe5B8JDsQnyhGo7vP.txt',
'byteFiles/bROPgBf8pYcSud4Ejem1.txt',
'byteFiles/imns4fDL9QbS7y60hwac.txt',
'byteFiles/CkOT69iQsaLAGq8mYEeV.txt',
'byteFiles/1eJx3418pcAFvMu0wrjB.txt',
'byteFiles/8UGK0mbY3De124d7wizL.txt',
'byteFiles/IweQoKz4HxGc8mAETNpM.txt',
'byteFiles/6QWNYDyohV3R11fMCZI5.txt',
'byteFiles/G29t8jq3S5usQkfeL4zn.txt',
'byteFiles/F0ScZ1qt1RyYEQTVK4XH.txt',
'byteFiles/9itfPnyplHzhTb8IMYq3.txt',
'byteFiles/9CVbzaWsLvcHjPqZu2Q3.txt',
'byteFiles/Huz3U7hA5kVnJ6q2QcgM.txt',
'byteFiles/dS6yJEktR1WQsBaIFYCp.txt',
'byteFiles/fEMdwkrczmt7j0RTL14J.txt',
'byteFiles/ITZQmUtjYaHDxMonFrKC.txt',
'byteFiles/FNAVfQgh4a6k3lUY5ZcK.txt',
'byteFiles/60of29sarJPj8GEbnkTS.txt',
'byteFiles/DwyNaEk0S4769ocfPOzn.txt',
'byteFiles/9pSnLG2tRJTQ6C1Hzqfg.txt',
'byteFiles/G31czXvpnwUfRtdJ4TFs.txt',
'byteFiles/6DmCqAs0tFKxGe5aPz1U.txt',
'byteFiles/0qNtRXvP8F2431bMmpse.txt',
'byteFiles/jWN2deu0SfDhU4cR6XVA.txt',
'byteFiles/J1hb9akq5NpXnvIg4DG8.txt',
'byteFiles/jCAmu6J5gIBHU3TxopMD.txt',
'byteFiles/A4xfRz73qWePOCnJ2Sya.txt',
'byteFiles/aC3qABWUp29F0XvZNkIY.txt',
'byteFiles/0pT04SVnWDehgU1YBAvq.txt',

'byteFiles/HRE9TjDJgZl6LcuV7yKn.txt',
'byteFiles/4WaGm3PhidQVRopgXz7J.txt',
'byteFiles/exSWamc4Hi0j2CG5Rq1u.txt',
'byteFiles/3v7fe0Diw9c0u6dAzsaE.txt',
'byteFiles/5dZSWcADY0fCq8hsuJ1L.txt',
'byteFiles/JajfhmMDxzBX80Kb9o4k.txt',
'byteFiles/19emSiJG8vrZgT2zXEIY.txt',
'byteFiles/8J6QdyfRgIeTnFsNPKDp.txt',
'byteFiles/3rf7KJiPd8hIZ645wBVN.txt',
'byteFiles/BxRn5YjcIH1qCEy67NP9.txt',
'byteFiles/BUGAnop1NDYTh0txE8d6.txt',
'byteFiles/2Z1LcMw4Vkim1pNYXj9v.txt',
'byteFiles/a6yPugSArOdWK1zQ0c7I.txt',
'byteFiles/2sf3AQd9MjP1xcLFRYyo.txt',
'byteFiles/8GilsX4Uk7xQdmz3wgyp.txt',
'byteFiles/Bnwd1vNx85KI7yzphZHs.txt',
'byteFiles/Fh9HCK50ZId4PbcQDTJr.txt',
'byteFiles/Jgtyz90K13M2CYkmZolp.txt',
'byteFiles/fvbowPJXxAKQsVt7MiRL.txt',
'byteFiles/ig4IwJckB1RaT3FUN6rV.txt',
'byteFiles/Epn4owBqFXM9YRyf31Ui.txt',
'byteFiles/5p2yY0MvAQGXSIuocmbn.txt',
'byteFiles/Jo4mF0en982QIKTAboXN.txt',
'byteFiles/KWg8T3NOYaQBF0tMqlm5.txt',
'byteFiles/gC2GsRX3ewpQyvj00Hb4.txt',
'byteFiles/9DVvn8GrbQRsJwXCWcSo.txt',
'byteFiles/D9QUL1KPYijhJezwSOFT.txt',
'byteFiles/4M9zLZjcgVbWiyINDEGp.txt',
'byteFiles/fSWvxIsOaV0YjkB185ut.txt',
'byteFiles/4HGx1XJ8Ncy2Ri5az3o7.txt',
'byteFiles/3iHvsmdx7DauhSPKfec4.txt',
'byteFiles/j9ViIqQ4RXTm7hSFwb1W.txt',
'byteFiles/hMyn4mwJjsLI5103PdV8.txt',
'byteFiles/BLUqHvpdaytR0E3rZMu6.txt',
'byteFiles/dLaQic9uU0G2SeRXDAHw.txt',
'byteFiles/0X6U3SfcmbuZNodnH9qL.txt',
'byteFiles/j04f0HCJN12TkB31PVyv.txt',
'byteFiles/bf6w4apWBzyGLZhCqK7s.txt',
'byteFiles/GBNdXiuAb7Cr05vVoK9I.txt',
'byteFiles/IyV1H08cxzfMLvjFbr0X.txt',
'byteFiles/2mYSX6J9Dd51kpLuNMc8.txt',
'byteFiles/D4wXBShNamu3JLz1p8M0.txt',
'byteFiles/EH7DRYeINXPfJUAG0Lqx.txt',
'byteFiles/igjezBrTh4SbGtIKZmlA.txt',
'byteFiles/IemKCvJBb7ia02owVYHN.txt',
'byteFiles/KpzMHfmGQJbXN92Sdsyx.txt',
'byteFiles/8A1wMUgpDrZ6HTcmQoCx.txt',
'byteFiles/Hov1QtF67wa9ELyuJmnS.txt',
'byteFiles/G1BR2jy6Nx1HUEag5ISm.txt',
'byteFiles/FyIcXN7Z6SQojdgntUD9.txt',
'byteFiles/HhFctSgYNdwDxKLnAJC.txt',
'byteFiles/8a3GEC7cPKtx01I6qYvL.txt',
'byteFiles/7dYo5p6fSV0rkqgLvzbT.txt',
'byteFiles/hnEAHVz5I7tQcjGWvRi1.txt',
'byteFiles/10H9detwuizByNSRMmxr.txt',
'byteFiles/i3HsbL10XJ2qTQPY9Cto.txt',
'byteFiles/DZYefdjFtAGU91KVkJHw.txt',
'byteFiles/ewsfhA3D unp49SFNM oVR.txt',
'byteFiles/iCPn8QUoEZ14IH p3jSkM.txt',
'byteFiles/FgHkyL3ot79E1TmIZcbS.txt',
'byteFiles/hTOJKwgGvBxb3dS8e9cn.txt',

'byteFiles/EdT3lQiGHbzmrX2L1NgD.txt',
'byteFiles/8d7zkmPAEtQrBiTRNlg6.txt',
'byteFiles/5mcXAfSQWlV43hb6Ns1d.txt',
'byteFiles/CLZ15SyTFMUBH3ldR427.txt',
'byteFiles/0hWRb28Umdgj7xcX0wtC.txt',
'byteFiles/6prXgmtMVs7BYvPL3ce5.txt',
'byteFiles/ElbIP3i6rgW41R19phnZ.txt',
'byteFiles/3o1eL8jrlyb7MY6NASQi.txt',
'byteFiles/e1KAXRxrOTc5snpuWvj.txt',
'byteFiles/73cMCHeUwStZjy2VoQRm.txt',
'byteFiles/g5NnsVCLJdb1lzPBEUeT.txt',
'byteFiles/6WKJVBfk5INU2q7RLDhY.txt',
'byteFiles/B0dnNijSbH3fmuCQIL18.txt',
'byteFiles/b4I670cxsfHSWlV8rJiEY.txt',
'byteFiles/hdYWDuVTnF3kPv9CI4Er.txt',
'byteFiles/f20PUJcTG4D0FRyh7Z1b.txt',
'byteFiles/502xXcWOnqet1TrUmPBi.txt',
'byteFiles/C7cZvXSQWwI0ftjA8xLN.txt',
'byteFiles/Ek1sLKWN2DoUFriJy7mC.txt',
'byteFiles/9VZ673sDcNIJKk8vP0TE.txt',
'byteFiles/CKIT85YvdjSMoQr09A2L.txt',
'byteFiles/GSYtwpD0MfegOyzNUHXr.txt',
'byteFiles/iypLG5f97Y2ZTKcBuUxW.txt',
'byteFiles/9FxgivRPntUB4YcV5IJa.txt',
'byteFiles/l5zY6vNZbOKF07uSn2Am.txt',
'byteFiles/JFiG1dVRDj6MrH3fAK7L.txt',
'byteFiles/gfSamebdC3Z4NFt5lonM.txt',
'byteFiles/alf0ZAc7xULvhP5Qsr4Y.txt',
'byteFiles/EYJryk0oKhPwRNbWuIUz.txt',
'byteFiles/hGJZKoMdCW9s6SEqNn0g.txt',
'byteFiles/I9wXdHpzEAqKRo2TUtDN.txt',
'byteFiles/23fQEb7B6W01VDIA4Sw5.txt',
'byteFiles/81IgT57cy1MUNGbLkiqm.txt',
'byteFiles/g1Uw2y9I4EX83msAMQcK.txt',
'byteFiles/H75RL96JxmMEfBvbe2jk.txt',
'byteFiles/gZRseTmhYfHPvEnVFDiq.txt',
'byteFiles/ADMm3H8ELRwcpUgjId6K.txt',
'byteFiles/jxOWbXpVtGuyQI5DSUCB.txt',
'byteFiles/84jKcbktfsh169TiErUG.txt',
'byteFiles/IjC7vkJo1dViZerbwxT0.txt',
'byteFiles/fYpFPV0nCAKS80qRtiH9.txt',
'byteFiles/C9hHuINUVJqk1zo5pTQX.txt',
'byteFiles/hNboYglQSecVZ1LntIGf.txt',
'byteFiles/fh1V8wnrOJAZGzsRQB60.txt',
'byteFiles/ANwkGK2YbhD3VjsvHrLQ.txt',
'byteFiles/1nyRDJYtOrpTEI3CmGHf.txt',
'byteFiles/ChOGatx1eA8JNEVqFKTI.txt',
'byteFiles/HpSJQCGdyYItq0eUuo3A.txt',
'byteFiles/hnCWyLtgETF1kJ84DqXS.txt',
'byteFiles/BMyRb1VjoWUws42eSLgc.txt',
'byteFiles/J058MAd10WiuFk9zhspH.txt',
'byteFiles/jEhpv08m7FOyKwXYoSs4.txt',
'byteFiles/c1QNt9kyP5Tw8SDVfBpE.txt',
'byteFiles/2oHOw4StzEV0hiKq3ByD.txt',
'byteFiles/iK503zvsIBnMQAwmgHux.txt',
'byteFiles/H7tC0wf5rzSluXiGsx8W.txt',
'byteFiles/6DvOLIs4i1Gk7P9dEeTo.txt',
'byteFiles/35rbwhMGTxqAWjuFt20m.txt',
'byteFiles/croUvXx8jI6DARFmhWJ1M.txt',
'byteFiles/an7uVbkw5MAPJf8oX9GQ.txt',
'byteFiles/BH561pgQn32eUyJDFImv.txt',

'byteFiles/d4lM2k7aq1RwJuTpg06X.txt',
'byteFiles/aZoe13jMRHwBEG4X768N.txt',
'byteFiles/0yvBxSufKXzM6CwpcSYU.txt',
'byteFiles/jsqeT23MYgH7ytdxDUki.txt',
'byteFiles/byKBxZ7dNkLanQtpiJHz.txt',
'byteFiles/7LbrmXchDZJ1ed5HnAfM.txt',
'byteFiles/6a7hF1ZnHrUAQXouv4x0.txt',
'byteFiles/0ASH2csN7k8jZyoRaqtN.txt',
'byteFiles/47GqvWtCkjShUVpE1oI0.txt',
'byteFiles/a3ychUzAELtHgxI7kSoC.txt',
'byteFiles/CM53GutBya9do7piSRe0.txt',
'byteFiles/kJ1SuTtgh6fP9vas7wCH.txt',
'byteFiles/5QAGY1NrgvSIeRxDH294.txt',
'byteFiles/cuNE2eL7XtkbJQw50zsW.txt',
'byteFiles/HpMdj30s7EDPkFLWvYy.txt',
'byteFiles/FDHjobQnCUSs6atXldJL.txt',
'byteFiles/0qPGt4cRVk9NoiJgubf2.txt',
'byteFiles/EZWPt7i0ocxdnqG3BvRD.txt',
'byteFiles/fEnNJU3pLTbSOeFxaQ0z.txt',
'byteFiles/jgpTUqt96Bv7i4xcIGzF.txt',
'byteFiles/Bv7fxdPDSZKmAXGp8baQ.txt',
'byteFiles/GD7JUv1IfRuHpneMKaYc.txt',
'byteFiles/5F4mbeBZsfwYCQpGgVLa.txt',
'byteFiles/ESqUten2YJiswjZodR40.txt',
'byteFiles/1PQFYMSBLA09TmKk2Zhj.txt',
'byteFiles/JiUyEfTOm6oasXKj7p1G.txt',
'byteFiles/ecvZowmJ09xHPGkBjCXV.txt',
'byteFiles/0aNj3qFgEZI6Akf4Kuv.txt',
'byteFiles/GqPZCULjirRngSyMwmh4.txt',
'byteFiles/fjbqGwVH6ANh5CRluirI.txt',
'byteFiles/Cdj3sjrGF0BaYTx7biRQ.txt',
'byteFiles/bye3ZNnLt4BfuqAJ08v9.txt',
'byteFiles/6SPDXoT29nxidIlhW1VC.txt',
'byteFiles/Fw3QCzKGsD9e86pgLiIm.txt',
'byteFiles/3Vhmj45EaPbB60rXoIMw.txt',
'byteFiles/e1Pxcv2KRUSyDWmAsid.txt',
'byteFiles/DNtJKyn0WV7PeS3q8Fv0.txt',
'byteFiles/efKbsRXGgy15NPHW8tI4.txt',
'byteFiles/EBghIvP2MTJU6GXkAR0m.txt',
'byteFiles/AbDuHJy7LUcZWh0aFEXh.txt',
'byteFiles/fQmbIDTrZBds8Uu45HR0.txt',
'byteFiles/Dw4btEPmk3reBuNsZ15S.txt',
'byteFiles/IwgDdxrbUkiXpfAMy6J1.txt',
'byteFiles/I2ZH0sW5hT4aDKwo9gR3.txt',
'byteFiles/0Fu9oETtMW4zlg1ZrUy6.txt',
'byteFiles/4wIxhDlpCbJVkZLfHKtP.txt',
'byteFiles/G2NCXj6ZW0yTFALv7Hqd.txt',
'byteFiles/dV8HAvgFer2naZEwu1Pk.txt',
'byteFiles/BGpm6FzEHcn5SujJikwy.txt',
'byteFiles/FBHvx5A1IchTEYDPkU0V.txt',
'byteFiles/7u180DTFRYWPrM0xfLd.txt',
'byteFiles/JuAyLarnUdx5cRYT1v7M.txt',
'byteFiles/4aMbnH1A9es8vKmOFR2Z.txt',
'byteFiles/HeBQriW1XgyVJcI31Y7f.txt',
'byteFiles/dZB602DGWj5pkLCtzgv0.txt',
'byteFiles/fcsNBbd6HYTaQAh5Mmq3.txt',
'byteFiles/IopqlNQASCdjT8e0WRHw.txt',
'byteFiles/0Hrfce4X5YGESJPjl9uL.txt',
'byteFiles/HJE15j7Rz00garybNnh3.txt',
'byteFiles/8djFkIgPK56XtYGDsb7R.txt',
'byteFiles/DwEQhKiobyXvg81rjpRe.txt',

'byteFiles/KIdHiQmorRG90vEcnbp7.txt',
'byteFiles/409FLy5cim8qgxjYobET.txt',
'byteFiles/Jy2Nw4xD3bPilQfRsW5m.txt',
'byteFiles/5XPfYvNdQRIBJowmj30s.txt',
'byteFiles/EzxVkJ2DCYKwZpng39AoS.txt',
'byteFiles/7aEbXfVnzhG6dRMLT58q.txt',
'byteFiles/DoEYOjCZuGA3x4JSgW5d.txt',
'byteFiles/0NXFnJyOEHBAISKfiU67.txt',
'byteFiles/0SaYqdFRA5G7rfktjT63.txt',
'byteFiles/awJH283VpYm7yBDxb0AI.txt',
'byteFiles/jM0aVyQfdqBhODSox1Ce.txt',
'byteFiles/1eCQkhvgDSFNXYLJbAdj.txt',
'byteFiles/hdJnQNkT6b4yDeVwaR3F.txt',
'byteFiles/5PoseIX8Ak7x24hKQuNE.txt',
'byteFiles/bMK0Xukv9xUAi2QSGJ6Z.txt',
'byteFiles/JFGA1Ypxt56ElPa3WUOQ.txt',
'byteFiles/aE6dSw0lys08Qn4cf2FM.txt',
'byteFiles/1xGM8wcYTV4pQFhJBd2f.txt',
'byteFiles/7pbVgdSnjUPTTrwcR5GXu.txt',
'byteFiles/KAЕеGNIw967T5pHFhMQd.txt',
'byteFiles/3L9crexw581VMXajnsmB.txt',
'byteFiles/Fkha6WrBdNbMje70q0T4.txt',
'byteFiles/bHT5pg8MuXLf2UZAwxBc.txt',
'byteFiles/BtLTcn2E0adFGRPsKqYz.txt',
'byteFiles/19zYbuW3XONcEedv7xU1.txt',
'byteFiles/4syoiE9gr1HVtAfwTQY8.txt',
'byteFiles/7PNjmAMet9183cVEDa1u.txt',
'byteFiles/K1cOh76tkCvbNxPfLmD.txt',
'byteFiles/cBTaW6V7hrS5fKgJvwjA.txt',
'byteFiles/e9Kji7QuZFdN1Lmc0Y03.txt',
'byteFiles/7dk9aUmW0HlvsbtNSFuZ.txt',
'byteFiles/bL63tFHUFNwMiAuvCV51.txt',
'byteFiles/KRNHAм094TC70JfEPp8h.txt',
'byteFiles/aXbDi405xeqQyn169Yk2.txt',
'byteFiles/AKzFf5dE2PxtL0YqXv7G.txt',
'byteFiles/A4oRW0sZJ3jkWazH5TLy.txt',
'byteFiles/DcgUNDMzXlah1kbxiYtv.txt',
'byteFiles/C6pcTZjgKOMtB0iREaSe.txt',
'byteFiles/JCzVMIhYuK13LZgo6Nct.txt',
'byteFiles/ABNUEpLD8G1K6iubzwWe.txt',
'byteFiles/1mStzMaQ8RUVBx7TkF4o.txt',
'byteFiles/3UV6Iek05Y2pKRB7D1g9.txt',
'byteFiles/FdPwzJ95fRp1LGceAEWI.txt',
'byteFiles/JrsqMwaubVkJ3vj9him0.txt',
'byteFiles/CZnRQNX0g8IGEcMq1lt.txt',
'byteFiles/b71ZgrI94Bp3XM0f6vuJ.txt',
'byteFiles/hjESsM1TOUgAXkLzzirr.txt',
'byteFiles/5LFWPi9qzEb6fhyVSYaw.txt',
'byteFiles/BXQNYRIMHrq8LZ7Ay9pc.txt',
'byteFiles/4ETQ9c3B0J7YKNpFt0Ln.txt',
'byteFiles/EjnR6Llw4p7sQkJWYiBq.txt',
'byteFiles/gtiofz9cBwuYCHjyJTdK.txt',
'byteFiles/dQbqMtOsVYrWDIVhv8k2S.txt',
'byteFiles/HNApyaXk5rm6g1tDcnoj.txt',
'byteFiles/2MgFXNYimTJpIzKVE5wn.txt',
'byteFiles/dAabnhJ37zorIS068q41.txt',
'byteFiles/25jQ4xItcnbB6Hfo0E9T.txt',
'byteFiles/HD1W2pcn973X5YaQVBtv.txt',
'byteFiles/bv60XkYZsCKHa3NgFUDr.txt',
'byteFiles/iks3Y5jQmILCpn806wR7.txt',
'byteFiles/iTEz2rB6XdsZIG8g071K.txt',

'byteFiles/dNw1IRg58zQ10oFVJS2q.txt',
'byteFiles/DHGMN38qz61rTmpSXsbn.txt',
'byteFiles/59npdcNYTD3qAvIPyQMs.txt',
'byteFiles/eQVmTuqlHYXLJ6DBCsg5.txt',
'byteFiles/9d3mFhrqe5gtuLE0AC0f.txt',
'byteFiles/2qZrQoxyYDeadjIp3mSH.txt',
'byteFiles/4gfR0YmZwOW8NvLSiXn1.txt',
'byteFiles/hgbBPkqK8S0I7v0FDYtw.txt',
'byteFiles/5YEafCpxyFcNqDG18TS.txt',
'byteFiles/deTXH9Zau7qmM0yfYsRS.txt',
'byteFiles/7waVJqpLnoxz4fKsZv39.txt',
'byteFiles/DQqSM2enGHopTfCbg90c.txt',
'byteFiles/h5uGZKiToFpXSrkcAR1B.txt',
'byteFiles/bNT4jk68t5lZghy3HPcr.txt',
'byteFiles/IVMFepgX5rASctxiL7y8.txt',
'byteFiles/7xwMNplsj3EhBJ0Zkd9Q.txt',
'byteFiles/4Ub1A5oWQLNmdVF8sDw0.txt',
'byteFiles/eKgQ10AwEzo12DBcs46I.txt',
'byteFiles/4RKTw3jA6pyNbV2lgUst.txt',
'byteFiles/GU32g0fBF91xp7JqNXHY.txt',
'byteFiles/FIVR08js5sZo6avbmPHk.txt',
'byteFiles/5ov2DWngteLM0GVEySrl.txt',
'byteFiles/4Qcv9hWIsmBX3PUNLltT.txt',
'byteFiles/6LJuRVBNsdWGlaC7KbqQ.txt',
'byteFiles/hmxjukz4UBqMPf9Jec6y.txt',
'byteFiles/1Flm7qajn0zdtwBZ29gL.txt',
'byteFiles/CfSF2urLkdAX6MbzQiIa.txt',
'byteFiles/gCBJQKMq14Atfe3ZSRX9.txt',
'byteFiles/FIxgyKhla8MGtcaf7Xur.txt',
'byteFiles/iI3TgAon2970drqY6CB0.txt',
'byteFiles/ADaWiK4BCZg7Gub09xjm.txt',
'byteFiles/40DnoUH5ets62Na1yVx.txt',
'byteFiles/6qLobnwUp07W1RyKBjvZ.txt',
'byteFiles/jw6dqJau5vzf7Vylsb9D.txt',
'byteFiles/JchT00vFpwRX9beS3i2k.txt',
'byteFiles/cqHlrY9oAVpyWMKJ8mOF.txt',
'byteFiles/2Z0Jfu1Pd8Mq7xXNSDct.txt',
'byteFiles/ETaX1605B0Hm7Ck3tGVF.txt',
'byteFiles/40as7FtBN0GZTu05XpWq.txt',
'byteFiles/2WHevmqLEZAb5ug3I7DN.txt',
'byteFiles/49RfyxpQP1CzvneWIG5q.txt',
'byteFiles/JvsW2yaUokpXVcbh10Mu.txt',
'byteFiles/AYVnhTKJPbW2eBjrQ3Cc.txt',
'byteFiles/6WKbReFUk18wYnCNQsyH.txt',
'byteFiles/IYFTBN6LCn7SEQgbPOaU.txt',
'byteFiles/HE9280WkP4RJtdsNL7bZ.txt',
'byteFiles/ASEw3L0HqWMePiy1uKxc.txt',
'byteFiles/7QPJ36MV4rktjX0H9GLq.txt',
'byteFiles/8BYckA0Zu47Rjm3bgUzp.txt',
'byteFiles/DVsfp6pU89nBdhZHqxI5.txt',
'byteFiles/ATr6NemOyI2KF0w48BR5.txt',
'byteFiles/j3A60bUBPWgX7Qf2wo8s.txt',
'byteFiles/izS2yBV3Z4nKch50mAuo.txt',
'byteFiles/JFiWOaubTAKCyRoBqzfn.txt',
'byteFiles/FZsf1RA4XawTS7Qj0t2N.txt',
'byteFiles/C5RGotXvrOxI3Np8g10A.txt',
'byteFiles/HsL2jYueFUSCJZbPKzkn.txt',
'byteFiles/gdMSLJ5UZmtPpN8W3ja6.txt',
'byteFiles/jnEoBxGXRksZqHw8euhd.txt',
'byteFiles/AUMy0gtvXYK74RT163ie.txt',
'byteFiles/d0iY69mAf1JzMSr3vkcr.txt',

'byteFiles/0FKerJl18x0c3jdoyg4A.txt',
'byteFiles/FjeNJD8d04AakcKu5qCp.txt',
'byteFiles/GEprqj7xuTa8m629zQP3.txt',
'byteFiles/iSaFcugoKwRrmPTDZ0yJ.txt',
'byteFiles/CQS2WYL40q1MpZnNFJU8.txt',
'byteFiles/gF1a0ZIBcAMKpJt0W5Sh.txt',
'byteFiles/6YexBrs0uzwSKTN5LpDn.txt',
'byteFiles/2Y8yNd7pfDX4l3qacwEk.txt',
'byteFiles/GKHQIcajdxPpbu8ltEvJ.txt',
'byteFiles/gYq9lthHL7vVfTSUBo2P.txt',
'byteFiles/DuAqvcFsUmZVMIEStakP.txt',
'byteFiles/CwtPLqOS2b3mg1rWJhx4.txt',
'byteFiles/1BYZ0XRoVfhtQGMJnkST.txt',
'byteFiles/58h7Q0perUzCyg3KHZPN.txt',
'byteFiles/72Aksjb4NMtae3whWUE6.txt',
'byteFiles/8g2wbZ4rox6KphNESDW7.txt',
'byteFiles/cw9MqhojugmlnC0iY386.txt',
'byteFiles/bnpVUwhFmj0AvQoL4ulx.txt',
'byteFiles/HmyrPzn4LNuxQetUM1fK.txt',
'byteFiles/e8uQB30WJNbqXEf9SaRI.txt',
'byteFiles/2rRK71xiQhSD4B6a9uEI.txt',
'byteFiles/76k80ovtGLQ5xBJKuUHb.txt',
'byteFiles/iRrWPIwXJje9usE2MvaD.txt',
'byteFiles/1Ch0ZFDqT9PzkGbSeyEV.txt',
'byteFiles/baDqk1ZjVdSAOPeBMF75.txt',
'byteFiles/6gHI9qzDjmLhd4nWS3oi.txt',
'byteFiles/g2pvrdbPVLKz408xBXef.txt',
'byteFiles/KRBDxPLfj0JQV5heGX4S.txt',
'byteFiles/IGlvaPi7VTjxkAmFb05C.txt',
'byteFiles/7ViTtBsPRgZQx1mqfHY0.txt',
'byteFiles/jJhBM9Xk6Hdu18c2yapZ.txt',
'byteFiles/dLgQ8Tpbc3o1anleXsKH.txt',
'byteFiles/7WkHb3T0oJ4XRYnBltNS.txt',
'byteFiles/K5tfseBz7vxwAhCPI6Rq.txt',
'byteFiles/E5CSXfWRKA8yhVNk0LPj.txt',
'byteFiles/abcOwUvuE9jGnpP5L241.txt',
'byteFiles/JPz2sCNRcjVUDrwLh8S3.txt',
'byteFiles/J2PF75IpyqAumjQVHZ6w.txt',
'byteFiles/6ZY9A10hSbQBGF1RcLVH.txt',
'byteFiles/gPICkhiGxzB9fea5wUWM.txt',
'byteFiles/bmrB20DR0ng6QwkCJXUH.txt',
'byteFiles/j6gIME8eK4T1mn9PFqhb.txt',
'byteFiles/7UEqLpy8zmkvfA2e5Nob.txt',
'byteFiles/hcnS2PgZufMNxtQ5GAz9.txt',
'byteFiles/dr2jSJMQ0sBozfZPLN18.txt',
'byteFiles/9Ix1He7TDc8iUFbGMum0.txt',
'byteFiles/DFJWGx1HYve5iBUp0K46.txt',
'byteFiles/9dcf2P1YWkyq7aZpge6E.txt',
'byteFiles/0sh9t0C5GprAyciTmaXW.txt',
'byteFiles/9YmT5IGJb0lxBD4h06Vy.txt',
'byteFiles/da4QzC51pXI0i3gy6Fm2.txt',
'byteFiles/eNRgx8kCGW9r64SbLhQt.txt',
'byteFiles/ARXaLhsFK8G1bSN2pjwz.txt',
'byteFiles/Cn09A1jfBmokT5qN7J6S.txt',
'byteFiles/G0mnSkldHD5ABFoNX7Jz.txt',
'byteFiles/BIoAC51ZjKvwIkWcigzs.txt',
'byteFiles/cWQjCi7fIb8BU0tY91Xp.txt',
'byteFiles/eBDOvM4fUQmy0XEhracH.txt',
'byteFiles/iYM4CwymksIGXDHT7uAT.txt',
'byteFiles/6LK1nVc8N7aHAyjbFWYT.txt',
'byteFiles/6U1ZRnLGvyg9iNFYwt7I.txt',

'byteFiles/194Iy5xv8QRz7XMTnmAk.txt',
'byteFiles/BK20sv0kr1lnCpSDUaIf.txt',
'byteFiles/gTzsRE5r0qPBDLWxt0mx.txt',
'byteFiles/7FQ15t8dsEISLkvBmqU2.txt',
'byteFiles/8GRqMaEyALrQtJfWIvk6.txt',
'byteFiles/1Q7DXjLweH09EokGNWiI.txt',
'byteFiles/j0ZMUJF5waGhqLQbWsYz.txt',
'byteFiles/6Ty0HzvEdYNFnmiafxgL.txt',
'byteFiles/ILkqxFMRCnTa465lyJtG.txt',
'byteFiles/iLdhAeC5DWmbVIBlTkV0.txt',
'byteFiles/CmoiKcBj6wySzXrJIfoW.txt',
'byteFiles/bB4krWM7TCiOD06x2AN1.txt',
'byteFiles/1PGWZqdnurksDEYgbB9X.txt',
'byteFiles/0pLrSeQMbfCc2IZHua0.txt',
'byteFiles/eq9gfGHRu1yE8CLlpARc.txt',
'byteFiles/FdSyWHUrZqeo9aE1mI70.txt',
'byteFiles/4jpx7Ah1w9f31YMikoCv.txt',
'byteFiles/15ADIikydX7eNq90WfcB.txt',
'byteFiles/BqXn7hw20xA6d3Di9QzY.txt',
'byteFiles/6GXdqZuQkcfbNyJA9gpx.txt',
'byteFiles/617nfOLPiCpmMhrVNI01.txt',
'byteFiles/DtnYlsgyQ654Xiw3dkoR.txt',
'byteFiles/8h7B6Zm4wfJGAHv3RpFW.txt',
'byteFiles/Ef1CXjtYxKousd7heFiw.txt',
'byteFiles/FmUz8pwN1XgbS7DW5yre.txt',
'byteFiles/kBoNGqDtARYyC6uz5V42.txt',
'byteFiles/BY2yoijLjKAPsgCvDXQW.txt',
'byteFiles/I0tk2asbMQNo9G8Bw6Sr.txt',
'byteFiles/afUoMhGcjJKL7kFTtINA.txt',
'byteFiles/7NhfmldVgt912wn8MYDs.txt',
'byteFiles/IRE32sreyaptx70cJTGb.txt',
'byteFiles/JBXrCHbwSvZ9c5y06701.txt',
'byteFiles/hczfxqdtwX61OPmZUJVS.txt',
'byteFiles/BMmCz4hZdgYjcbQD5oai.txt',
'byteFiles/a3BszoA7hXZc0EuHjmwi.txt',
'byteFiles/495fFS8xhen07WqZX6uw.txt',
'byteFiles/FmA5rcLou4t3f97bDvIK.txt',
'byteFiles/I0aTjkCwgMY8P6h9vWSt.txt',
'byteFiles/i1QF5bIjKc6AYl7u2fxe.txt',
'byteFiles/BfZpbY7j0PIw30rygsvq.txt',
'byteFiles/fRGWEv9TLIUOnhj0VV35.txt',
'byteFiles/cudBgeQFL8z4prNGIsXH.txt',
'byteFiles/I2W3i5wHhq9UpyFVRgzE.txt',
'byteFiles/5C0Ke7rXRtb9smU3Qkho.txt',
'byteFiles/a2UF0Vkeh5P74psmiR6v.txt',
'byteFiles/eajbXFgi4VTRHPM26dsh.txt',
'byteFiles/Kkd8eR0wGQLfXIrs53E.txt',
'byteFiles/2v1WwXSQbf3RrN4eTG81.txt',
'byteFiles/GkMB1vTidr5oyY0VJCwc.txt',
'byteFiles/9kc11e37zbpLWTPhSoiH.txt',
'byteFiles/0BEsCP7NAUy8XmkenHWG.txt',
'byteFiles/LNAnljesrEDma5oUfVZI.txt',
'byteFiles/cL8foiVrhFI21PQ5ARUt.txt',
'byteFiles/4kaSNXsgRmT6Mc0xrA5o.txt',
'byteFiles/eRLxnaXMK0zUtEjvi9H2.txt',
'byteFiles/B2Vfj0zqxgmZb7E8QpGK.txt',
'byteFiles/5MrYiQtEc4lwYBRWh8GJ.txt',
'byteFiles/1t03IYLhvDx6CKoHamQb.txt',
'byteFiles/3Ce9LP6ngx2uBtSyczRG.txt',
'byteFiles/198PdXRpGVZAyXcRF0e.txt',
'byteFiles/7JjQBd2CAyvtgY1SL0e3.txt',

'byteFiles/2ZI9SdMvc7WAXK03nPyu.txt',
'byteFiles/FQzrcxDH57n4XYwUltvb.txt',
'byteFiles/8WTumgpMQzadctY6vqZ7.txt',
'byteFiles/HBq8GsJYNS4Zenkxt1vp.txt',
'byteFiles/3QvRB7xD6n9kNutKOYAe.txt',
'byteFiles/JzTXaNDCqKGuiEvrHOYS.txt',
'byteFiles/i1hsFQal0fbxepMygR3C.txt',
'byteFiles/EXb2I0iutSyZ1V8r69wk.txt',
'byteFiles/3t50wjAi2VgW6vrnIQGe.txt',
'byteFiles/idQxka9c57Il1oBzjsgC.txt',
'byteFiles/hLfB8arbs4nmcx6tPAJz.txt',
'byteFiles/JI3ESrZoxtwm8kfdQ2UC.txt',
'byteFiles/cxCf1UQEJT7mNFg08KpM.txt',
'byteFiles/csEMJY7Rt0IGCKF2j5m.txt',
'byteFiles/dB4bEcHq15Xsnr7NQofV.txt',
'byteFiles/gJQKCF6t81Ge005PykWH.txt',
'byteFiles/2v819FWbJZg4G7MCNB6r.txt',
'byteFiles/6PG57gUmE2rLebCBTMno.txt',
'byteFiles/16cTMTKIjH5SbyovWuBq.txt',
'byteFiles/KBfcF596eUdtViub3ZrR.txt',
'byteFiles/bcqIoNTEXGPDUwzCWRx0.txt',
'byteFiles/0glscKoNakWL84EpunPe.txt',
'byteFiles/Haw0t4hzSPJyk6egrKfv.txt',
'byteFiles/DRw7V8B0HGF1mpjkeadN.txt',
'byteFiles/eNvtBzPm6MnGE8rJQcLF.txt',
'byteFiles/9ugQ8rA0fRaSOdVHbsYM.txt',
'byteFiles/3RrCi87ovZ1Qw6WuygGB.txt',
'byteFiles/egHnxmYdqACWaPsOBzRc.txt',
'byteFiles/AhEcNQy9nYHdfzZreoX1.txt',
'byteFiles/HaxgVnrq14bk7FEot9Sm.txt',
'byteFiles/jmH0bgfyxVa9pzCQT1WD.txt',
'byteFiles/ieCL1F3nH0Ky8rxzSGtQ.txt',
'byteFiles/FYPXNW2S6e09zMVR3HaK.txt',
'byteFiles/bsQcpVBgjyPva7I0kzwx.txt',
'byteFiles/DwNQdfBUbjI2WRGn1vLE.txt',
'byteFiles/JoDLpkCx5ziwtAU73NVF.txt',
'byteFiles/K5WyCRrPOn9bkGfxgY8v.txt',
'byteFiles/IVRTstGWhAPOdgY8UZiL.txt',
'byteFiles/dXB9U7M1pgt8lPGhyNuo.txt',
'byteFiles/Am7nEVXHePkCtqsdfzG5.txt',
'byteFiles/EXJvSznUcQjwyD4rGCNq.txt',
'byteFiles/dWnFKsSrTwUPpziR896B.txt',
'byteFiles/8z6VcnraQIhU0Yk1tg4e.txt',
'byteFiles/fPV9HUmlW2kwINly8YhOS.txt',
'byteFiles/5Um2hXGgeCVPwdtInObq.txt',
'byteFiles/foaW2s9e160A5kjJVcKx.txt',
'byteFiles/fQs8jHrye30NE6KMP9R0.txt',
'byteFiles/ft4RC3LwYxDib5K1EOX6.txt',
'byteFiles/C60EvgQ2n5F93Tjhoemt.txt',
'byteFiles/JuVXk5yQvxpOqPh3d9lI.txt',
'byteFiles/i1kqNORFDrEeHIdKuSy6.txt',
'byteFiles/I0wz6TmBL92MSGgxUboV.txt',
'byteFiles/faOmpIAWF9ZURkYqg8vr.txt',
'byteFiles/46ywoHZAMYPuqcbVQpUT.txt',
'byteFiles/ALuyE2Q9mjGdOkq8HWoz.txt',
'byteFiles/CVZd31hxBF0Tqu9K8ymE.txt',
'byteFiles/0XrH26wcU1ASsvK5ZqhD.txt',
'byteFiles/AuWm52oTYEKXDlgr8S00.txt',
'byteFiles/cnms0zSvd3FqDexiwZo6.txt',
'byteFiles/GLQ3Z104wCebKFV8A5z0.txt',
'byteFiles/2FpVtbaHJWAG0ImwTzxU.txt',

'byteFiles/0Cq4wfhLrKBJiut1lYAZ.txt',
'byteFiles/0sdr2cfHLzN09Ep4RSjT.txt',
'byteFiles/FRvg9Ji5jnkDWqm4sMcA.txt',
'byteFiles/Fhdbq2nojZ3pKW1HgrIQ.txt',
'byteFiles/imjGPLrpKkesMR764noI.txt',
'byteFiles/kc01Wwz52sCaEV1FXgpG.txt',
'byteFiles/IR2aS8pG4mbN59ZjCixk.txt',
'byteFiles/AogVYidsWqnUS4I6Q1cN.txt',
'byteFiles/gqm9zEewrY2hDRBTixs6.txt',
'byteFiles/Dmd0CialNW4R0kqXuvwE.txt',
'byteFiles/fr197DA51ehQgZULyW2z.txt',
'byteFiles/AHIPR4avZ1J9GiXsD3Wf.txt',
'byteFiles/iPqe0av3fHD2VJrt9GoE.txt',
'byteFiles/KivjcOQFy2PmDhodWJxC.txt',
'byteFiles/6g0bYqkiQC8R5ptdr9UN.txt',
'byteFiles/AK8ZjCh1nXIzygf2s3Tv.txt',
'byteFiles/65cjJpPCUQiLDRyXfWd4.txt',
'byteFiles/kh6DYLinI27EHPNSJXp.txt',
'byteFiles/CX7o16EG05dcu8pPirjI.txt',
'byteFiles/F30eacGAomRBgzErX4nL.txt',
'byteFiles/ECp9uB10cNkKGslaZnjU.txt',
'byteFiles/jNUWxT9ZE3AzQHJcf1Rg.txt',
'byteFiles/cwPNTvn3zIrJhXGK1aDS.txt',
'byteFiles/IkFwMHPbmDB0VvxEWOGX.txt',
'byteFiles/1QwaNAKMEpI9kb0oYzuJ.txt',
'byteFiles/1mCvIc0K8Gd3P5suXTRM.txt',
'byteFiles/36mPaSFkA7gVLnCbhqTt.txt',
'byteFiles/93qwE6DU817h5JxGBn4d.txt',
'byteFiles/EVYzBCh5WsU7RbHMgvFK.txt',
'byteFiles/jcxNbYug1rMTmGlQz904.txt',
'byteFiles/93tGL7eRYkUqXnDw8P0J.txt',
'byteFiles/Ftqd7rZi148x0zIJMQPY.txt',
'byteFiles/2NakDPd4GEWzuHhTx5o.txt',
'byteFiles/1gYyAWCIEbGBJaMHL0k2.txt',
'byteFiles/kJNLP6CyAKfgVMFwUInG.txt',
'byteFiles/0ZTEyLXaWReMK3rYVCjv.txt',
'byteFiles/Hai3lhUAXN6Qv0jW9FMg.txt',
'byteFiles/kMd0a9cZ3jw8GUfu52Tq.txt',
'byteFiles/IQk0AncFSwEV0Zh2efJy.txt',
'byteFiles/Czdn32NEL4YfxARaJuM9.txt',
'byteFiles/CmkrzR9pL7si010q35IM.txt',
'byteFiles/dt3fR14c06CpjsLHNWgn.txt',
'byteFiles/6K7qtMej5oVChsbfiuxH.txt',
'byteFiles/j1f907g3wFBd5GqlKNML.txt',
'byteFiles/JxCiRH4M6F5Kw8rDfjq.txt',
'byteFiles/kg24YRJTB8DNdKMxpOH.txt',
'byteFiles/bgP2EJeKxofuGBTTrqaNs.txt',
'byteFiles/iPtRTr5EZ4bQljzKXYcJ.txt',
'byteFiles/jg03mkayeMX8CuLbxFpV.txt',
'byteFiles/e5fn4GutJO601CgUSZF8.txt',
'byteFiles/d3zMqD8Fehx6EvRUmBuZ.txt',
'byteFiles/EmVYF0ZrIBWCGL5168XH.txt',
'byteFiles/2Ep7f84CvJDn9NwMHBPk.txt',
'byteFiles/JonIk6ScNpZM7WgFzm1D.txt',
'byteFiles/EPGmbKVeLBNAdnkpT5sM.txt',
'byteFiles/iknpnx8RDFJdFmIcQGzha.txt',
'byteFiles/dRYGQSA8TrmIPtabWNeu.txt',
'byteFiles/4GLY2XFSEMmB6CVsAJ9f.txt',
'byteFiles/BdgfJiVGChknbcSDUI4Q.txt',
'byteFiles/cN908Z6lWsCeQT5U2Kvh.txt',
'byteFiles/Aoa1lgwOyETuRGsIHCE3.txt',

'byteFiles/FY8XcAmQ9V51jwKTNPdq.txt',
'byteFiles/7rneDC6Ng10vEPuicRWV.txt',
'byteFiles/0xkbZDBOAicQJK5wC7Pm.txt',
'byteFiles/GAFXJNzhQ50rMLdt8p39.txt',
'byteFiles/iSD0zmncxkoV39KgPQLy.txt',
'byteFiles/FNk5ufg4GxHDLUpeOBmo.txt',
'byteFiles/6cJ40V1TGifn92SCPNEA.txt',
'byteFiles/c6aJtveQr9l5gsinVXRY.txt',
'byteFiles/h8a31y9YAc5CUst6JISG.txt',
'byteFiles/EhSAMWFg7Uk5oqBfN1cC.txt',
'byteFiles/dbMYm7QKarliB0hEDHU0.txt',
'byteFiles/fL1Dy8PFZHJwIU7vT312.txt',
'byteFiles/1hfkVAzp9jybt7YKEBrF.txt',
'byteFiles/I21HaMnK9rSlpNvoyUhg.txt',
'byteFiles/JNr1KDpRFXhP6wcQjMna.txt',
'byteFiles/3MbT7ePSqZuEWgXpiKDm.txt',
'byteFiles/KCxI1ZA3oiEqc8Xe4Mk0.txt',
'byteFiles/AIBO0E1uSDz3vhtwkofq.txt',
'byteFiles/hpntqKQa1RVmfJGu1B48.txt',
'byteFiles/4JEYXLGnodxWIcv6qZ9w.txt',
'byteFiles/6Y7KC0nLeOUjH9mpEtFh.txt',
'byteFiles/bFZTc8dX7Ui5JAVLvBx.txt',
'byteFiles/fR7HG9tZAaI1q62XnKyB.txt',
'byteFiles/2nBdD7s8Kklrtvm3zSow.txt',
'byteFiles/91HyVeXERTIaZcL4kWqM.txt',
'byteFiles/JEbogFC9kPnGSdi2KDaZ.txt',
'byteFiles/6gHQkuhNC1iso84zvqMG.txt',
'byteFiles/8jYo6TrpmwN2BGyEgsKd.txt',
'byteFiles/5Zrij7tnYzFRaxI9edvf.txt',
'byteFiles/fS6uGVWmT29r4N7qc5hp.txt',
'byteFiles/ezHTxiry7cw5dsY6NjsZ.txt',
'byteFiles/hpqvBuROHrncA5E0fxjy.txt',
'byteFiles/5H1DjF2dx8sJfcbeR6BW.txt',
'byteFiles/8s5WLuf0pivyOC7YI2er.txt',
'byteFiles/indtboM4Vuay6ZkYIQNw.txt',
'byteFiles/fZgza0Rsom2wvBLktWn3.txt',
'byteFiles/EPZ1kC7e34xH9RNVGAtW.txt',
'byteFiles/j5HNabSQY1ZivMrW2gFL.txt',
'byteFiles/asoPA4pgUtHm0dQzn9Tb.txt',
'byteFiles/ezUIPLcrOXBlw1mkTtQZ.txt',
'byteFiles/5VehuFCfk08kZx2DH1Wv.txt',
'byteFiles/CJyopa7dqfBG1R6wUgSV.txt',
'byteFiles/ehGfVZy2QXtFzJk01ESb.txt',
'byteFiles/dsB2UgC4XFQ1cj8DMyhT.txt',
'byteFiles/BxLpsG41I0qTgXKJPrtC.txt',
'byteFiles/gbYsj2FL7BQpaxX9310P.txt',
'byteFiles/8rKwm13nGVhWAiUjDH19.txt',
'byteFiles/6Jdm89ABg0zRfitPVQ3M.txt',
'byteFiles/DjV4Jdp805avuimhTBNX.txt',
'byteFiles/F3hnUSIGZ4Aqya8zvm1B.txt',
'byteFiles/IDWtdR130GLyVmPQs8uX.txt',
'byteFiles/c8ePf7t6VQAWNGOB40Zk.txt',
'byteFiles/dBuCJTrSQLaR8cmVjD1M.txt',
'byteFiles/2qJP9B5Q1xeDECIVzLo4.txt',
'byteFiles/HJ3G8CS9f4yVmadQ0t0P.txt',
'byteFiles/DsXzfyoVU1SAM8PwtZOB.txt',
'byteFiles/h7eLbvRP0mF8GWnSV0td.txt',
'byteFiles/Cq1FyksTYLVw0eQfJZGD.txt',
'byteFiles/hIDK0p1PAiJdL87CzG69.txt',
'byteFiles/0MmZ8j5pn2R3VG9wlxYi.txt',
'byteFiles/DNPAd4wCIHh59p2nbqJT.txt',

'byteFiles/58nzDrNVjhdog7QYkxe4.txt',
'byteFiles/5c9nbNjd6TxmDfP4Avp3.txt',
'byteFiles/4ZVUhT8gufvReG2XS51w.txt',
'byteFiles/jUY2ipvXO9E5SPuLwHGs.txt',
'byteFiles/gsqykQJmjZ9FdXMeYoNa.txt',
'byteFiles/Dc9mlVIqbHf03CX8gP7s.txt',
'byteFiles/bFwlXdoNrsPiWEkZB0va.txt',
'byteFiles/bDuyza0fBOtQ1PNFHnAJ.txt',
'byteFiles/fpeCgovE4zqPdDLmhxA.txt',
'byteFiles/2AiNV4msBIWd51whz9Xp.txt',
'byteFiles/3aisrY0cQetkUqnfp5K.txt',
'byteFiles/HkMSw1B8a963R0nPtihu.txt',
'byteFiles/7tDpNPXUIS1eJCkZQjr9.txt',
'byteFiles/6Cvyc7zMbzq5DfJprmXP.txt',
'byteFiles/AccWbv1hrPl0uNps5mgG.txt',
'byteFiles/Bmf4CXo9651TrigSb0jk.txt',
'byteFiles/8S9KmFuU2iXolxvsdnq0.txt',
'byteFiles/hXGizHjZ12qRYLw073eT.txt',
'byteFiles/BAyEjNmU0FSDwtZiWf4X.txt',
'byteFiles/eCFQTgkih9lsKtoqIAj1.txt',
'byteFiles/13YpdP5vTL0azSQFRgJn.txt',
'byteFiles/8LrpcC15nHjUDGu7Xio0.txt',
'byteFiles/DoVLn5KN3BZeMrzPCmhT.txt',
'byteFiles/7AOuExQwmYLvhqZ2T1cN.txt',
'byteFiles/g01pdR9EGwDQCK5f4LhJ.txt',
'byteFiles/6qjSVHW1IFQOTALiRzca.txt',
'byteFiles/57w4VeFy3QZx8PNk1rv2.txt',
'byteFiles/9IFO1DNV0XYi4mr3yzUs.txt',
'byteFiles/7P3A5vpdRa86qQ1LHJUt.txt',
'byteFiles/hHnpTaPYU4StmeC6yNxr.txt',
'byteFiles/kMALu9InmCyvPpEDL46i.txt',
'byteFiles/cwLo2p34lqzMi0rT5PAD.txt',
'byteFiles/5nKHxu3hfqJYvS2DQi7Z.txt',
'byteFiles/7c13Hk2SFD8NJeBYTUZs.txt',
'byteFiles/C23GNv6fP5JjLdrphgnc.txt',
'byteFiles/Hf2W0klmDrv5G6cEba7h.txt',
'byteFiles/8MSg6ycte4mZs0QL9bVN.txt',
'byteFiles/6fmGMN2PtKoAR5XkwZaj.txt',
'byteFiles/bMCSSqgxuwh1fOXUKLFm.txt',
'byteFiles/6aTHENeI5PBBygzJndM7V.txt',
'byteFiles/EXapH4diUyGeJ95vf1YT.txt',
'byteFiles/Gx3dKXvNTwWPFM8oZisj.txt',
'byteFiles/e7JJZ2kUbKHzhfpMRctqE.txt',
'byteFiles/a63qScpsxBFJtAXjd1nV.txt',
'byteFiles/baoYfR8qZ1cnAQ3dH6t5.txt',
'byteFiles/D0FvmACsaLoJVi0k3ybX.txt',
'byteFiles/7zp9WvQJTgrOnyjl4wGi.txt',
'byteFiles/HFyT41U2KtEZfXP3xYY6.txt',
'byteFiles/2s8QiutzMXLBcIlyHwkG.txt',
'byteFiles/ECYBut7y9fqKpw2g4xed.txt',
'byteFiles/bA1NQWCd1IBP3kuxYFM9.txt',
'byteFiles/JIPoNu820rbAGE4zcqaT.txt',
'byteFiles/ChfEv4KV6BQi0koFWMRd.txt',
'byteFiles/8H34IhqQL21V0at7CSdK.txt',
'byteFiles/BcuE9gIzJoh70dAlbVqw.txt',
'byteFiles/9Rwu8Es7InNiLjDkQSXd.txt',
'byteFiles/dbDTpWY0c8Jz6tuPfeoq.txt',
'byteFiles/gJQLAK983BsIvVwCHXRD.txt',
'byteFiles/hmxWlvCrP5zufkQb8Ljn.txt',
'byteFiles/hFVEPajwWYZU61Sb4Dyx.txt',
'byteFiles/FSZsJwdxRCkPAeB2Trf4.txt',

'byteFiles/jKxs6EIpRPOWuef8XykS.txt',
'byteFiles/dp4u8GQi5gJTHY12PsNq.txt',
'byteFiles/hrguW0CqMt2mJdfvX8IB.txt',
'byteFiles/CE2cZKbPRmYuOGwn5jDa.txt',
'byteFiles/H3teMfnQ1LGR5Sjr4p7W.txt',
'byteFiles/EsBRxv5MTp296yaKkwjW.txt',
'byteFiles/BmaZRd8oub0pFTYgN2Dn.txt',
'byteFiles/iXvTztNsUj9dJw5rqAWo.txt',
'byteFiles/3DPkmSKTydJUFQ19O14s.txt',
'byteFiles/d7czUVqJApx6WEj4ownT.txt',
'byteFiles/5wkJpsMIWzDaQ4njKuy7.txt',
'byteFiles/7AhtpZDSxMjVEIs8l09W.txt',
'byteFiles/hSNv2TgwqX0Ip3HnVFEj.txt',
'byteFiles/dtcBSuoU8MZ0wilHFe7n.txt',
'byteFiles/7Qx5U3phNBIMglLwV9PS.txt',
'byteFiles/jrAH5I14hpeni2TbcPOR.txt',
'byteFiles/1pf8sevSo30YJnhcZX6M.txt',
'byteFiles/c06HYvVGxLNWKZyCTgAf.txt',
'byteFiles/0Iv6U2hbCP1xeBitW50o.txt',
'byteFiles/KcfgCTNF3aMdqJ0jyWSG.txt',
'byteFiles/1pfyGlxvBcaMTNj3VgE9.txt',
'byteFiles/0pqtVjvLk89PdYr5MwZ0.txt',
'byteFiles/aJtX1TdUzZGoEwB3i5Pm.txt',
'byteFiles/2ISirUkfPhltjXVQ86Ge.txt',
'byteFiles/Ao2B0aDcrdZ8xLJt4Xes.txt',
'byteFiles/jAz2bdEis1hUHxcgNL15.txt',
'byteFiles/39gowdFUK2kGxcBnbALH.txt',
'byteFiles/BL4IWnsa2AbXyoPHQ0kJ.txt',
'byteFiles/60iPhkAuq2Zo1MxfDBNT.txt',
'byteFiles/iMIVeqXKnkP3NAc4uC1d.txt',
'byteFiles/KNP2R0q6J8YEcmmyrtSjV.txt',
'byteFiles/65nIFXYRgmzdsZq0kPox.txt',
'byteFiles/9zbitD6NE10I4FjdgJvV.txt',
'byteFiles/a7ugQM5sWX6R4fibefJFj.txt',
'byteFiles/JbaZC7Fdx0rzD83LEghm.txt',
'byteFiles/BCw58jIVumFW43RExqJk.txt',
'byteFiles/deou4UEvA59Vi2DZySL8.txt',
'byteFiles/BIDzw09f5oL1Q0tJXH2N.txt',
'byteFiles/k31UV2ChNYupKI8vZxme.txt',
'byteFiles/bSKx5RYr13dPteiJFgmU.txt',
'byteFiles/E8hFMgxcp9DAb1qGwQtB.txt',
'byteFiles/K05YJpcfytCU4EdShQI1.txt',
'byteFiles/4mAycV9nGEikubgUXteL.txt',
'byteFiles/EQ9XCpb8NMh2ql6Ki4ky.txt',
'byteFiles/Gtg83nF6ri2sIjAPMVhC.txt',
'byteFiles/f0jhPJ2913c56yY4mwEi.txt',
'byteFiles/BOeTymAhdfoYM23aj5Xr.txt',
'byteFiles/6nLCYyu3cTMDSXWxbj4o.txt',
'byteFiles/c1nP81vwZ6xB3jLdirMW.txt',
'byteFiles/E4SvkGO5CuB8LehadFo6.txt',
'byteFiles/eiaL8AoqubhkJxGWRFDM.txt',
'byteFiles/hAtDMOrJazquYw7Xk0Zj.txt',
'byteFiles/3kpwYGsKyF8HlnjDm6aA.txt',
'byteFiles/gYZueKo9wDtQ4FRbB2lf.txt',
'byteFiles/gkMcsUACjoBPnFYObQfa.txt',
'byteFiles/5mVLtUwgRfWi3yeCFJsZ.txt',
'byteFiles/EKeIjqfa10JmzunTxciN.txt',
'byteFiles/g00h7udtqs4yaXNQAUPM.txt',
'byteFiles/Drg6AkshP8zYZ1JoOwxf.txt',
'byteFiles/g89DUAmuqr4fMC02XSjf.txt',
'byteFiles/5fEXZxI1vB08nbFdS2Mi.txt',

'byteFiles/89UZDSqlnhazB0CJrk3t.txt',
'byteFiles/K2QNGrnBZP7I5MiF9kUb.txt',
'byteFiles/aitshTGMD1EnVb8oAIzf.txt',
'byteFiles/Czif6LUN1YaQOHmsX0T7.txt',
'byteFiles/c5hXZVENsBA4RvDfjJPw.txt',
'byteFiles/i2Jkh6cyaOwnjlwg5VN8.txt',
'byteFiles/cRsXVg26e15oik7JNje9.txt',
'byteFiles/Bnyr5gF3Q0SLeh7CEdoD.txt',
'byteFiles/HCbqVZF15800xptuS7sm.txt',
'byteFiles/IndbLEkgMT0ZXa7DS4YH.txt',
'byteFiles/8nc3AGwk7qMWrEJhdSxP.txt',
'byteFiles/7PtaAnTZi0QCphfNy6ms.txt',
'byteFiles/GZndsp4oXvf5Wzb96CQu.txt',
'byteFiles/gdjL7Gu2qxoX30HrSkWQ.txt',
'byteFiles/FqbIWSJMyRK1Ni7tVAu0.txt',
'byteFiles/kczYF5vSxW3ZrilfJesd.txt',
'byteFiles/C0zcsP1IiGEXyLkx6bp2.txt',
'byteFiles/a0vhxHWSPFtLfpUm793Z.txt',
'byteFiles/HhyutfnaJMeAN5wEsZ0g.txt',
'byteFiles/JYEvhBDFS2aNocqUwkMb.txt',
'byteFiles/gHUeRsYtF14BaX7So8W0.txt',
'byteFiles/hECVdzWU83LJkc46iGMR.txt',
'byteFiles/7Z54Kq6E2GTvgbyXunm8.txt',
'byteFiles/CTKbnhZkRy7AGstxePSd.txt',
'byteFiles/ebI3LHRqBUwco8aQtW15.txt',
'byteFiles/fQRdkX2ahe0t0c5ysNU8.txt',
'byteFiles/Fbdqw108uxIcZ6rDMgmY.txt',
'byteFiles/CYEUPAcItFWBs2VQTyKb.txt',
'byteFiles/h1nZq0BX5LruizPw3EUm.txt',
'byteFiles/lC4sdkhcZS6zYaE53D1v.txt',
'byteFiles/gDYeNAplMujiwPtU0qnJW.txt',
'byteFiles/dVSQ8bXUG0Jol3ZmDqOp.txt',
'byteFiles/e0ujNRMU2Txqd4rfAXaF.txt',
'byteFiles/H6q1L1JXFrbITo8MQdmO.txt',
'byteFiles/btaERgnQGZMzrDWvkSoU.txt',
'byteFiles/4m1Gsx5JvMFgATY7L0z0.txt',
'byteFiles/5z6qBep3fIQuW7dZhYAP.txt',
'byteFiles/80rsVW03Jp4LYekTAmZU.txt',
'byteFiles/biWpRxwyFrmuLSzK9ZIv.txt',
'byteFiles/2qk8R3jeQvKzsMp0iGT1.txt',
'byteFiles/JSoqeuhRdwMs3G76lQp8.txt',
'byteFiles/EKH1PBcSCXqnJzNRi8uj.txt',
'byteFiles/BqhT4IuJponYPSKRxk2Q.txt',
'byteFiles/BMjARP1c6GI70NfitDJd.txt',
'byteFiles/atGqPui5x1RbN6FDWhsZ.txt',
'byteFiles/joFHWBAZySvl8Ea1LD4N.txt',
'byteFiles/3JI42LUuZET9tgvlYDqx.txt',
'byteFiles/jGqYz9DaFMOWU1SwJBge.txt',
'byteFiles/74SYbIefnAEJpGHZq6md.txt',
'byteFiles/hSNuWGB9KZteamA7tF4R.txt',
'byteFiles/3nZODLNAGgEF76dck14t.txt',
'byteFiles/DhWInBmVuLfEF18YRPHZ.txt',
'byteFiles/gCJk0uacmv1BAosn74b6.txt',
'byteFiles/IR1CwpWFqsnd8gJZtbca.txt',
'byteFiles/gKvtaoZ7Lk11V6yxq3Fc.txt',
'byteFiles/gRNQ2Kw6dH10XrhBa0Vb.txt',
'byteFiles/iCN4j126qfaQxTAe5Jzh.txt',
'byteFiles/iqSDCrEjpfoRdaXHOIG4.txt',
'byteFiles/d7aHABNQz41Sni2XyVUJ.txt',
'byteFiles/jxDGoTHUWQeLfakytd05.txt',
'byteFiles/aud9U7Xni4qS0cTrmj0R.txt',

'byteFiles/gRZ9nxcONPftvzDCyKTQ.txt',
'byteFiles/JsTBaYEw5hQAC8kzcUMx.txt',
'byteFiles/k9LJAopKrhzt8H3iIDuf.txt',
'byteFiles/cXuqeLQGf6dyh4ZIUgFV.txt',
'byteFiles/IZeQVP4mMEwxvrzfUd6R.txt',
'byteFiles/HZ1WSp9dEVFumCsKg8Rv.txt',
'byteFiles/2JfpcM9ACq4XatFLho8z.txt',
'byteFiles/BwN6PaW0ZG1QA3DgICEx.txt',
'byteFiles/dW4sDcJMQCV0EgFjxXU0.txt',
'byteFiles/bRoM4EpjukSGXwAlZ6zn.txt',
'byteFiles/0DM3hS6Gg2QVKb1fZydv.txt',
'byteFiles/h1sIMzUVPxaBEX1nbdpS.txt',
'byteFiles/ieUqYbvXQJhLB4NDm097.txt',
'byteFiles/IQJlrKuUAyj4seftWL2b.txt',
'byteFiles/3suymN9GSHAtvYZ71Kaw.txt',
'byteFiles/fojWIsY21tQG8K6BCkgV.txt',
'byteFiles/bfRpQsFeqZN5BVrhkD6S.txt',
'byteFiles/7mGrwqj2tdfEuNDn4gVR.txt',
'byteFiles/1NvKF7pMREf4iVyzZewr.txt',
'byteFiles/gZyCkLS1Hf0x35REWjhjm.txt',
'byteFiles/7t2LuUCxAsgdqwWGIXY6.txt',
'byteFiles/kQsiVxDbAXt23wRWa157.txt',
'byteFiles/dTzmQ0cYijuhqwRtsPNe.txt',
'byteFiles/jsSM7FgazoEt1V80ynmY.txt',
'byteFiles/czBpmyW2hKAiGLDbv0Uq.txt',
'byteFiles/4f91UESFtRDrnLdyZOCN.txt',
'byteFiles/k4U7bfepOEJzgK8S231X.txt',
'byteFiles/K0ij9XYNHeJR37b6naD2.txt',
'byteFiles/J0Wv6HNY1xwuqPrICEdR.txt',
'byteFiles/BN76utiov1bnXza58ZTJ.txt',
'byteFiles/8sNQwJpT43ohCx2MjHaD.txt',
'byteFiles/hZBge41CJMxA60blif7P.txt',
'byteFiles/4S5AJVNXiIrfZG7na0KQ.txt',
'byteFiles/BKDwzzfquvWotCTXOM15.txt',
'byteFiles/EgJw2tGHKzz3YeR85aID.txt',
'byteFiles/BJxWK1M3u740AEkSoysr.txt',
'byteFiles/BFmbPXnMt19oVyI1SYvT.txt',
'byteFiles/eoA1Nd32LZGSPITn6mTM.txt',
'byteFiles/d0j39TQiauSXbpsINeU1.txt',
'byteFiles/16KuXyrdDOCWGhc7P4kL.txt',
'byteFiles/cUKzvqVwXx914BprdFgP.txt',
'byteFiles/K9kNdLS83clsP0wDF0bZ.txt',
'byteFiles/4sCh00eqidfLtS83JMVE.txt',
'byteFiles/bkyRaZK7X1lNBeWmDJGV.txt',
'byteFiles/8026Dh4VpfjPekaCgAYQ.txt',
'byteFiles/Fe4yEv8HmfLTBQuigKC6.txt',
'byteFiles/57v6nKMC81hLV34je2RW.txt',
'byteFiles/6hUFQK1g0x72GzzjJyI5.txt',
'byteFiles/hpRHqfonKgzAbd83XVUI.txt',
'byteFiles/HgcNkyo38ZnXM5ePvLQT.txt',
'byteFiles/3n1LAd1CtkErg090IfFx.txt',
'byteFiles/7P5gA4ZdSNWuYCi8lqpa.txt',
'byteFiles/dQ5kFDwP4ByqU3tnLTZf.txt',
'byteFiles/B0iZ1J2qwUe9RvbOSzGC.txt',
'byteFiles/aleOGLX5tBWIAMk2qTcY.txt',
'byteFiles/cVjkXNizAryxYd1R5oBL.txt',
'byteFiles/C4ViP Ej5K1vuetjhLdX3.txt',
'byteFiles/9HAkGKpDrq6hecQZysNd.txt',
'byteFiles/I9zqNg8Bj7RECDreiUX4.txt',
'byteFiles/K7AwxRfWPnc3hzXLatQE.txt',
'byteFiles/AsdcT125YI6a8SxzPbhf.txt',

```
'byteFiles/a1tHxq97wypZ3A8QjhK1.txt',
'byteFiles/4pehWH6b9IgL3ZF7uJCG.txt',
'byteFiles/A4i9wPlygmMQK2ftbXaI.txt',
'byteFiles/BFuieglUYJ934Dh0kw8.txt',
'byteFiles/doEyIsH3nRDTjfZ5lVY0.txt',
'byteFiles/hTe1goq5BG2YU7kQiuFm.txt',
'byteFiles/5qEodvVi3umUTpZfskhP.txt',
'byteFiles/27o93DUCHcGPXdTx1Nwg.txt',
'byteFiles/eS81zGXw3KYrBUnhm65a.txt',
'byteFiles/9mHqfWkgcJb7utKZ5rCV.txt',
'byteFiles/gybTWSU3EjBGpCVI1JA4.txt',
'byteFiles/6HcXhAzQiBMUK5en92FY.txt',
'byteFiles/B3F57b1RXv8cuGZifynU.txt',
'byteFiles/iHZnpXNuLTUQtElFz3jr.txt',
'byteFiles/83uOwngTAXtNQPvD2CxY.txt',
'byteFiles/dsuIp14XKzEt6kU29yeV.txt',
'byteFiles/ImBUE37WLcf9PJ4kgZnl.txt',
'byteFiles/7schXv12FHSujB3aVK9e.txt',
'byteFiles/dvVcMR0BzWxa8hAyZNHg.txt',
'byteFiles/CrKzvf0VWbetkShcayF8.txt',
'byteFiles/4S0KT7xLGJ91CzYfMNpI.txt',
'byteFiles/k4tKCVH9XwxuyieA2cG3.txt',
'byteFiles/4skKNvbjShmzaWHIRQXi.txt',
'byteFiles/0LAXajqhQy7po16dw8Tx.txt',
'byteFiles/7KVtWPn30iJpvBkIej4o.txt',
'byteFiles/iED1fIxngyWRLuYNvl9s.txt',
...]
```

In [8]:

```
num_samples = len(flist)
print("Number of samples: %i" %(num_samples))
```

Number of samples: 10868

In [9]:

```
import os
ids=[]
for i in os.listdir('byteFiles'):
    ids.append(i[:-4])
```

In [10]:

```
ids
```

Out[10]:

```
[ '9MW5Nuf0ogCEcR1YJeKG',
  'GFb1ksovaPYLI5XeC8RU',
  'kjQFMnf7vADCqtX4ai2u',
  'EKmgShQsf6a9vY0znN1U',
  '1WCXg8qEtJFITMilaf6k',
  '40dBpWdhOUzz5yVAKJF3',
  'bLNRehaM1SUd19V68ugD',
  'FfYe8u5S2rh1B0sDXIaZ',
  '0u3DIJTqnpkg187VSbOL',
  'jHfu9AM11GkNWch0Yv53',
  '4flnV7Xk0PyZ5ojuIiQc',
  'HoOEWQVSZyThrFj1tf2b',
  'b5qZNQRzLp6MULfHisoF',
  'bPx08C4QDMk9nFzyvHTB',
  '1D6IwN4EdfHta1YJuA95',
  'cPo7FESsJQvVNqmaOiiY',
  'iJpgS9hMDuFrCTItyok1',
  'ByjYPtU1Q0T76FekSJH5',
  '6WauKLy80DTdQf5hjZx7',
  'G7DQnb0fcE62usXeAmIN',
  '6QWh10fPM4EYicpGBv1L',
  'Esg6ijhbT8W9tQmcaH41',
  'kMDrx5LS2P864CIwNiVs',
  'KawMHneL7VYj8oQ4m1cZ',
  'gaYrbUuHQn9w5q10dMcm',
  '6tIbeDRa0jfnZYcFqVdL',
  '3WfGisHe71kxVFN8cjXz',
  'bG1zA9Iwqf6gCVSoJ0ZN',
  '2XiGfo95mCALRbu6KzVY',
  'gpQPsBn4Hc6uo xmFbYU3',
  '1JiaRhjycfZb8BpmX4TM',
  'EI4Jgyofj8CHR0ts59Ur',
  'CDFibvHczS19akfOEBwo',
  'kKEvpQMViqzLfeTWDtF2',
  'FNMk3wv1iVuQLCe90TDg',
  'bF5KEHN7mXfvktq0Qxzd',
  'JV PnwpfeExcNDt2ZTAQm',
  'i0JebqWZLA71XSktgc2N',
  'fs9oYJKeZ4BzXS1IWy3E',
  'Aghn90rSzad2DvLYPtqv',
  'dI89k371NKgURwFqsOBy',
  'FNKr36t4qTZoy5veUDRm',
  '2r4503Hczp6t0gVfAJBI',
  '2zrDU8t1pNjoWaisTeK9',
  'E5FzXZIDnTmau6VWepvC',
  'd9o6Is371a5eh0cQXtvY',
  'JBwr0MNAVT4z13kvifYb',
  'g6C1WcYBrAHskpjGXeFE',
  '9gpKWS6osX5Fq8E4UCLr',
  'fXOBiTSSxvUeIMnVha3y',
  '31KjQiNWATfPaozHtY01',
  '0P6tEopzr2mIhkuOKiCD',
  '7twBm0QdZXenW3vcUFJ2',
  '2eSoCjQzDri3E0cuxUYb',
  'GmPLUX1Y2DQhan563EtF',
  'baz0p8SLuGgisveCP0QI',
  'HRhfKEv12kTVqm95NQ8o',
  'bfcBD4va0tUJG6K5imqx',
  'gXLSZjPsvM1Bd0nJGQhN',
```

'J1YNWgvo2LCkQRV8d74u',
'fdGV3HjRu4nBSMNibIkX',
'jZxB2zDWto8UTu6hPENM',
'a0tVNG6x3empWgK2Zcyd',
'Jtp04ZqNWST65MdH7xhD',
'euAfyh0HkFiaKndptmsJ',
'Fa87QkVGAHec3PfuwbMm',
'2e5P1iKzXugGRBFf0Enc',
'2xX0y0GfquKYWIpsPa3n',
'CG6sL3Aug4yomNnztEIw',
'e9M8Qjno0aqN3HhLZkY4',
'jNUAeECmnssxGJ1QrTubv',
'0bjN3Kgw50ATsreRmEdi',
'gVzzU8WGxyoR07FesaJ1',
'esQpnxYJW1PVrT8ovu43',
'Ga17B5qZe6VoM0kCxTAF',
'g46jFCHzY1oqtws0yx05',
'c6tLU5OoSs7hqI9kmQVM',
'9Ur08SqELYFnQAHGva7z',
'3G6jJKRVLz1poimkvfhA',
'FsHojWp48ESXDbxczK0Q',
'50iq3CaS7H1PDWetAMLR',
'H6pGwSYfNmZVd80aeWg3',
'0u4QcXSvNvYnMa10BCygh',
'0sMxy5SNLu1YCkKJzaB3',
'Iec3SYD4Gimhob1xQ7sV',
'3EdJ5pSORwmeINFjaglQ',
'Dsg960KfZAq5ynarhiwT',
'6teIPpCrYKw0zFusDSix',
'ieFMEkfnc1pmS8taKHo9',
'juPCJRHOoLngNAvBZ7x1',
'ICleA7H05Kaf4jMXZG2y',
'6tmqe5B8JDsQnyhGo7vP',
'bROPgBf8pYcSud4Ejem1',
'imns4fDL9QbS7y60hwac',
'CkOT69iQsaLAGq8mYEeV',
'1eJx3418pcAFvMu0wrjB',
'8UGK0mbY3De124d7wizL',
'IweQoKz4HxGc8mAETNpM',
'6QWNYDyohV3R11fMCZI5',
'G29t8jq3S5usQkfeL4zn',
'F0ScZ1qtiRyYEQTVK4XH',
'9itfPnyp1HzhTb8IMYq3',
'9CVbzaWsLvcHjPqZu2Q3',
'Huz3U7hA5kVnJ6q2QcgM',
'dS6yJEktR1WQsBaIFYCp',
'fEMdwkrczmt7j0RTL14J',
'ITZQmUtjYaHDxMonFrKC',
'FNAVfQgh4a6k31UY5ZcK',
'60of29sarJPj8GEbnkTS',
'DwyNaEk0S4769ocfPOzn',
'9pSnLG2tRJTQ6C1Hzqfg',
'G31czXvpnwUfRtdJ4TFs',
'6DmCqAs0tFKxGe5aPz1U',
'0qNtRXvP8F2431bMmpse',
'jWN2deuOSfDhU4cR6XVA',
'J1hb9akq5NpXnvIg4DG8',
'jCAmu6J5gIBHU3TxopMD',
'A4xfRz73qWePOCnJ2Sya',
'aC3qABWUp29F0XvZNkIY',
'0pT04SVnWDehgU1YBAvq',

'HRE9TjDJgZl6LcuV7yKn',
'4WaGm3PhidQVRopgXz7J',
'exSWamc4Hi0j2CG5Rq1u',
'3v7fe0Diw9c0u6dAzsaE',
'5dZSwcADY0fCq8hsuJlL',
'JajfhmMDxzBX80Kb9o4k',
'19emSiJG8vrZgT2zXEIY',
'8J6QdyfRgIeTnFsNPKDp',
'3rf7KJiPd8hIZ645wBVN',
'BxRn5YjcIH1qCEy67NP9',
'BUGAnop1NDYTh0txE8d6',
'2Z1LcMw4Vkim1pNYXj9v',
'a6yPugSAr0dWK1zQ0c7I',
'2sf3AQd9MjP1xcLFRYYo',
'8GilsX4Uk7xQdmz3wgyp',
'Bnwd1vNx85KI7yzphZHs',
'Fh9HCK50ZId4PbcQDTJr',
'Jgtyz9OK13M2CYkmZolp',
'fvbowPJXxAKQsVt7MiRL',
'ig4IwJckB1RaT3FUN6rV',
'EpN4owBqFXM9YRyf3lUi',
'5p2yY0MvAQGXSIuocmbn',
'Jo4mF0en982QIKTabOXN',
'KWg8T3NOYaQBF0tMqLm5',
'gC2GsRX3ewpQyvj00Hb4',
'9DVvn8GrbQRsJwXCWcSo',
'D9QUL1KPYijhJezwSOFT',
'4M9zLzjcgVbwiyINDEgp',
'fSWvxIsOaV0YjkB185ut',
'4HGxlXJ8Ncy2Ri5az3o7',
'3iHvsmdx7DauhSPKfec4',
'j9ViIqQ4RXTm7hSFwb1W',
'hMyN4mwJjsLI5103PdV8',
'BLUqHvpdaytR0E3rZMu6',
'dLaQic9uU0G2SeRXDAHw',
'0X6U3SfcmbuZNodnH9qL',
'j04f0HCJN12TkB31PVyv',
'bf6w4apWBzyGLzhCqK7s',
'GBNdXiuAb7Cr05vVoK9I',
'IyV1H08cx fzMLvjFbrOX',
'2mYSX6J9Dd51kpLuNMc8',
'D4wXBShNamu3JLz1p8M0',
'EH7DRYeINXPfJUAG0Lqx',
'igjezBrTh4SbGtIKZm1A',
'IemKCvJBb7ia02owVYHN',
'KpzMHfmGQJbXN92Sdsyx',
'8A1wMUgpDrZ6HTcmQoCx',
'Hov1QtF67wa9ELyuJmnS',
'G1BR2jy6Nx1HUEag5ISm',
'FyIcXN7Z6SQojdgntUD9',
'HhFctSgYNdwYDxKLnAJC',
'8a3GEC7cPKtx01I6qYvL',
'7dYo5p6fSV0rkqgLvzbT',
'hnEAVHz5I7tQcjGWvRi1',
'10H9detwuiZByNSRMmx',
'i3HsbL10XJ2qtQPY9Cto',
'DZYefdjFtAGU91KVkJHw',
'ewsfhA3Dunp49SFNMoVR',
'iCPn8QUoEZ14IHp3jSkM',
'FgHkyL3ot79E1TmIZcbS',
'hTOJKwgGvBxb3dS8e9cn',

'EdT3lQiGHbzmrX2L1NgD',
'8d7zkmpAEtQrBiTRN1g6',
'5mcXAfSQW1V43hb6Ns1d',
'CLZ15SyTFMUBH31dR427',
'0hWRb28Umdgj7xcX0wtC',
'6prXgmtMVs7BYvPL3ce5',
'ELbIP3i6rgW41R19phnZ',
'3o1eL8jrlyb7MY6NASQi',
'e1KAXRxrOTc5snpuWVvj',
'73cMCHeUwStZjy2VoQRm',
'g5NnsVCLJdb1lzPBEUeT',
'6WKJVBfk5INU2q7RLDhY',
'B0dnNijSbH3fmuCQIL18',
'b4I670cxsfHSW1V8rJiEY',
'hdYWDuVTnF3kPv9CI4Er',
'f20PUJcTG4DOFRyh7Z1b',
'502xXcWOnqet1TrUmPBi',
'C7cZvXSQWwI0ftjA8xLN',
'Ek1sLKWN2DoUFriJy7mC',
'9VZ673sDcNIJKk8vP0TE',
'CKIT85YvdjSMoQr09A2L',
'GSYtwpD0MfegOyzNUHXr',
'iypLG5f97Y2ZTKcBuUxW',
'9FxgivRPntUB4YcV5IJa',
'15zY6vNZbOKF07uSn2Am',
'JFiG1dVRDj6MrH3fAK7L',
'gfSamebdC3Z4NFt5lonM',
'a1F0ZAc7xULvhP5Qsr4Y',
'EYJryk0oKhPwRNbWuIUz',
'hGJZKoMdcW9s6SEqNnOg',
'I9wXdHpzEAqKRo2TUtDN',
'23fQEb7B6W01VDIA4Sw5',
'81IgT57cylMUNGbLkiqm',
'g1Uw2y9I4EX83msAMQcK',
'H75RL96JxmMEfBvbe2jk',
'gZRseTmhYfHPvEnVFDiq',
'ADMm3H8ELRwcpUgjId6K',
'jx0WbXpVtGuyQI5DSUCB',
'84jKcbktfsh169TiErUG',
'IjC7vkJo1dViZerbzwT0',
'fYpFPV0nCAKS80qRtiH9',
'C9hHuINUVJqk1zo5pTQX',
'hNboYg1QSecVZ1LntIGf',
'fh1V8wnrOJAZGzsRQB60',
'ANwkGK2YbhD3VjsvHrLQ',
'1nyRDJYtOrpTEI3CmGHf',
'ChOGatx1eA8JNEVqFKTI',
'HpSJQCGdyYItq0eUuo3A',
'hnCWyLtgETF1kJ84DqXS',
'BMyRb1VjoWUws42eSLgc',
'J058MAd10WiuFk9zhspH',
'jEhpv08m7F0yKwXYoSs4',
'c1QNt9kyP5Tw8SDVfBpE',
'2oH0w4StzEV0hiKq3ByD',
'iK503zvsIBnMQAwmgHux',
'H7tC0wf5rzSluXiGsx8W',
'6DvOLIs4i1Gk7P9dEeTo',
'35rbwhMGTxqAWjuFt20m',
'crOuVx8jI6DARFmhWJ1M',
'an7uVbkw5MAPJf8oX9GQ',
'BH561pgQn32eUyJDFImv',

'd4lM2k7aq1RwJuTpg06X',
'aZoel3jMRHwBEG4X768N',
'0yvBxsufKXzM6CwpcSYU',
'jsqeT23MYgH7ytdxDUki',
'byKBxZ7dNkLanQtpiJHz',
'7LbrmXchDZJ1ed5HnAfM',
'6a7hF1ZnHrUAQXouv4x0',
'0ASH2csN7k8jZyoRaqtN',
'47GqvWtCkjShUVpE1oI0',
'a3ychUzAELtHgxI7kSoC',
'CM53GutBya9do7piSRe0',
'kJ1SuTtgh6fP9vas7wCH',
'5QAGY1NrvgSIeRxDH294',
'cuNE2eL7XtkbJQw50zsW',
'HpMdj30s7EDPkifLWvYy',
'FDHjobQnCUSs6atXldJL',
'0qPGt4cRVk9NoiJgubf2',
'EZWPt7i0ocxdnqG3BvRD',
'fEnNJU3pLTbSOeFxaQ0z',
'jgpTUqt96Bv7i4xcIGzF',
'Bv7fxdPDSZKmAXGp8baQ',
'GD7JUv1IfRuHpneMKaYc',
'5F4mbeBZsfwYCQpGgVLa',
'ESqUten2YJiswjZodR40',
'1PQFYMSBLA09TmKk2Zhj',
'JiUyEfT0m6oasXKj7p1G',
'ecvZowmJ09xHPGkBjCXV',
'0aVNj3qFgEZI6Akf4Kuv',
'GqPZCULjirRngSyMwmh4',
'fjbqGwVH6ANh5CR1uiri',
'CdJ3sjrGFOBaYTx7biRQ',
'bye3ZNnLt4BfuqAJ08v9',
'6SPDXoT29nxidIlhW1VC',
'Fw3QCzKGsD9e86pgLiIm',
'3Vhmj45EaPbB60rXoIMw',
'e1Pxcv2KRUSyDWmAasid',
'DNtJKyn0WV7PeS3q8FvO',
'efKbsRXGgy15NPHW8tI4',
'EBghIvP2MTJU6GXkAR0m',
'AbDuHJy7LUcZWnOaFEXh',
'fQmbIDTrZBds8Uu45HR0',
'Dw4btEPMk3reBuNsZ15S',
'IwgDdxrbUkiXpfAMy6J1',
'I2ZH0sW5hT4aDKwo9gR3',
'0Fu9oETtMW4z1g1ZrUy6',
'4wIxhD1pCbJVkZLfHKtP',
'G2NCXj6ZW0yTFALv7Hqd',
'dV8HAVzFer2naZEwu1Pk',
'BGpm6FzEHcn5SujJikwy',
'FBHvx5A1IchTEYDPkU0V',
'7u180DTFRYWPryM0xfLd',
'JuAyLarnUdx5cRYT1v7M',
'4aMbnH1A9es8vKmOFR2Z',
'HeBQriW1XgyVJcI31Y7f',
'dZB602DGWj5pkLCtzgv0',
'fcsNBbd6HYTaQAh5Mmq3',
'IopqlNQASCdjT8e0WRHw',
'0Hrfce4X5YGESJPjl9uL',
'HJE15j7Rz00garybNnh3',
'8djFkIgPK56XtYGDsb7R',
'DwEQhKiobyXvg81rjpRe',

'KIdHiQmorRG90vEcnbp7',
'409FLy5cim8qgxjYobET',
'Jy2Nw4xD3bPi1QfRsW5m',
'5XPfYvNdQRIBJoWmj3Os',
'EzxVkJ2DCYKWZpng39AoS',
'7aEbXfVnzhG6dRMLT58q',
'DoEY0jCZuGA3x4JSgW5d',
'0NXFnJyOEhBAISKfiU67',
'0SaYqdFRA5G7rfktjT63',
'awJH283VpYm7yBDxbOAI',
'jM0aVyQfdqBhODSox1Ce',
'1eCQkhvgDSFNXYLJbAdj',
'hdJnQNkT6b4yDeVwaR3F',
'5PoseIX8Ak7x24hKQuNE',
'bMK0Xukv9xUAi2QSGJ6Z',
'JFGA1Ypxt56ElPa3WUOQ',
'aE6dSw0llys08Qn4cf2FM',
'1xGM8wcYTV4pQFhJBd2f',
'7pbVgdSnjUPTrwcR5GXu',
'KAЕeГNIw967T5pHFhMQd',
'3L9crexw581VMXajnsmB',
'Fkha6WrBdNbMje70q0T4',
'bHT5pg8MuXLf2UZAwxBc',
'BtLTcn2E0adFGRPsKqYz',
'19zYbuW3XONcEedv7xU1',
'4syoiE9gr1HvtAfwtQY8',
'7PNjmAMet9183cVEDa1u',
'K1cOh76tkCvbNxPfLmD',
'cBTaW6V7hrS5fKgJvwja',
'e9Kji7QuZFdn1Lmc0Y03',
'7dk9aUmW0H1vsbtNSFuZ',
'bL63tFHUFNwMiAuvcV51',
'KRNHAm094TC70JfEPp8h',
'aXbDi405xeqQyn169Yk2',
'AKzFf5dE2PxtL0YqXv7G',
'A4oRW0sZJ3jkWazH5TLy',
'DcgUNDmzXlah1kbxiYtv',
'C6pcTZjgKOMtB0iREaSe',
'JCzVMIhYuK13Lzgo6Nct',
'ABNUEpLD8G1K6iubzwWe',
'1mStzMaQ8RUVBx7TkF4o',
'3UV6Iek05Y2pKRB7D1g9',
'FdPwzJ95fRp1LGceAEWI',
'JrsqMwaubVkJ3vj9him0',
'CZnRQNX0g8IGEcdMq1lt',
'b71ZgrI94Bp3XM0f6vuJ',
'hjESSm1TOUgAXkLzzirr',
'5LFWPi9qzEb6fhyVSYaw',
'BXQNYRIMHrq8Lz7Ay9pc',
'4ETQ9c3BOJ7YKNpFt0Ln',
'EjnR6L1w4p7sQkJWYiBq',
'gtiofz9cBwuYCHjyJTdK',
'dQbqMtOsVYrWDIVh8k2S',
'HNApayaXk5rm6g1tDcnoj',
'2MgFXNYimTJpIzKVE5wn',
'dAabnhJ37zorIS068q41',
'25jQ4xItcnbB6Hfo0E9T',
'HD1W2pcn973X5YaQVBtv',
'bv60XkYZsCKHa3NgFUDr',
'iks3Y5jQmILCpn806wR7',
'iTEz2rB6XdsZIG8g071K',

'dNw1IRg58zQ10oFVJS2q',
'DHGMN38qz61rTmpSXsbn',
'59npdcNYTD3qAvIPyQMs',
'eQVmTuqlHYXLJ6DBCsg5',
'9d3mFhrqe5gtuLE0ACOf',
'2qZrQoxyYDeadjIp3mSH',
'4gfR0YmZw0W8NvLSiXn1',
'hgbBPkqK8SOI7v0FDYtw',
'5YEafCpxyFcNqDG118TS',
'deTXH9Zau7qmM0yfYsRS',
'7waVJqpLnoxz4fKsZv39',
'DQqSM2enGHopTfCbg90c',
'h5uGZKiToFpXSrkAR1B',
'bNT4jk68t51zghy3HPcr',
'IVMFepgX5rASctxiL7y8',
'7xwMnp1sj3EhBJ0Zkd9Q',
'4Ub1A5oWqLNmdVF8sDw0',
'eKgQ10AwEzo12DBcs46I',
'4RKTw3jA6pyNbV21gUst',
'GU32g0fBF91xp7JqNXHY',
'FIVR08js5sZo6avbmPHk',
'5ov2DWngteLM0GVEySr1',
'4Qcv9hWIsmBX3PUNLltT',
'6LJuRVBNSdWGlaC7KbqQ',
'hmxjukz4UBqMPf9Jec6y',
'1Flm7qajn0zdtwBZ29gL',
'CfSF2urLkdAX6MbZQiIa',
'gCBJQKMq14Atfe3ZSRX9',
'FIxgyKh1a8MGTcAf7Xur',
'iI3TgAon2970drqY6CB0',
'ADaWiK4BCZg7Gub09xjm',
'40DnoUH5ets62Na11yVx',
'6qLobnwUp07W1RyKBjvZ',
'jw6dqJau5vzf7Vylsb9D',
'JchT00vFpwRX9beS3i2k',
'cqHlrY9oAVpyWMKJ8mOF',
'2Z0Jfu1Pd8Mq7xXNSDct',
'ETaX1605B0Hm7Ck3tGVF',
'40as7FtBN0GZTu05XpWq',
'2WHevmqLEZAb5ug3I7DN',
'49RfyxpQP1CzvneWIG5q',
'JvsW2yaUokpXVcbh1OMu',
'AYVnhTKJPbW2eBjrQ3Cc',
'6WKbReFUk18wYnCNQsyH',
'IYFTBN6LCn7SEQgbPOaU',
'HE9280Wkp4RJtdsNL7bZ',
'ASEw3L0HqWMePiy1uKxc',
'7QPJ36MV4rktjX0H9GLq',
'8BYcka0Zu47Rjm3bgUzp',
'DVsfP6pU89nBdhZHqxI5',
'ATr6NemOyI2KF0w48BR5',
'j3A60bUBPWgX7Qf2wo8s',
'izS2yBV3Z4nKch50mAuo',
'JFiWOaubTAKCyRoBqzfn',
'FZsflRA4XawTS7Qj0t2N',
'C5RGotXVr0xI3Np8g10A',
'HsL2jYueFUSCJZbPKzkn',
'gdMSLJ5UZmtPpN8W3ja6',
'jnEoBxGXRksZqHw8euhd',
'AUMy0gtvXYK74RT163ie',
'd0iY69mAf1JzMSr3vkcr',

'0FKerJ118x0c3jdoyg4A',
'FjeNJD8d04AakcKu5qCp',
'GEprqj7xuTa8m629zQP3',
'iSaFcugoKwRrmPTDZOyJ',
'CQS2WYL40q1MpZnNfJU8',
'gF1a0ZIBcAMKpJtOW5Sh',
'6YexBrs0uzwSKTN5LpDn',
'2Y8yNd7pfDX4l3qacwEk',
'GKHQIcajdxPpbu8ltEvJ',
'gYq91thHL7vVfTSUBo2P',
'DuAqvcFsUmZVMIEStakP',
'CwtPLqOS2b3mg1rWJhx4',
'1BYZ0XRoVfhtQGMJnkST',
'58h7Q0perUzCyg3KHZPN',
'72Aksjb4NMtae3whWUE6',
'8g2wbZ4rox6KphNESDW7',
'cw9MqhojugmlnC0iY386',
'bnpVUwhFmj0AvQoL4ulx',
'HmyrPzn4LNuxQetUM1fK',
'e8uQB30WJNbqXEf9SaRI',
'2rRK71xiQhSD4B6a9uEI',
'76k80ovtGLQ5xBJKuUhb',
'iRrWPiWxJje9usE2MvaD',
'1Ch0ZFDqT9PzkGbSeyEV',
'baDqk1ZjVdSAOPeBMF75',
'6gHI9qzDjmLhd4nWS3oi',
'g2pvrdbPVLKz408xBXef',
'KRBDxPLfj0JQV5heGX4S',
'IG1vaPi7VTjxkAmFb05C',
'7ViTtBsPRgZQx1mqfHYO',
'jJhBM9Xk6Hdu18c2yapZ',
'dLgQ8Tpbc3o1anleXsKH',
'7WkHb3TOoJ4XRYnBltNS',
'K5tfseBz7vxwAhCPI6Rq',
'E5CSXfWRKA8yhVNk0LPj',
'abcOwUvuE9jGnpP5L241',
'JPz2sCNRcjVUDrwLh8S3',
'J2PF75IpyqAumjQVHZ6w',
'6ZY9A10hSbQBGF1RcLVH',
'gPICkhiGxzB9fea5wUWM',
'bmrB2ODR0ng6QwkCJXUH',
'j6gIME8eK4T1mn9PFqhb',
'7UEqLpy8zmkvfA2e5Nob',
'hcnS2PgZufMNxtQ5GAz9',
'dr2jSJMQ0sBozfZPLN18',
'9Ix1He7TDc8iUFbGMumO',
'DFJWGx1HYve5iBUp0K46',
'9dcf2P1Ywkyq7aZpge6E',
'0sh9tOC5GprAyciTmaXW',
'9YmT5IGJb01xBD4h06Vy',
'da4Qzc51pXI0i3gy6Fm2',
'eNRgx8kCGW9r64SbLhQt',
'ARXaLhsFK8G1bSN2pjwz',
'Cn09A1jfBmokT5qn7J6S',
'G0mnSkldHD5ABFoNX7Jz',
'B1oAC51ZjKvw1kWcigzs',
'cWQjCi7fIb8BU0tY91Xp',
'eBD0Vm4fUQmy0XEhrach',
'iYM4CwymksIGXDht7uAT',
'6LK1nVc8N7aHAyjbFWYT',
'6U1ZRnLGvyg9iNFYwt7I',

'194Iy5xv8QRz7XMTnmAk',
'BK20sv0kr1lnCpSDUaIf',
'gTzsRE5r0qPBDLWXtOmx',
'7FQ15t8dsEISLkvBmqU2',
'8GRqMaEyALrQtJfWIVk6',
'1Q7DXjLweH09EokGNWii',
'j0ZMUJF5waGhqLQbWsYz',
'6Ty0HzvEdYNFnmiafxgL',
'ILkqxFMRCnTa465lyJtG',
'iLdhAeC5DWmbVIB1TkvO',
'CmoiKcBj6wySzXrJIff0W',
'bB4krWM7TCiOD06x2AN1',
'1PGWZqdnurksDEYgbB9X',
'0pLrSeQMbfC2IZHUaO',
'eq9gfGHRU1yE8CLlpARc',
'FdSyWHUrZqeo9aE1mI70',
'4jpx7Ah1w9f31YMikoCv',
'15ADIikydX7eNq90WfcB',
'BqXn7hw20xA6d3Di9QzY',
'6GXdqZuQkcfbNyJA9gpx',
'617nfOLPiCpmMhrVNI01',
'DtnYlsgyQ654Xiw3dkoR',
'8h7B6Zm4wfJGAHv3RpFW',
'Ef1CXjtYxKousd7heFiw',
'FmUz8pwN1XgbS7DW5yre',
'kBoNGqDtARYyC6uz5V42',
'BY2yoijLjkAPsgCvDXQW',
'I0tk2asbMQNo9G8Bw6Sr',
'afUoMhGcjJKL7kFTtINA',
'7NhfmldVgt912wn8MYDs',
'IRE32sreyaptx70cJTGb',
'JBXrCHbwSvZ9c5y06701',
'hczfxqdtwX61OPmZUJVS',
'BMmCz4hZdgYjcbQD5oai',
'a3BszoA7hXZc0EuHjmwi',
'495ffS8xhen07WqZX6uw',
'FmA5rcLou4t3f97bDvIK',
'I0aTJkCwgMY8P6h9vWSt',
'i1QF5bIjKc6AY17u2fxe',
'BfZpbY7j0PIw30rygsvq',
'fRGWEv9TLIUOnhj0YV35',
'cudBgeQFL8z4prNGIsXH',
'I2W3i5wHhq9UpyFVRgzE',
'5C0Ke7rXRtb9smU3Qkho',
'a2UF0Vkeh5P74psmiR6v',
'eajbXFgi4VTRHPM26dsh',
'Kkd8eR0wGQSLfXIrs53E',
'2v1WwXSQbf3RrN4eTG81',
'GkMB1vTidr5oyYOVJCwc',
'9kc11e37zbpLWTPhSoiH',
'0BEsCP7NAUy8XmkenHWG',
'LNAn1jesrEDma5oUfVZI',
'cL8foiVrhFI21PQ5ARUt',
'4kaSNXsgRmT6Mc0xrA5o',
'eRLxnaXMKOzUtEjvi9H2',
'B2Vfj0zqxgmZb7E8QpGK',
'5MrYiQtEc4lwypBRWh8GJ',
'1t03IYLhvDx6CKoHamQb',
'3Ce9LP6ngx2uBtSyczRG',
'198PdXRpGVZAyXCrF0e',
'7JjQBd2CAyvtgY1SL0e3',

'2ZI9SdMvc7WAXK03nPyu',
'FQzrcxDH57n4XYwU1tvb',
'8WTumgpMQzadctY6vqZ7',
'HBq8GsJYNS4Zenkxt1vp',
'3QvRB7xD6n9kNutKOYAe',
'JzTXaNDcQkGUiEvrHOYS',
'i1hsFQa10fbxepMygR3C',
'EXb2I0iutSyZ1V8r69wk',
'3t50wjAi2VgW6vrnIQGe',
'idQxka9c57I11oBzjsgC',
'hLfB8arbs4nmcx6tPAJz',
'JI3ESrZoxtwm8kfdQ2UC',
'cxCf1UQEJT7mNFg08KpM',
'csEMJY7Rt0IGCTKF2j5m',
'dB4bEcHq15Xsnr7NQofV',
'gJQKCF6t81Ge005PykWH',
'2v819FWbJZg4G7MCNB6r',
'6PG57gUmE2rLebCBTMno',
'16cT MtKIjH5SbyovWuBq',
'KBfcF596eUdtViub3ZrR',
'bcqIonTEXGPDUwzCWRxO',
'0glscKoNakWL84EpunPe',
'Haw0t4hzSPJyk6egrKfv',
'DRw7V8B0HGF1mpjkeadN',
'eNvtBzPm6MnGE8rJQcLF',
'9ugQ8rA0fRaSOdVHbsYM',
'3RrCi87ovZ1Qw6WuygGB',
'egHnxmYdqACWaPsOBzRc',
'AhEcNQy9nYHdfzZreoX1',
'HaxgVnrql4bk7FEot9Sm',
'jmH0bgfyxVa9pzCQT1WD',
'ieCL1F3nH0Ky8rxzSGtQ',
'FYPXNW2S6e09zMVR3HaK',
'bsQcpVBgjyPva7I0kzwx',
'DwNQdfBUbj12WRGn1vLE',
'JoDLpkCx5ziwtAU73NVF',
'K5WyCRrPOn9bkGfxgY8v',
'IVRTstGWhAP0dgY8UZiL',
'dXB9U7M1pgt81PGhyNuo',
'Am7nEVXHePkCtqsdfzG5',
'EXJvSznUcQjwyD4rGCNq',
'dwnFKsSrTwUPpziR896B',
'8z6VcnraQIhU0Yk1tg4e',
'fPV9HUmlW2kwINly8YhOS',
'5Um2hXGgeCPWdtInObq',
'foaW2s9e160A5kjJVcKx',
'fQs8jHrye30NE6KMP9RO',
'fT4RC3LwYxDib5K1EOX6',
'C60EvgQ2n5F93Tjhoemt',
'JuVXk5yQvxpoqPh3d91I',
'i1kqNORFDreEhIdKuSy6',
'I0wz6TmBL92MSGgxUboV',
'fa0mpIAWF9ZURkYqg8vr',
'46ywoHZAMYPuqcbVQpUT',
'ALuyE2Q9mjGd0kq8Hw0z',
'CVZd31hxBF0Tqu9K8ymE',
'0XrH26wcU1ASsvK5ZqhD',
'AuWm52oTYEKXDigr8S00',
'cnms0zSvd3FqDexiwZo6',
'GLQ3Z104wCebKFV8A5z0',
'2FpVtbaHJWAG0ImwTzxU',

'0Cq4wfhLrKBJiut1lYAZ',
'0sdr2cfHLzN09Ep4RSjt',
'FRvg9Ji5jnkDWqm4sMcA',
'Fhdbq2nojZ3pKW1HgrIQ',
'imjGPLrpKkesMR764noI',
'kc01Wwz52sCaEV1FXgpG',
'IR2aS8pG4mbN59ZjCixk',
'AogVYidsWqnUS4I6Q1cN',
'gqm9zEewrY2hDRBTixs6',
'Dmd0CialNW4R0kqXuvwE',
'fr197DA51ehQgZULyW2z',
'AHIPR4avZ1J9GiXsD3Wf',
'iPqe0av3fHD2VJrt9GoE',
'KivjcOQFy2PmDhodWJxc',
'6g0bYqkiQC8R5ptdr9UN',
'AK8ZjCh1nXIzygf2s3Tv',
'65cjJpPCUQiLDRyXfWd4',
'kh6DYLinI27EHWPNSJxp',
'CX7o16EG05dcu8pPirji',
'F30eacGAomRBgzErX4nL',
'ECp9uB1OcNkKGs1aZnjU',
'jNUWxT9ZE3AzQHJcf1Rg',
'cwPNTvn3zIrJhXGK1aDS',
'IkFwMHPbmDB0VvxEWOGX',
'1QwaNAKMEpI9kb0oYzuJ',
'1mCvIc0K8Gd3P5suXTRM',
'36mPaSFkA7gVLnCbhqTt',
'93qwE6DU817h5JxGBn4d',
'EVYzBCh5WsU7RbHMgvFK',
'jcxNbYug1rMTmG1Qz904',
'93tGL7eRYkUqXnDw8POJ',
'Ftqd7rZi148x0zIJMQPY',
'2NakDPd4GWEwzuHhTx5o',
'1gYyAWCIEbGBJaMHL0k2',
'kJNLp6CyAKfgVMFwUInG',
'0ZTEyLXaWReMK3rYVCjv',
'HaI3lhUAXN6Qv0jW9FMg',
'kMdOa9cZ3jw8GUfu52Tq',
'IQk0AncFSwEVOZh2efJy',
'Czdn32NEL4YfxARaJuM9',
'CmkrzR9pL7si010q35IM',
'dt3fR14c06CpjlsLHNWgn',
'6K7qtMej5oVChsbfiuxH',
'j1f907g3wFBd5GqlKNML',
'JxCiRH4M6F5Kw8rDfjq',
'kg24YRJTB8DNdKMxpwOH',
'bgP2EJeKxofuGBTRqaNs',
'iPtRTt5EZ4bQ1jzKXYcJ',
'jg03mkayeMX8CuLbxFpV',
'e5fn4GutJ0601CgUSZF8',
'd3zMqD8Fehx6EvRUmBuZ',
'EmVYF0ZrIBWCGL5168XH',
'2Ep7f84CvJDn9NwMHBPK',
'JonIk6ScNpZM7WgFzm1D',
'EPGmbKVeLBNAdnkpT5sM',
'iknpx8RDFJdFmIcQGzha',
'dRYGQSA8TrmIPtabWNeu',
'4GLY2XFSEMmB6CVsAJ9f',
'BdgfJiVGChknbcSDUI4Q',
'cN908Z6lWsCeQT5U2Kvh',
'Aoa1lgwOyETuRGsIHCe3',

'FY8XcAmQ9V51jwKTNPdq',
'7rneDC6Ng10vEPuicRWV',
'0xkbZDBOAicQJK5wC7Pm',
'GAFXJNzhQ50rMLdt8p39',
'iSD0zmncxkoV39KgPQLy',
'FNk5ufg4GxHDLUpeOBmo',
'6cJ40V1TGifn92SCPNEA',
'c6aJtveQr915gsinVXRY',
'h8a31y9YAc5CUst6JISG',
'EhSAMWFg7Uk5oqBfNlcC',
'dbMYm7QKarliBohEDHU0',
'fL1Dy8PFZHJwIU7vT312',
'1hfkVAzp9jybt7YKEBrF',
'I21HaMnK9rS1pNvoyUhg',
'JNr1KDpRFXhP6wcQjMna',
'3MbT7ePSqZuEWgXpiKDm',
'KCxI1ZA3oiEqc8Xe4MkO',
'AIB00E1uSDz3vhtwkofq',
'hpntqKQa1RVmfJGu1B48',
'4JEYXLGnodxWIcv6qZ9w',
'6Y7KC0nLeOUjh9mpEtFh',
'bFZTc8dX7Ui5JAVLvbX',
'fR7HG9tZAaI1q62XnKyB',
'2nBdD7s8Kklrtvm3zSow',
'91HyVeXERtIaZcL4kwqm',
'JEbogFC9kPnGSdi2KDaZ',
'6gHQkuhNC1iso84zvqMG',
'8jYo6TrpmwN2BGyEgsKd',
'5Zrij7tnYzFRaxI9edvf',
'fS6uGVWmT29r4N7qc5hp',
'ezHTxiry7cw5dsY6NjSZ',
'hpqvBuROHrncA5E0fxjy',
'5H1DjF2dx8sJfcbeR6BW',
'8s5WLuf0pivyOC7YI2er',
'indtboM4Vuay6ZkYIQNw',
'fZgza0Rsom2wvBLktWn3',
'EPZ1kC7e34xH9RNVGAtW',
'j5HNabSQY1ZivMrW2gFL',
'asoPA4pgUtHm0dQzn9Tb',
'ezUIPLcrOXBlw1mkTtQZ',
'5VehuFCfk08kZx2DH1Wv',
'CJyopa7dqfBG1R6wUgSV',
'ehGfVZy2QXtFzJk01ESb',
'dsB2UgC4XFQ1cj8DMyhT',
'BxLpsG41I0qTgXKJPrtC',
'gbYsj2FL7BQpaxX9310P',
'8rKwm13nGVhWAiUjDH19',
'6Jdm89ABg0zRfitPVQ3M',
'DjV4Jdp805avuimhTBNX',
'F3hnUSIGZ4Aqya8zvm1B',
'IDWtdR130GLyVmPQs8uX',
'c8ePf7t6VQAwnGOB40Zk',
'dBuCJTrSQLaR8cmVjD1M',
'2qJP9B5Q1xeDECIVzLo4',
'HJ3G8CS9f4yVmAdQ0t0P',
'DsXzfyoVU1SAm8PwtZOB',
'h7eLbvRPOmF8GWnSV0td',
'CqlFyksTYLVw0eQfJZGD',
'hIDK0p1PAiJdl87CzG69',
'0MmZ8j5pn2R3VG9wlxYi',
'DNPAd4wCIHh59p2nbqJT',

'58nzDrNVjhdog7QYkxe4',
'5c9nbNjd6TxmDfP4Avp3',
'4ZVUhT8gufvReG2XS51w',
'jUY2ipvXO9E5SPuLwHGs',
'gsqykQJmjZ9FdXMeYoNa',
'Dc9mlVIqbHf03CX8gP7s',
'bFw1XdoNrsPiWEkZBOva',
'bDuyza0fB0tQ1PNFHnAJ',
'fpeCgovE4zqPdDLmhxNA',
'2AiNV4msBIWd51whz9Xp',
'3aisrY0cQetkUqnfvp5K',
'HkMSw1B8a963ROnPtihu',
'7tDpNPXUIS1eJCKZQjr9',
'6Cvyc7zMbZq5DfJprmXP',
'AcCwbv1hrPl0uNps5mgG',
'Bmf4CXo9651TrigSb0jk',
'8S9KmFuU2iXo1xvsdnq0',
'hXGizHjZ12qRYLw073eT',
'BAyEjNmU0FSDwTZiWf4X',
'eCFQTgkih9lsKtoqIAj1',
'13YpdP5vTL0azSQFRgJn',
'8LrpcC15nHjUDGu7Xio0',
'DoVLn5KN3BZeMrzPCmhT',
'7AOuExQwmYLvhqZ2T1cN',
'g01pdR9EGwDQCK5f4LhJ',
'6qjSVHW1IFQOTALiRzca',
'57w4VeFy3QZx8PNk1rv2',
'9IF01DNV0XYi4mr3yzUs',
'7P3A5vpdRa86qQ1LHJUt',
'hHnpTaPYU4StmeC6yNxr',
'kMALu9InmCyvPpEDL46i',
'cwLo2p34lqzMiOrT5PAD',
'5nKHxu3hfqJYvS2DQi7Z',
'7c13Hk2SFD8NJeBYTUzs',
'C23GNv6fP5JjLdrphgnc',
'Hf2W0k1mDrv5G6cEba7h',
'8MSg6ycte4mZs0QL9bVN',
'6fmGMN2PtKoAR5XkwZaj',
'bMCSSqgxuwh1fOXUKLFm',
'6aTHENeI5PBBygzJndM7V',
'EXapH4diUyGeJ95vf1YT',
'Gx3dKXvNTwWPFM8oZisj',
'e7JZ2kUbKHzhfpMRctqe',
'a63qScpsxBFJtAXjd1nV',
'baoYfR8qZ1cnAQ3dH6t5',
'D0FvmACsaLoJVi0k3ybX',
'7zp9WvQJTgrOnyj14wGi',
'HFyT41U2KtEZfXP3xYV6',
'2s8QiutzMXLBcIlyHwkG',
'ECYBut7y9fqKpw2g4xed',
'bA1NQWCd1IBP3kuxYFM9',
'JIPoNu820rbAGE4zcqaT',
'ChfEv4KV6BQi0koFWMRd',
'8H34IhqQL21V0at7CSdK',
'BcuE9gIzJoh70dAlbVqw',
'9Rwu8Es7InNiLjDkQSXd',
'dbDTpWY0c8Jz6tuPfeoq',
'gJQLAK983BsIvVwCHXRD',
'hmxW1vCrP5zufkQb8Ljn',
'hFVEPajwWYZU61Sb4Dyx',
'FSZsJwdxRCKPAeB2Trf4',

'jKxs6EIpRPOWuef8XykS',
'dp4u8GQi5gJTHY12PsNq',
'hrguW0CqMt2mJdfvX8IB',
'CE2cZKbPRmYuOGwn5jDa',
'H3teMfnQ1LGR5Sjr4p7W',
'EsBRxv5MTp296yaKkwjW',
'BmaZRd8oub0pFTYgN2Dn',
'iXvTztNsUj9dJw5rqAWo',
'3DPkmSKTydJUFQ19014s',
'd7czUVqJApX6WEj4ownT',
'5wkJpsMIWzDaQ4njKuy7',
'7AhtpZDSxMjVEIs8l09W',
'hSNv2TgwqX0Ip3HnVFEj',
'dtcBSuoU8MZ0wilHF7n',
'7Qx5U3phNBIMglLwV9PS',
'jrAH5I14hpeni2TbcPOR',
'1pf8sevSo3OYJnhcZX6M',
'c06HYvVGxLNWKzyCTgAf',
'0Iv6U2hbcP1xeBitW50o',
'KcfgCTNF3aMdqJ0jyWSG',
'1pfyG1xvBcaMTNj3VgE9',
'0pqtVjvLk89PdYr5MwZO',
'aJtXlTdUzzGoEwB3i5Pm',
'2ISirUkfPhltjXVQ86Ge',
'Ao2B0aDcrdZ8xLJt4Xes',
'jAz2bdEis1hUhxcgNL15',
'39gowdFUK2kGxcBnbALH',
'BL4IWnsa2AbXyoPHQ0kJ',
'60iPhkAuq2Zo1MxfDBNT',
'iM1veqXKnkP3NAc4uC1d',
'KNP2R0q6J8YEcmmyrtSjV',
'65nIFXYRgmzdsZq0kPox',
'9zbitD6NE10I4FjdgJvV',
'a7ugQM5sWX6R4fibeJFj',
'JbaZC7Fdx0rzD83LEghm',
'BCw58jIVumFW43RExqJk',
'deou4UEvA59Vi2DZySL8',
'BIDzw09f5oL1Q0tJXH2N',
'k31UV2ChNYupKI8vZxme',
'bSKx5RYr13dPteiJFgmU',
'E8hFMgxc9DAb1qGwQtB',
'K05YJpcfyTCU4EdShQI1',
'4mAycV9nGEikubgUXteL',
'EQ9XCpb8NMh2ql6Ki4ky',
'Gtg83nF6ri2sIjAPMVhc',
'f0jhPJ2913c56yY4mwEi',
'BOeTymAhdf0YM23aj5Xr',
'6nLCYyu3cTMDSXWxbj4o',
'c1nP81vwZ6xB3jLdirMW',
'E4SvkGO5CuB8LehadFo6',
'eiaL8AoqubhkJxGWRFDM',
'hAtDMOrJazquYw7Xk0Zj',
'3kpwYGsKyF8HlnjDm6aA',
'gYZueKo9wDtQ4FRbB2lf',
'gkMcsUACjoBPnFY0bQfa',
'5mVLtUwgRfWi3yeCFJsZ',
'EKeIjqfa10JmzunTXciN',
'g00h7udtqs4yaXNQAUPM',
'Drg6AkshP8zYZ1Jo0wxr',
'g89DUAmuqr4fMC02XSjf',
'5fEXzxI1vB08nbFdS2Mi',

'89UZDSqlnhazB0CJrk3t',
'K2QNGrnBZP7I5MiF9kUb',
'aitshTGMD1EnVb8oAIzf',
'CziF6LUN1YaQOHmsX0T7',
'c5hXZVENsBA4RvDfjJPw',
'i2Jkh6cyaOwnjltg5VN8',
'cRsXVg26e15oik7JNjE9',
'Bnyr5gF3Q0SLeh7CEdoD',
'HCbgVZF15800xptuS7sm',
'IndbLEkgMT0ZXa7DS4YH',
'8nc3AGwk7qMWrEJhdSxP',
'7PtaAnTZi0QCphfNy6ms',
'GZndsp4oXvf5Wzb96CQu',
'gdjL7Gu2qxoX30HrSkWQ',
'FqbIWSJMyRK1Ni7tVAuO',
'kczyF5vSxW3ZrilfJesd',
'C0zcsP1IiGEXyLkx6bp2',
'a0vhxHwSPFtlfpUm793Z',
'HhyutfnajMeAN5wEsZ0g',
'JYEvhBDFS2aNocqUwkMb',
'gHUeRsYtF14BaX7So8W0',
'hECVdzWU83LJkc46iGMR',
'7Z54Kq6E2GTvgbyXunm8',
'CTKbnhZkRy7AGstxePSd',
'ebI3LHRqBUwco8aQtW15',
'fQRdkX2ahē0t0c5ysNU8',
'Fbdqw108uxIcZ6rDMgmY',
'CYEUPAcItFWBs2VQTyKb',
'h1nZq0BX5LruizPw3EUm',
'1C4sdkhcZS6zYaE53D1v',
'gDYeNAplMujwPtU0qnJW',
'dVSQ8bXUG0Jo13ZmDqOp',
'e0ujNRMU2Txqd4rfAXaF',
'H6q1L1JXFrbITo8MQdmO',
'btaERgnQGZMzrDWvkSoU',
'4m1Gsx5JvMFgATY7Loz0',
'5z6qBep3fIQuW7dZhYAP',
'80rsVW03Jp4LYekTAmZU',
'biWpRxwyFrmULSzK9ZIv',
'2qk8R3jeQvKzsMp0iGT1',
'JSoqueuhRdwMs3G76lQp8',
'EKH1PBcSCXqnJzNRi8uj',
'BqhT4IuJponYPSKRxk2Q',
'BMjARP1c6GI7ONfitDJd',
'atGqPui5x1RbN6FDWhsZ',
'joFHWBAZySv18Ea1LD4N',
'3JI42LUuZET9tgvlYDqx',
'jGqYz9DaFMOWU1SwJBge',
'74SYbIefnAEJpGHZq6md',
'hSNuWGB9KZteamA7tF4R',
'3nZODLNAGgEF76dck14t',
'DhWInBmVuLfEF18YRPHZ',
'gCJk0uacmv1BAosn74b6',
'IR1CwpWFqsnd8gJZtbca',
'gKvtaoZ7Lk11V6yxq3Fc',
'gRNQ2Kw6dH10XrhBa0Vb',
'iCN4j126qfaQxTAe5Jzh',
'iqSDCrEjpfoRdaXHOIG4',
'd7aHABNQz41Sni2XyVUJ',
'jxDGoTHUWQeLfakytd05',
'aud9U7Xni4qS0cTrmjOR',

'gRZ9nxcONPftvzDCyKTQ',
'JstBaYEw5hQAC8kzcUMx',
'k9LJAopKrhzt8H3iIDuf',
'cXuqeLQGf6dyh4ZIUgFV',
'IZeQP4mMEwxvrzfUd6R',
'HZ1WSp9dEVFumCsKg8Rv',
'2JfpcM9ACq4XatFLho8z',
'BwN6PaW0ZG1QA3DgICEx',
'dw4sDcJMQCVOEgFjxXU0',
'bRoM4EpjukSGXwA1Z6zn',
'0DM3hS6Gg2QVKb1fZydv',
'h1sIMzUVPx aBEX1nbdpS',
'ieUqYbvXQJhLB4NDm097',
'IQJlrKuUAyj4seftWL2b',
'3suymN9GSHAtvYZ71Kaw',
'fojWI sY21tQG8K6BCkgV',
'bfRpQsFeqZN5BVrhkD6S',
'7mGrwqj2tdfEuNDn4gVR',
'1NvKF7pMREf4iVyzZewr',
'gZyCkLS1HF0x35REWhjm',
'7t2LuUCxAsgdq wWGIXY6',
'kQsiVxDbAXt23wRWa157',
'dTzmQ0cYijuhqwRtsPNe',
'jsSM7FgazoEt1V80ynmY',
'czBpmyW2hKAiGLDbv0Uq',
'4f91UESFtRDrnLdyZOCN',
'k4U7bfepOEJzgK8S231X',
'K0ij9XYNHeJR37b6naD2',
'JOWv6HNY1xwuqPrICEdR',
'BN76utiov1bnXza58ZTJ',
'8sNQwJpT43ohCx2MjHaD',
'hZBge41CJM uA60blif7P',
'4S5AJVNXiIrfZG7na0KQ',
'BKDWzzfquvWotCTXOM15',
'EgJw2tGHKzz3YeR85aID',
'BJxWK1M3u740AEkSoysr',
'BFmbPXnMt19oVyI1SYvT',
'eoA1Nd32LZGSPITn6mTM',
'd0j39TQiauSXbpsINeU1',
'16KuXyrdDOCW Ghc7P4kL',
'cUKzvqVwXx914BprdFgP',
'K9kNdLS83clsP0wDFObZ',
'4sCh00eqidfLtS83JMVE',
'bkyRaZK7X11NBewmDJGV',
'8026Dh4VpfjPekaCgAYQ',
'Fe4yEv8HmfLTBQuigKC6',
'57v6nKMC81hLV34je2RW',
'6hUFQK1g0x72GzZjJyI5',
'hpRHqfonKgzAbd83XVUI',
'HgcNkyo38ZnXM5ePvLQT',
'3n1LAd1CtkErg090IFFx',
'7P5gA4ZdSNwuYCi8lqpa',
'dQ5kFDwP4ByqU3tnLTZf',
'B0iZlJ2qwUe9RvbOSzGC',
'a1EOGLX5tBWIAMk2qTcY',
'cVjkXNizAryxYd1R5oBL',
'C4VIpEJ5K1vuetjhLdX3',
'9HAkGKpDrq6hecQZysNd',
'I9zqNg8Bj7RECDreiUX4',
'K7AwxRfWPnc3hzXLatQE',
'AsdcT125YI6a8SxzPbhf',

```
'a1tHxq97wypZ3A8QjhK1',
'4pehWH6b9IgL3ZF7uJCG',
'A4i9wPlygmMQK2ftbXaI',
'BFuiieglsUYJ934Dh0kw8',
'doEyIsH3nRDTjfZ5lvY0',
'hTe1goq5BG2YU7kQiuFm',
'5qEodvVi3umUTpZfskhP',
'27o93DUCHcGPXdTx1Nwg',
'eS81zGXw3KYrBUnhm65a',
'9mHqfWkgcJb7utKZ5rCV',
'gybTWSU3EjBGPcVI1JA4',
'6HcXhAzQiBMUK5en92FY',
'B3F57b1RXv8cuGZifynU',
'iHZnpXNuLTUQtE1Fz3jr',
'83uOwngTAXtNQPvD2CxY',
'dsuIp14XKzEt6kU29yeV',
'ImBUE37WLcf9PJ4kgZn1',
'7schXv12FHSujB3aVK9e',
'dvVcMR0BzWxa8hAyZNHg',
'CrKzvf0VWbetkShcayF8',
'4S0KT7xLGJ91CzYfMNpI',
'k4tKCVH9XwxuyieA2cG3',
'4skKNvbjShmzaWHIRQXi',
'0LAXajqhQy7po16dw8Tx',
'7KVtWPn30iJpvBkIej4o',
'iED1fIxnxGyWRLuYNv19s',
...]
```

In [17]:

```
vectorizer = CountVectorizer(input='filename', ngram_range=(2, 2), stop_words=None, token_pattern=r'\b\w\w\b', min_df=1)
```

In [18]:

```
X = vectorizer.fit_transform(flist)
```

In [19]:

```
X.shape
```

Out[19]:

```
(10868, 65536)
```

In [11]:

```
a=np.load('byte2gram.npy')
```

In [12]:

```
import numpy as np
from scipy import sparse
b=sparse.csr_matrix(a)
```

In [13]:

```
from sklearn.preprocessing import normalize
X_trasnfomed=normalize(b, axis = 0)
```

In [14]:

```
p=pd.read_csv('Final_normalized_bigram.csv')
```

In [15]:

```
col=list(p.columns)
col.remove('Unnamed: 0')
col.remove('ID')
```

In [16]:

```
Final_bigram_normalized=pd.DataFrame.sparse.from_spmatrix(X_trasnformed,columns=col)
```

In [17]:

```
Final_bigram_normalized['ID']=ids
```

In [18]:

```
del a
del b
del p
del X_trasnformed
```

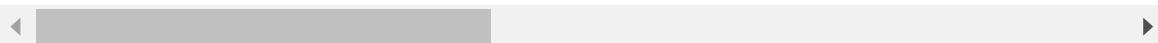
In [19]:

```
Final_bigram_normalized.head()
```

Out[19]:

	00 00	00 01	00 02	00 03	00 04	00 05	00 06	00 07	00 08
0	0.008655	0.007061	0.002831	0.002624	0.000860	0.001050	0.001247	0.003206	0.002425
1	0.017421	0.008077	0.008150	0.008867	0.009652	0.010989	0.010693	0.009714	0.011313
2	0.004986	0.002313	0.001460	0.002982	0.003311	0.001295	0.001021	0.001119	0.001182
3	0.000462	0.000201	0.000063	0.001770	0.000051	0.000105	0.000060	0.000057	0.000046
4	0.001982	0.001197	0.000203	0.000637	0.000616	0.000157	0.000135	0.000550	0.000184

5 rows × 65537 columns



In [20]:

```
byte_features_with_bigram = Final_bigram_normalized.merge(data_size_byte, on='ID')
byte_features_with_bigram.head()
```

Out[20]:

	00 00	00 01	00 02	00 03	00 04	00 05	00 06	00 07	00 08
0	0.008655	0.007061	0.002831	0.002624	0.000860	0.001050	0.001247	0.003206	0.002425
1	0.017421	0.008077	0.008150	0.008867	0.009652	0.010989	0.010693	0.009714	0.011313
2	0.004986	0.002313	0.001460	0.002982	0.003311	0.001295	0.001021	0.001119	0.001182
3	0.000462	0.000201	0.000063	0.001770	0.000051	0.000105	0.000060	0.000057	0.000046
4	0.001982	0.001197	0.000203	0.000637	0.000616	0.000157	0.000135	0.000550	0.000184

5 rows × 65539 columns

In [21]:

```
byte_features_with_bigram.drop('size',axis=1,inplace=True)
```

In [22]:

```
col=byte_features_with_bigram.drop(['ID','Class'],axis=1).columns
```

In [23]:

```
len(col)
```

Out[23]:

65536

For selecting top 300 byte features

In [30]:

```
clf = RandomForestClassifier(n_estimators = 100, n_jobs =4)
clf.fit(byte_features_with_bigram.drop(['ID','Class'],axis=1), byte_features_with_bigram['Class'])
```

Out[30]:

```
RandomForestClassifier(n_jobs=4)
```

In [36]:

```
imp_feature_indx = np.argsort(clf.feature_importances_)[-1][:300]
```

In [37]:

```
col_big=np.take(col,imp_feature_indx)
```

In [40]:

```
final_bigram=byte_features_with_bigram.loc[:,['ID','Class']+list(col_big)]
```

Custom ByteFiles Bigram

I used it and was getting accurate result with this implementation but when i was try to normalize the data,it was showing 48 hrs to complete so i used sciketlearn countvectorizer

In []:

```

def secondprocess():
    #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
    byte_vocab = "00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,
17,18,19,1a,1b,1c,1d,1e,1f,20,21,22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,
34,35,36,37,38,39,3a,3b,3c,3d,3e,3f,40,41,42,43,44,45,46,47,48,49,4a,4b,4c,4d,4e,4f,50,
51,52,53,54,55,56,57,58,59,5a,5b,5c,5d,5e,5f,60,61,62,63,64,65,66,67,68,69,6a,6b,6c,6d,
6e,6f,70,71,72,73,74,75,76,77,78,79,7a,7b,7c,7d,7e,7f,80,81,82,83,84,85,86,87,88,89,8a,
8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,9a,9b,9c,9d,9e,9f,a0,a1,a2,a3,a4,a5,a6,a7,
a8,a9,aa,ab,ac,ad,ae,af,b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,ba,bb,bc,bd,be,bf,c0,c1,c2,c3,c4,
c5,c6,c7,c8,c9,ca,cb,cc,cd,ce,cf,d0,d1,d2,d3,d4,d5,d6,d7,d8,d9,da,db,dc,dd,de,df,e0,e1,
e2,e3,e4,e5,e6,e7,e8,e9,ea,eb,ec,ed,ee,ef,f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,
ff,??"
    vocab = []
    for i, v in enumerate(byte_vocab.split(',')):
        for j in ((byte_vocab.split(','))):
            vocab.append((v,j))

    files = os.listdir('second')

    file=open('result2.csv', 'w+', newline='')
    writer = csv.DictWriter(file,fieldnames=[ 'ID']+vocab)
    writer.writeheader()

    for f in files:
        f2=f.split('.')[0]
        t=dict.fromkeys(vocab,0)
        d=dict()
        with open('second/'+ f) as fli:
            for lines in fli:
                value=lines.rstrip().lower()
                for i, v in enumerate(value.split()):
                    for j in value.split():
                        t[(v,j)]+=1

    #pushing the values into the file after reading whole file
    d['ID']=f2
    for z,w in t.items():
        d[z]=w

    writer.writerow(d)
    print('wrote successfully')

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
    byte_vocab = "00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,
17,18,19,1a,1b,1c,1d,1e,1f,20,21,22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,
```

```

34,35,36,37,38,39,3a,3b,3c,3d,3e,3f,40,41,42,43,44,45,46,47,48,49,4a,4b,4c,4d,4e,4f,50,
51,52,53,54,55,56,57,58,59,5a,5b,5c,5d,5e,5f,60,61,62,63,64,65,66,67,68,69,6a,6b,6c,6d,
6e,6f,70,71,72,73,74,75,76,77,78,79,7a,7b,7c,7d,7e,7f,80,81,82,83,84,85,86,87,88,89,8a,
8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,9a,9b,9c,9d,9e,9f,a0,a1,a2,a3,a4,a5,a6,a7,
a8,a9,aa,ab,ac,ad,ae,af,b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,ba,bb,bc,bd,be,bf,c0,c1,c2,c3,c4,
c5,c6,c7,c8,c9,ca,cb,cc,cd,ce,cf,d0,d1,d2,d3,d4,d5,d6,d7,d8,d9,da,db,dc,dd,de,df,e0,e1,
e2,e3,e4,e5,e6,e7,e8,e9,ea,eb,ec,ed,ee,ef,f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,
ff,??"
```

```

vocab = []
for i, v in enumerate(byte_vocab.split(',')):
    for j in ((byte_vocab.split(','))):
        vocab.append((v,j))

files = os.listdir('first')

file=open('result1.csv', 'w+', newline='')
writer = csv.DictWriter(file,fieldnames=[ 'ID']+vocab)
writer.writeheader()

for f in files:
    f2=f.split('.')[0]
    t=dict.fromkeys(vocab,0)
    d=dict()
    with open('first/'+ f) as fli:
        for lines in fli:
            value=lines.rstrip().lower()
            for i, v in enumerate(value.split()):
                for j in value.split():
                    t[(v,j)]+=1

#pushing the values into the file after reading whole file
d['ID']=f2
for z,w in t.items():
    d[z]=w

writer.writerow(d)
print('wrote successfully')

def thirdprocess():
    #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
    byte_vocab = "00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,
17,18,19,1a,1b,1c,1d,1e,1f,20,21,22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,
34,35,36,37,38,39,3a,3b,3c,3d,3e,3f,40,41,42,43,44,45,46,47,48,49,4a,4b,4c,4d,4e,4f,50,
51,52,53,54,55,56,57,58,59,5a,5b,5c,5d,5e,5f,60,61,62,63,64,65,66,67,68,69,6a,6b,6c,6d,
6e,6f,70,71,72,73,74,75,76,77,78,79,7a,7b,7c,7d,7e,7f,80,81,82,83,84,85,86,87,88,89,8a,
8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,9a,9b,9c,9d,9e,9f,a0,a1,a2,a3,a4,a5,a6,a7,
a8,a9,aa,ab,ac,ad,ae,af,b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,ba,bb,bc,bd,be,bf,c0,c1,c2,c3,c4,
c5,c6,c7,c8,c9,ca,cb,cc,cd,ce,cf,d0,d1,d2,d3,d4,d5,d6,d7,d8,d9,da,db,dc,dd,de,df,e0,e1,
e2,e3,e4,e5,e6,e7,e8,e9,ea,eb,ec,ed,ee,ef,f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,
ff,??"
```

```

vocab = []
```

```

for i, v in enumerate(byte_vocab.split(',')):
    for j in ((byte_vocab.split(','))):
        vocab.append((v,j))

files = os.listdir('third')

file=open('result3.csv', 'w+', newline='')
writer = csv.DictWriter(file,fieldnames=['ID']+vocab)
writer.writeheader()

for f in files:
    f2=f.split('.')[0]
    t=dict.fromkeys(vocab,0)
    d=dict()
    with open('third/'+ f) as fli:
        for lines in fli:
            value=lines.rstrip().lower()
            for i, v in enumerate(value.split()):
                for j in value.split():
                    t[(v,j)]+=1

#pushing the values into the file after reading whole file
d['ID']=f2
for z,w in t.items():
    d[z]=w

writer.writerow(d)
print('wrote successfully')

def fourthprocess():
    #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
    byte_vocab = "00,01,02,03,04,05,06,07,08,09,0a,0b,0c,0d,0e,0f,10,11,12,13,14,15,16,
17,18,19,1a,1b,1c,1d,1e,1f,20,21,22,23,24,25,26,27,28,29,2a,2b,2c,2d,2e,2f,30,31,32,33,
34,35,36,37,38,39,3a,3b,3c,3d,3e,3f,40,41,42,43,44,45,46,47,48,49,4a,4b,4c,4d,4e,4f,50,
51,52,53,54,55,56,57,58,59,5a,5b,5c,5d,5e,5f,60,61,62,63,64,65,66,67,68,69,6a,6b,6c,6d,
6e,6f,70,71,72,73,74,75,76,77,78,79,7a,7b,7c,7d,7e,7f,80,81,82,83,84,85,86,87,88,89,8a,
8b,8c,8d,8e,8f,90,91,92,93,94,95,96,97,98,99,9a,9b,9c,9d,9e,9f,a0,a1,a2,a3,a4,a5,a6,a7,
a8,a9,aa,ab,ac,ad,ae,af,b0,b1,b2,b3,b4,b5,b6,b7,b8,b9,ba,bb,bc,bd,be,bf,c0,c1,c2,c3,c4,
c5,c6,c7,c8,c9,ca,cb,cc,cd,ce,cf,d0,d1,d2,d3,d4,d5,d6,d7,d8,d9,da,db,dc,dd,de,df,e0,e1,
e2,e3,e4,e5,e6,e7,e8,e9,ea,eb,ec,ed,ee,ef,f0,f1,f2,f3,f4,f5,f6,f7,f8,f9,fa,fb,fc,fd,fe,
ff,??""
    vocab = []
    for i, v in enumerate(byte_vocab.split(',')):
        for j in ((byte_vocab.split(','))):
            vocab.append((v,j))

files = os.listdir('fourth')

```

```
file=open('result4.csv', 'w+', newline=' ')
writer = csv.DictWriter(file,fieldnames=[ 'ID']+vocab)
writer.writeheader()

for f in files:
    f2=f.split('.')[0]
    t=dict.fromkeys(vocab,0)
    d=dict()
    with open('fourth/'+ f) as fli:
        for lines in fli:
            value=lines.rstrip().lower()
            for i, v in enumerate(value.split()):
                for j in value.split():
                    t[(v,j)]+=1

    #pushing the values into the file after reading whole file
    d['ID']=f2
    for z,w in t.items():
        d[z]=w

    writer.writerow(d)
print('wrote successfully')
```

In []:

In []:

In []:

```
#initially create five folders
#first
#second
#third
#fourth
#fifth
#this code tells us about random split of files into five folders
folder_1 ='first'
folder_2 ='second'
folder_3 ='third'
folder_4 ='fourth'

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

source='byteFiles/'
files = os.listdir('byteFiles')

for i in range(0,10868):
    if i <2717:
        shutil.move(source+files[i],'first')
    elif 2717<=i<5434:
        shutil.move(source+files[i],'second')
    elif 5434<=i<8151:
        shutil.move(source+files[i],'third')
    elif 8151<=i<10868:
        shutil.move(source+files[i],'fourth')
```

In []:

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

if __name__=="__main__":
    main()
```

In []:

```
df1=pd.read_csv('result1.csv')
df2=pd.read_csv('result2.csv')
df3=pd.read_csv('result3.csv')
df4=pd.read_csv('result4.csv')
```

In []:

```
frames = [df1, df2, df3, df4]
```

In []:

```
result = pd.concat(frames, ignore_index=True)
```

In []:

```
result.to_csv(r'Final_bigram.csv')
```

Custom asm file implementation

In []:

In []:

```
#initially create five folders
#first
#second
#third
#fourth
#fifth
#this code tells us about random split of files into five folders
folder_1 ='first'
folder_2 ='second'
folder_3 ='third'
folder_4 ='fourth'

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

source='asmFiles/'
files = os.listdir('asmFiles')

for i in range(0,10868):
    if i <2717:
        shutil.move(source+files[i],'first')
    elif 2717<=i<5434:
        shutil.move(source+files[i],'second')
    elif 5434<=i<8151:
        shutil.move(source+files[i],'third')
    elif 8151<=i<10868:
        shutil.move(source+files[i],'fourth')
```

In []:

```
os.makedirs(str('asm_image'))
```

Got this from :[\(https://github.com/kunwar-vikrant/Microsoft-Malware-Detection/blob/master/A%2317.pdf\)](https://github.com/kunwar-vikrant/Microsoft-Malware-Detection/blob/master/A%2317.pdf)

In []:

```

import array

import imageio

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

```

In []:

```

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=collect_img_asm,args=('first',))
    p2=Process(target=collect_img_asm,args=('second',))
    p3=Process(target=collect_img_asm,args=('third',))
    p4=Process(target=collect_img_asm,args=('fourth',))
    # p5=Process(target=fifthprocess)
    #p1.start() is used to start the thread execution
    p1.start()
    p2.start()
    p3.start()
    p4.start()
    #p5.start()
    #After completion all the threads are joined
    p1.join()
    p2.join()
    p3.join()
    p4.join()
    #p5.join()

if __name__=="__main__":
    main()

```

In []:

```

source='first/'
files = os.listdir('first')

for i in range(0,len(os.listdir('first'))):
    shutil.move(source+files[i], 'asmFiles')

```

In []:

```
source='second/'
files = os.listdir('second')

for i in range(0,len(os.listdir('second'))):
    shutil.move(source+files[i], 'asmFiles')
```

In []:

```
source='third/'
files = os.listdir('third')

for i in range(0,len(os.listdir('third'))):
    shutil.move(source+files[i], 'asmFiles')
```

In []:

```
source='fourth/'
files = os.listdir('fourth')

for i in range(0,len(os.listdir('fourth'))):
    shutil.move(source+files[i], 'asmFiles')
```

In []:

```
import cv2
image_features = 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]
    image_features[i]= img_arr
```

In []:

```
from sklearn.preprocessing import normalize
img_feat = []
for i in range(200):
    img_feat.append('pix' + str(i))
img_final = pd.DataFrame(normalize(image_features, axis = 0), columns = img_feat)
```

In []:

```
Id=[]
for i, asmfile in enumerate(os.listdir("asmFiles")):
    Id.append(asmfile.split('.')[0])
img_final['ID'] = Id
```

final merge

In [2]:

```
data_1=pd.read_csv('final_bigram.csv')
```

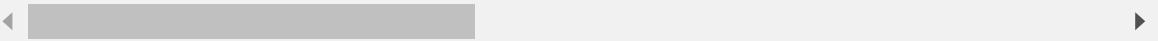
In [3]:

```
data_1
```

Out[3]:

	Unnamed: 0		ID	Class	1e 6e	e3 bc	d7 ed	29 86
0	0	9MW5Nuf0ogCEcRIYJeKG	8	0.000000	0.000000	0.000000	0.000000	0.000000
1	1	GFblksovaPYLI5XeC8RU	9	0.000090	0.005745	0.006436	0.000114	
2	2	kjQFMnf7vADCqtX4ai2u	6	0.000030	0.003830	0.002145	0.000160	
3	3	EKmgShQsf6a9vY0znNIU	4	0.000090	0.001915	0.003218	0.000068	
4	4	1WCXg8qEtJFITMilaf6k	6	0.000060	0.001915	0.000000	0.000000	
...
10863	10863	ExiOZLUCAAsSm3rogNy2D	3	0.000195	0.013405	0.011799	0.000228	
10864	10864	GWBugsn5QySAo76NDE0H	2	0.000000	0.000000	0.000000	0.000023	
10865	10865	A1Bn4Cb86GzdTkf0xo5W	1	0.000000	0.000958	0.001073	0.000046	
10866	10866	FvobzqOYfnQ0klx1tZLc	3	0.000180	0.011490	0.012872	0.000319	
10867	10867	hf8ecIPDq0WrSU6JZjoy	3	0.000286	0.009575	0.017163	0.000319	

10868 rows × 303 columns



In [4]:

```
data_2=pd.read_csv('img_features.csv')
```

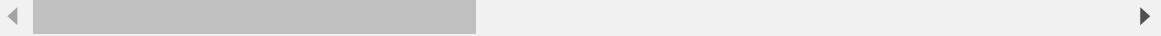
In [5]:

```
data_2
```

Out[5]:

	Unnamed: 0	pix0	pix1	pix2	pix3	pix4	pix5	pix6	
0	0	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
1	1	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
2	2	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
3	3	0.006560	0.006560	0.006560	0.013504	0.013504	0.013504	0.012927	0.0
4	4	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
...
10863	10863	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
10864	10864	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
10865	10865	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
10866	10866	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0
10867	10867	0.010268	0.010268	0.010268	0.008033	0.008033	0.008033	0.008320	0.0

10868 rows × 202 columns



In [6]:

```
result_temp = pd.merge(data_1,data_2,on='ID', how='left')
```

In [8]:

```
result_temp
```

Out[8]:

	Unnamed: 0_x	ID	Class	1e 6e	e3 bc	d7 ed	29 86
0	0	9MW5Nuf0ogCEcRlYJeKG	8	0.000000	0.000000	0.000000	0.000000
1	1	GFblksovaPYLl5XeC8RU	9	0.000090	0.005745	0.006436	0.000114
2	2	kjQFMnf7vADCqtX4ai2u	6	0.000030	0.003830	0.002145	0.000160
3	3	EKmgShQsf6a9vY0znNIU	4	0.000090	0.001915	0.003218	0.000068
4	4	1WCXg8qEtJFITMilaf6k	6	0.000060	0.001915	0.000000	0.000000
...
10863	10863	ExiOZLUCAsSm3rogNy2D	3	0.000195	0.013405	0.011799	0.000228
10864	10864	GWBugsn5QySAo76NDE0H	2	0.000000	0.000000	0.000000	0.000023
10865	10865	A1Bn4Cb86GzdTkf0xo5W	1	0.000000	0.000958	0.001073	0.000046
10866	10866	FvobzqOYfnQ0klx1tZLc	3	0.000180	0.011490	0.012872	0.000319
10867	10867	hf8ecIPDq0WrSU6JZjoy	3	0.000286	0.009575	0.017163	0.000319

10868 rows × 504 columns

In [32]:

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

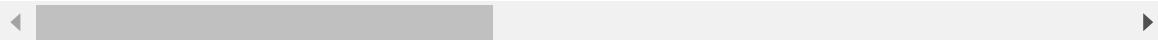
In [33]:

```
result_x
```

Out[33]:

	ID	0	1	2	3	4	5
0	01azqd4InC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835
1	01IsoiSMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873
2	01jsnpXSAlg6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354
4	01SuzwMJEIXsK7A8dQbI	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232
...
10863	IoIP1tiwELF9YNZQjSUO	0.002300	0.001657	0.000596	0.000659	0.000758	0.000656
10864	LOP6HaJKXpkic5dyuVnT	0.001324	0.000420	0.000138	0.000158	0.000168	0.000121
10865	LOqA6FX02GWguYrl1Zbe	0.002476	0.000311	0.000150	0.000174	0.000192	0.000088
10866	LoWgaidpb2lUM5ACcSGO	0.001588	0.000615	0.000252	0.000273	0.000313	0.000221
10867	IS0IVqXeJrN6Dzi9Pap1	0.001543	0.000525	0.000214	0.000233	0.000303	0.000226

10868 rows × 312 columns



In [34]:

```
result_x=pd.merge(result_temp,result_x,on='ID', how='left')
```

In [35]:

result_x

Out[35]:

		Unnamed: 0_x	ID	Class_x	1e 6e	e3 bc	d7 ed	29 86
0	0	9MW5Nuf0ogCEcRlYJeKG	8	0.000000	0.000000	0.000000	0.000000	0.000000
1	1	GFblksovaPYLl5XeC8RU	9	0.000090	0.005745	0.006436	0.000111	0.000000
2	2	kjQFMnf7vADCqtX4ai2u	6	0.000030	0.003830	0.002145	0.000166	0.000000
3	3	EKmgShQsf6a9vY0znNIU	4	0.000090	0.001915	0.003218	0.000000	0.000000
4	4	1WCXg8qEtJFITMilaf6k	6	0.000060	0.001915	0.000000	0.000000	0.000000
...
10863	10863	ExiOZLUCAAsSm3rogNy2D	3	0.000195	0.013405	0.011799	0.000212	0.000000
10864	10864	GWBugsn5QySAo76NDE0H	2	0.000000	0.000000	0.000000	0.000000	0.000000
10865	10865	A1Bn4Cb86GzdTkf0xo5W	1	0.000000	0.000958	0.001073	0.000042	0.000000
10866	10866	FvobzqOYfnQ0klx1tZLc	3	0.000180	0.011490	0.012872	0.000311	0.000000
10867	10867	hf8ecIPDq0WrSU6JZjoy	3	0.000286	0.009575	0.017163	0.000311	0.000000

10868 rows × 815 columns

◀ ▶

In [36]:

```
result_y = result_x['Class_x']
result_x = result_x.drop(['ID', 'rtn', '.BSS:', '.CODE', 'Class_x', 'Unnamed: 0_x'], axis=1)
result_x.head()
```

Out[36]:

	1e 6e	e3 bc	d7 ed	29 86	82 cd	fd 87	4e 47	f9 08	41 44
0	0.00000	0.000000	0.000000	0.000000	0.000751	0.000626	0.000000	0.000215	0.000227
1	0.00009	0.005745	0.006436	0.000114	0.005255	0.005010	0.000398	0.007326	0.000604
2	0.00003	0.003830	0.002145	0.000160	0.004504	0.001879	0.000133	0.001724	0.000604
3	0.00009	0.001915	0.003218	0.000068	0.000751	0.000626	0.000133	0.000215	0.000076
4	0.00006	0.001915	0.000000	0.000000	0.001501	0.000626	0.013235	0.000215	0.023039

5 rows × 809 columns

◀ ▶

In [37]:

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

Final model

In [39]:

```
x_cfl=XGBClassifier()

params={
    '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=params,verbose=10,n_jobs=-1,)
random_cfl.fit(X_train_merge, y_train_merge)
```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done   2 tasks    | elapsed:  1.1min
[Parallel(n_jobs=-1)]: Done   9 tasks    | elapsed:  5.3min
[Parallel(n_jobs=-1)]: Done  16 tasks    | elapsed: 13.1min
[Parallel(n_jobs=-1)]: Done  25 tasks    | elapsed: 21.2min
[Parallel(n_jobs=-1)]: Done  34 tasks    | elapsed: 23.0min
[Parallel(n_jobs=-1)]: Done  41 out of  50 | elapsed: 25.1min remaining: 5.5min
[Parallel(n_jobs=-1)]: Done  47 out of  50 | elapsed: 30.3min remaining: 1.9min
[Parallel(n_jobs=-1)]: Done  50 out of  50 | elapsed: 30.9min finished
```

Out[39]:

```
RandomizedSearchCV(estimator=XGBClassifier(base_score=None, booster=None,
                                         colsample_bylevel=None,
                                         colsample_bynode=None,
                                         colsample_bytree=None, gamma=None,
                                         gpu_id=None, importance_type='gain',
                                         interaction_constraints=None,
                                         learning_rate=None,
                                         max_delta_step=None, max_depth=None,
                                         min_child_weight=None, missing=None,
                                         monotone_constraints=None,
                                         n_estimators=100, n_jobs=-1,
                                         random_state=None, reg_alpha=None,
                                         reg_lambda=None,
                                         scale_pos_weight=None,
                                         subsample=None, tree_method=None,
                                         validate_parameters=None,
                                         verbosity=None),
                     n_jobs=-1,
                     param_distributions={'colsample_bytree': [0.1, 0.3, 0.5, 1],
                                         'learning_rate': [0.01, 0.03, 0.05, 0.1,
                                                          0.15, 0.2],
                                         'max_depth': [3, 5, 10],
                                         'n_estimators': [100, 200, 500, 1000,
                                                          2000],
                                         'subsample': [0.1, 0.3, 0.5, 1]},
                     verbose=10)
```

In [40]:

```
print(random_cfl.best_params_)

{'subsample': 0.3, 'n_estimators': 200, 'max_depth': 3, 'learning_rate': 0.1, 'colsample_bytree': 0.5}
```

In []:

```
# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
# -----
# 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_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
# s)

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


x_cfl=XGBClassifier(n_estimators=200,max_depth=3,learning_rate=0.1,colsample_bytree=0.5
,subsample=0.3,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)
```

In [42]:

```
predict_y = sig_clf.predict_proba(X_train_merge)
print ("The train log loss is:",log_loss(y_train_merge, predict_y))
predict_y = sig_clf.predict_proba(X_cv_merge)
print("The cross validation log loss is:",log_loss(y_cv_merge, predict_y))
predict_y = sig_clf.predict_proba(X_test_merge)
print("The test log loss is:",log_loss(y_test_merge, predict_y))
plot_confusion_matrix(y_test_merge,sig_clf.predict(X_test_merge))
```

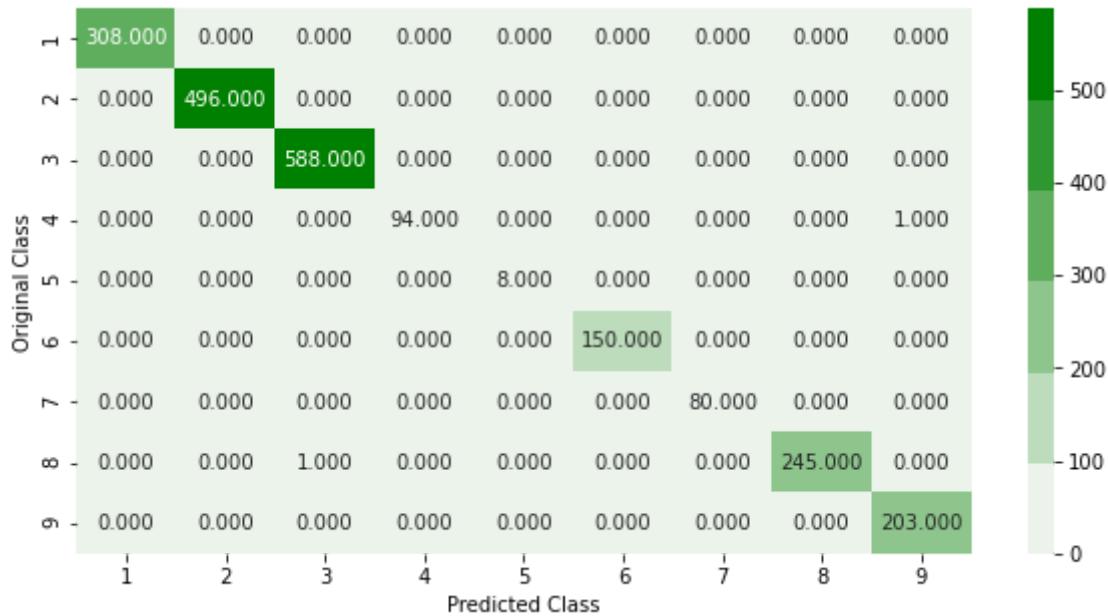
The train log loss is: 0.007798047108288545

The cross validation log loss is: 0.007759574356459539

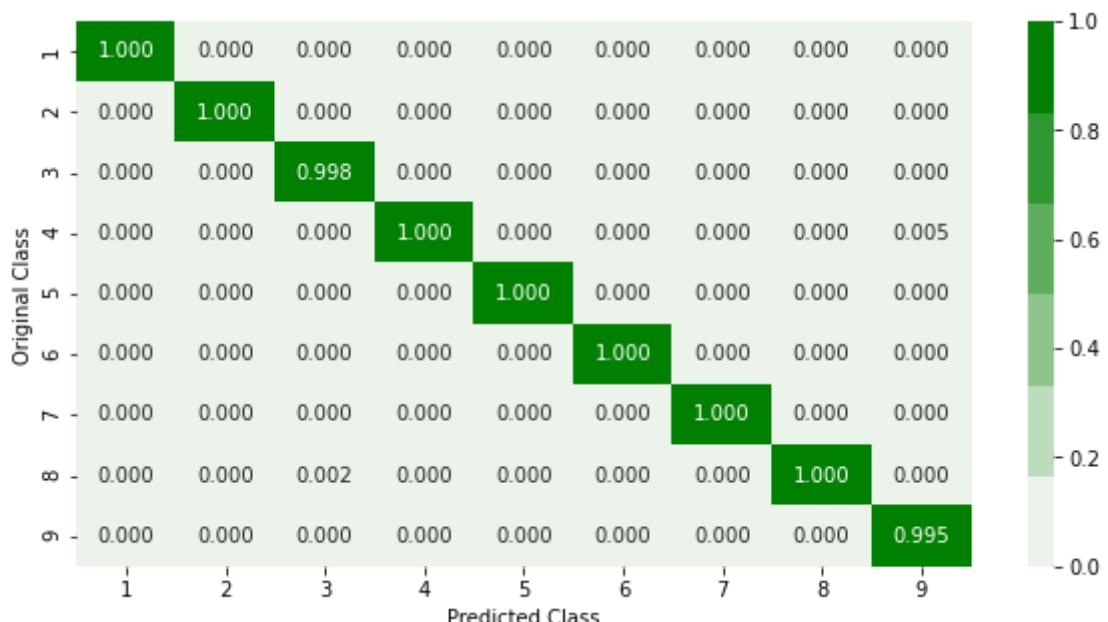
The test log loss is: 0.010942646776749578

Number of misclassified points 0.09199632014719411

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

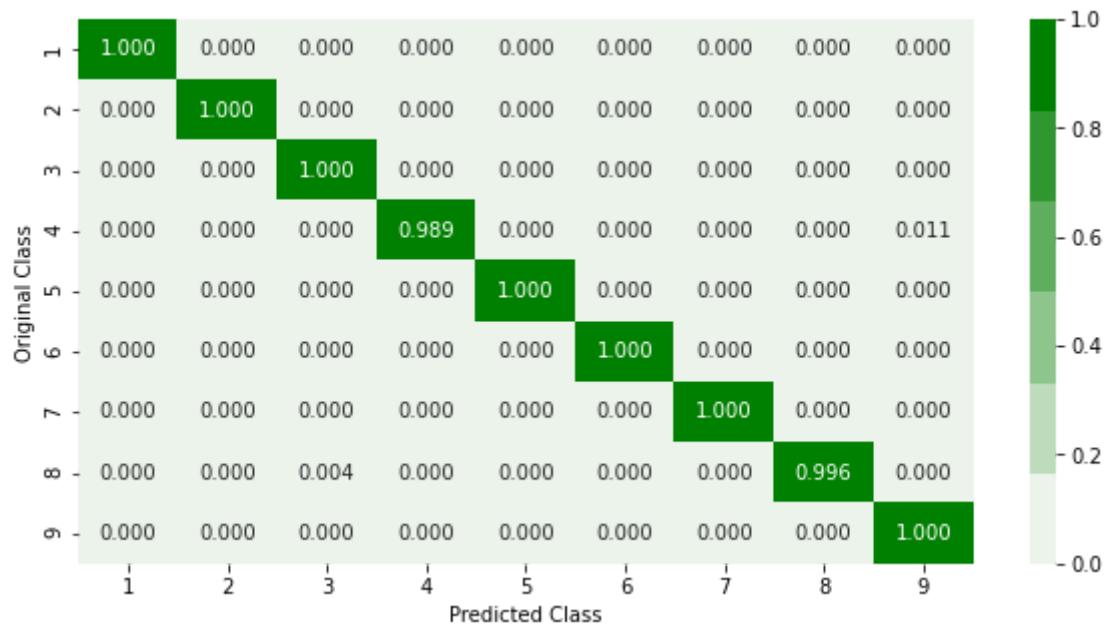


----- Precision matrix -----



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

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



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

In []: