

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

```
In [2]: 1 def plot_confusion_matrix(test_y, predict_y):
2     C = confusion_matrix(test_y, predict_y)
3     print("Number of misclassified points ",(len(test_y)-np.trace(C))/len(test_y)*100)
4     # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j
5
6     A =(((C.T)/(C.sum(axis=1))).T)
7     #divid each element of the confusion matrix with the sum of elements in that column
8
9
10    # C = [[1, 2],
11          # [3, 4]]
12    # C.T = [[1, 3],
13            # [2, 4]]
14    # C.sum(axis = 1) axis=0 corresonds to columns and axis=1 corresponds to rows in two dimensional array
15    # C.sum(axix =1) = [[3, 7]]
16    # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
17                                # [2/3, 4/7]]
18
19    # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
20                                   # [3/7, 4/7]]
21    # sum of row elements = 1
22
23    B =(C/C.sum(axis=0))
24    #divid each element of the confusion matrix with the sum of elements in that row
25    # C = [[1, 2],
26          # [3, 4]]
27    # C.sum(axis = 0) axis=0 corresonds to columns and axis=1 corresponds to rows in two dimensional array
28    # C.sum(axix =0) = [[4, 6]]
29    # (C/C.sum(axis=0)) = [[1/4, 2/6],
30                           # [3/4, 4/6]]
31
32    labels = [1,2,3,4,5,6,7,8,9]
33    cmap=sns.light_palette("green")
34    # representing A in heatmap format
35    print("-"*50, "Confusion matrix", "-"*50)
36    plt.figure(figsize=(10,5))
37    sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
38    plt.xlabel('Predicted Class')
39    plt.ylabel('Original Class')
40    plt.show()
41
42    print("-"*50, "Precision matrix", "-"*50)
43    plt.figure(figsize=(10,5))
44    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
45    plt.xlabel('Predicted Class')
46    plt.ylabel('Original Class')
47    plt.show()
48    print("Sum of columns in precision matrix",B.sum(axis=0))
49
50    # representing B in heatmap format
51    print("-"*50, "Recall matrix" , "-"*50)
52    plt.figure(figsize=(10,5))
53    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
54    plt.xlabel('Predicted Class')
55    plt.ylabel('Original Class')
56    plt.show()
57    print("Sum of rows in precision matrix",A.sum(axis=1))
```

1.Image Extraction from ASM -files

Image Extraction using : <https://github.com/dchad/malware-detection/blob/master/mmcc/feature-extraction.ipynb> (<https://github.com/dchad/malware-detection/blob/master/mmcc/feature-extraction.ipynb>)

```
In [3]: 1 def entropy(p,n):
2     '''Returns entrophy'''
3     p_ratio = float(p)/(p+n)
4     n_ratio = float(n)/(p+n)
5     return -p_ratio*math.log(p_ratio) - n_ratio * math.log(n_ratio)

In [4]: 1 def info_gain(p0,n0,p1,n1,p,n):
2     '''Return information gain'''
3     return entropy(p,n) - float(p0+n0)/(p+n)*entropy(p0,n0) - float(p1+n1)/(p+n)*entropy(p1,n1)
```

```
In [5]: 1 def read_image(filename):
2     '''
3     Read the image data
4     '''
5     f = open(filename,'rb')
6     ln = os.path.getsize(filename) # Length of file in bytes
7     width = 256
8     rem = ln%width
9     a = array.array("B") # uint8 array
10    a.fromfile(f,ln-rem)
11    f.close()
12    g = np.reshape(a,(int(len(a)/width), width))
13    g = np.uint8(g)
14    g = np.resize(g, (1000,))
15    return list(g)
```

```
In [6]: 1 # Do asm image extraction
2 def extract_asm_image_features(tfiles):
3     asm_files = [i for i in tfiles if '.asm' in i]
4     ftot = len(asm_files)
5
6     pid = os.getpid()
7     print('Process id:', pid)
8     feature_file = str(pid) + '-image-features-asm.csv'
9     print('feature file:', feature_file)
10
11    outrows = []
12    with open(feature_file,'w') as f:
13        fw = writer(f)
14        column_names = ['filename'] + [("ASM_{:s}").format(str(x))] for x in range(1000)]
15        fw.writerow(column_names)
16        for idx, fname in enumerate(asm_files):
17            file_id = fname.split('.')[0]
18            image_data = read_image('asmFiles/' + fname)
19            outrows.append([file_id] + image_data)
20
21            # Print progress
22            if (idx+1) % 10 == 0:
23                print(pid, idx + 1, 'of', ftot, 'files processed.')
24                fw.writerows(outrows)
25                outrows = []
26
27            # Write remaining files
28            if len(outrows) > 0:
29                fw.writerows(outrows)
30            outrows = []
```

```
In [7]: 1 # TRAIN FILES ASM
2 # Now divide the train files into four groups for multiprocessing
3 import time
4 from multiprocessing import Pool
5 start_time = time.time()
6 tfiles = os.listdir('asmFiles')
7 quart = len(tfiles)//4
8 train1 = tfiles[:quart]
9 train2 = tfiles[quart:(2*quart)]
10 train3 = tfiles[(2*quart):(3*quart)]
11 train4 = tfiles[(3*quart):]
12 print(len(tfiles), quart, (len(train1)+len(train2)+len(train3)+len(train4)))
13 trains = [train1, train2, train3, train4]
14 p = Pool(4)
15 p.map(extract_asm_image_features, trains)
16 print("Elapsed time: {:.2f} hours.".format((time.time() - start_time)//3600.0))
```

10868 2717 10868
Process id:Process id:Process id:Process id: 1139114011381137

feature file:feature file:feature file:feature file: 1139-image-features-asm.csv1140-image-features-asm.csv1138-image-features-asm.csv1137-image-features-asm.csv

1138 10 of 2717 files processed.
1137 10 of 2717 files processed.
1140 10 of 2717 files processed.
1138 20 of 2717 files processed.
1139 10 of 2717 files processed.
1137 20 of 2717 files processed.
1138 30 of 2717 files processed.
1140 20 of 2717 files processed.
1139 20 of 2717 files processed.
1137 30 of 2717 files processed.
1140 30 of 2717 files processed.
1138 40 of 2717 files processed.
1138 50 of 2717 files processed.
1139 30 of 2717 files processed.
1140 40 of 2717 files processed.
1137 40 of 2717 files processed.
1139 40 of 2717 files processed.
1138 60 of 2717 files processed.
1137 50 of 2717 files processed.
1138 70 of 2717 files processed.
1139 50 of 2717 files processed.
1140 50 of 2717 files processed.
1140 60 of 2717 files processed.
1138 80 of 2717 files processed.
1139 60 of 2717 files processed.
1137 60 of 2717 files processed.
1140 70 of 2717 files processed.
1139 70 of 2717 files processed.
1137 70 of 2717 files processed.
1138 90 of 2717 files processed.
1139 80 of 2717 files processed.
1138 100 of 2717 files processed.
1140 80 of 2717 files processed.
1137 80 of 2717 files processed.
1140 90 of 2717 files processed.
1138 110 of 2717 files processed.
1137 901139 90 of 2717 2717files processed.
files processed.
1140 100 of 2717 files processed.
1138 120 of 2717 files processed.
1139 100 of 2717 files processed.
1137 100 of 2717 files processed.
1140 110 of 2717 files processed.
1139 110 of 2717 files processed.
1138 130 of 2717 files processed.
1137 110 of 2717 files processed.
1139 120 of 2717 files processed.
1140 120 of 2717 files processed.
1137 120 of 2717 files processed.
1139 130 of 2717 files processed.
1140 130 of 2717 files processed.
1138 140 of 2717 files processed.
1138 150 of 2717 files processed.
1140 140 of 2717 files processed.
1138 160 of 2717 files processed.
1137 130 of 2717 files processed.
1139 140 of 2717 files processed.
1138 170 of 2717 files processed.
1137 140 of 2717 files processed.
1140 150 of 2717 files processed.
1139 150 of 2717 files processed.
1138 180 of 2717 files processed.
1139 160 of 2717 files processed.
1138 190 of 2717 files processed.
1137 150 of 2717 files processed.
1139 170 of 2717 files processed.
1140 160 of 2717 files processed.
1139 180 of 2717 files processed.1140 170
of 2717 files processed.
1137 160 of 2717 files processed.
1139 190 of 27171138 files processed.200
of 2717 files processed.
1137 170 of 2717 files processed.
1140 180 of 2717 files processed.
1140 190 of 2717 files processed.
1139 200 of 2717 files processed.
1137 180 of 2717 files processed.
1138 210 of 2717 files processed.
1140 200 of 2717 files processed.
1138 220 of 2717 files processed.
1139 210 of 2717 files processed.
1137 190 of 2717 files processed.
1138 230 of 2717 files processed.
1140 210 of 2717 files processed.
1137 2001138 240 of 2717 files processed.2717
files processed.
1139 220 of 2717 files processed.
1140 220 of 2717 files processed.
1137 210 of 2717 files processed.
1139 230 1140of 2717 files processed.
230 of 2717 files processed.
1138 250 of 2717 files processed.
1139 240 of 2717 files processed.
1137 220 of 2717 files processed.
1138 260 of 2717 files processed.
1140 240 of 2717 files processed.
1139 250 of 2717 files processed.
1137 230 of 2717 files processed.
1140 250 of 2717 files processed.
1138 270 of 2717 files processed.
1137 240 of 2717 files processed.
1139 260 of 2717 files processed.
1139 270 of 2717 files processed.
1140 260 of 2717 files processed.
1138 280 of 2717 files processed.
1139 280 of 2717 files processed.
1139 290 of 2717 files processed.
1138 290 of 2717 files processed.
1137 250 of 2717 files processed.
1140 270 of 2717 files processed.
1137 260 of 2717 files processed.
1139 300 of 2717 files processed.
1138 300 of 2717 files processed.
1137 270 of 2717 files processed.
1139 310 of 2717 files processed.
1138 310 of 2717 files processed.
1140 280 of 2717 files processed.
1137 280 of 2717 files processed.
1140290 of 2717 files processed.
1138 320 of 2717 files processed.
1137 290 of 2717 files processed.
1140 300 of 2717 files processed.
1139 320 of 2717 files processed.
1137 300 of 2717 files processed.
1138 330 of 2717 files processed.
1139 330 of 2717 files processed.
1138 340 of 2717 files processed.
1137 310 of 2717 files processed.
1139 340 of 2717 files processed.
1140 310 of 2717 files processed.
1138 350 of 2717 files processed.
1137 320 of 2717 files processed.
1140 320 of 2717 files processed.
1138 360 of 2717 files processed.
1140 330 of 2717 files processed.
1139 350 of 2717 files processed.
1137 330 of 2717 files processed.
1138 370 of 2717 files processed.
1140 340 of 2717 files processed.
1137 340 of 2717 files processed.
1139 360 of 2717 files processed.
1139 370 of 2717 files processed.
1138 380 of 2717 files processed.
1140 350 of 2717 files processed.
1137 350 of 2717 1139files processed.
380 of 2717 files processed.
11391138 390 390of of 27172717 files processed. files processed.

```
1140 630 of 1138 670 2717 of files processed. 2717
files processed.
```


1138 820 of 2717 files processed.
1140 780 of 2717 files processed.
1137 830 of 2717 files processed.
1138 830 of 2717 files processed.
1140 790 of 2717 files processed.
1137 840 of 2717 files processed.
1139 840 of 2717 files processed.
1137 850 of 2717 files processed.
1138 840 of 2717 files processed.
1140 800 of 2717 files processed.
1139 850 of 2717 files processed.
1140 810 of 2717 files processed.
1137 860 of 2717 files processed.
1138 850 of 2717 files processed.
1140 820 of 2717 files processed.
1139 860 of 2717 files processed.
1138 860 of 2717 files processed.
1137 870 of 2717 files processed.
1138 870 of 2717 files processed.
1139 870 of 2717 files processed.
1140 830 of 2717 files processed.
1137 880 of 2717 files processed.
1139 880 of 2717 files processed.
1139 890 of 2717 files processed.
1138 880 of 2717 files processed.
1140 840 of 2717 files processed.
1138 890 of 2717 files processed.
1137 890 of 2717 files processed.
1140 850 of 2717 files processed.
1138 900 of 2717 files processed.
1140 860 of 2717 files processed.
1137 900 of 2717 files processed.
1139 900 of 2717 files processed.
1137 910 of 2717 files processed.
1140 870 of 2717 files processed.
1139 910 of 2717 files processed.
1138 910 of 2717 files processed.
1139 920 of 2717 files processed.
1140 880 of 2717 files processed.
1137 920 of 2717 files processed.
1139 930 of 2717 files processed.
1138 920 of 2717 files processed.
1137 930 of 2717 files processed.
1140 890 of 2717 files processed.
1138 930 of 2717 files processed.
1139 940 of 2717 files processed.
1137 940 of 2717 files processed.
1137 950 of 2717 files processed.
1139 950 of 2717 files processed.
1140 900 of 2717 files processed.
1138 940 of 2717 files processed.
1137 960 of 2717 files processed.
1139 960 of 2717 files processed.
1137 970 of 2717 files processed.
1139 970 of 2717 files processed.
1139 980 of 2717 files processed.
1137 980 of 2717 files processed.
1140 910 of 2717 files processed.
1140 920 of 2717 files processed.
1138 950 of 2717 files processed.
1137 990 of 2717 files processed.
1139 990 of 2717 files processed.
1138 960 of 2717 files processed.
1140 930 of 2717 files processed.
1139 1000 of 2717 files processed.
1138 970 of 2717 files processed.
1139 1010 of 2717 files processed.
1137 1000 of 2717 files processed.
1140 940 of 2717 files processed.
1138 980 of 2717 files processed.
1139 1020 of 2717 files processed.
1140 950 of 2717 files processed.
1139 1030 of 2717 files processed.
1140 960 of 2717 files processed.
1137 1010 of 2717 files processed.
1139 1040 of 2717 files processed.
1137 1020 of 2717 files processed.
1138 990 of 2717 files processed.
1140 970 of 2717 files processed.
1139 1050 of 2717 files processed.
1138 1000 of 2717 files processed.
1137 1030 of 2717 files processed.
1140 980 of 27171138 files processed.
1010 of 2717 files processed.
1137 1040 of 2717 files processed.
1138 1020 of 2717 files processed.
1139 1060 of 2717 files processed.
1138 1030 of 2717 files processed.
1140 990 of 2717 files processed.
1139 1070 of 2717 files processed.
1138 1040 of 2717 files processed.
1138 1050 of 2717 files processed.
1137 1050 of 2717 files processed.
1139 1080 of 2717 files processed.
1140 1000 of 2717 files processed.
1140 1010 of 2717 files processed.
1137 1060 of 2717 files processed.
1139 1090 of 2717 files processed.
1139 1100 of 2717 files processed.
1138 1060 of 2717 files processed.
1140 1020 of 2717 files processed.
1139 1110 of 2717 files processed.
1137 1070 of 2717 files processed.
1137 1080 of 2717 files processed.
1139 1120 of 2717 files processed.
1140 1030 of 2717 files processed.
1137 1090 of 2717 files processed.
1138 1070 of 2717 files processed.
1138 1080 of 2717 files processed.
1138 1090 of 2717 files processed.
1140 1040 of 2717 files processed.
1139 1130 of 2717 files processed.
1137 1100 of 2717 files processed.
1139 1140 of 2717 files processed.
1138 1100 of 2717 files processed.
1137 1110 of 2717 files processed.
1138 1110 of 2717 files processed.
1140 1050 of 2717 files processed.
1139 1150 of 2717 files processed.
1137 1120 of 2717 files processed.
1138 1120 of 2717 files processed.
1140 1060 of 2717 files processed.
1138 1130 of 2717 files processed.
1139 1160 of 2717 files processed.
1137 1130 of 2717 files processed.
1140 1070 of 2717 files processed.
1138 1140 of 2717 files processed.
1140 1080 of 2717 files processed.
1138 1150 of 2717 files processed.
1137 1140 of 2717 files processed.
1140 1090 of 2717 files processed.
1139 1170 of 2717 files processed.
1139 1180 of 2717 files processed.
1140 1100 of 2717 files processed.
1137 1150 of 2717 files processed.
1138 1160 of 2717 files processed.
1139 1190 of 2717 files processed.
1138 1170 of 2717 files processed.
1140 1110 of 27171139 1200 of files processed.2717 files processed.
1137 1160 of 2717 files processed.
1139 1210 of 2717 files processed.
1137 1170 of 2717 files processed.
1139 1220 of 2717 files processed.
1140 1120 of 2717 files processed.
1140 1130 of 2717 files processed.
1138 1180 of 2717 files processed.
1137 1180 of 2717 files processed.
1139 1230 of 2717 files processed.
1140 1140 of 2717 files processed.
1137 1190 of 2717 files processed.
1138 1190 of 2717 files processed.
1139 1240 of 2717 files processed.
1137 1200 of 2717 files processed.
1138 1200 of 2717 files processed.
1140 1150 of 2717 files processed.
1139 1250 of 2717 files processed.
1137 1210 of 2717 files processed.
1138 1210 of 2717 files processed.
1140 1160 of 2717 files processed.
1137 1220 of 2717 files processed.
1138 1220 of 2717 files processed.
1140 1170 of 2717 files processed.
1140 1180 of 2717 files processed.
1139 1260 of 2717 files processed.
1137 1230 of 2717 files processed.
1140 1190 of 2717 files processed.
1138 1230 of 2717 files processed.
1139 1270 of 2717 files processed.
1137 1240 of 2717 files processed.
1140 1200 of 2717 files processed.
1139 1280 of 2717 files processed.
1138 1240 of 2717 files processed.
1137 1250 of 2717 files processed.
1139 1290 of 2717 files processed.

```

1138 1420 of 2717 files processed.
1137 1420 of 2717 files processed.
1139 1480 of 2717 files processed.
1140 1400 of 2717 files processed.
1138 1430 of 2717 files processed.
1137 1430 of 2717 files processed.
1139 1490 of 2717 files processed.
1140 1410 of 2717 files processed.
1138 1440 of 2717 files processed.
1139 1500 of 2717 files processed.
1140 1420 of 2717 files processed.
1139 1510 of 2717 files processed.
1137 1440 of 2717 files processed.
1140 1430 of 2717 files processed.
1138 1450 of 2717 files processed.
1137 1450 of 2717 files processed.
1140 1440 of 2717 files processed.
1140 1450 of 2717 files processed.
1137 1460 of 2717 files processed.
1138 1460 of 2717 files processed.
1137 1470 of 2717 files processed.
1139 1520 of 2717 files processed.
1137 1480 of 2717 files processed.
1139 1530 of 2717 files processed.
1137 1490 of 2717 files processed.
1138 1470 of 2717 files processed.
1137 1510 of 2717 files processed.
1139 1550 of 2717 files processed.
1138 1490 of 2717 files processed.
1137 1520 of 2717 files processed.
1139 1560 of 2717 files processed.
1140 1480 of 2717 files processed.
1138 1500 of 2717 files processed.
1138 1510 of 2717 files processed.
1137 1530 of 2717 files processed.
1139 1570 of 2717 files processed.
1138 1520 of 2717 files processed.
1140 1490 of 2717 files processed.
1138 1530 of 2717 files processed.
1140 1500 of 2717 files processed.
1139 1580 of 2717 files processed.
1139 1590 of 2717 files processed.
1138 1540 of 2717 files processed.
1140 1510 of 2717 files processed.
1137 1540 of 2717 files processed.
1138 1550 of 2717 files processed.
1139 1600 of 2717 files processed.
1140 1520 of 2717 files processed.
1138 1560 of 2717 files processed.
1139 1610 of 2717 files processed.
1137 1550 of 2717 files processed.
1138 1570 of 2717 files processed.
1139 1620 of 2717 files processed.
1137 1560 of 2717 files processed.
1140 1530 of 2717 files processed.
1138 1580 of 2717 files processed.
1137 1570 of 2717 files processed.
1138 1590 of 2717 files processed.
1630 of 2717 files processed.
1140 1540 of 2717 files processed.
1140 1550 of 2717 files processed.
1137 1580 of 2717 files processed.
1139 1640 of 2717 files processed.
1138 1600 of 2717 files processed.
1137 1590 of 2717 files processed.
1140 1560 of 2717 files processed.
1139 1650 of 2717 files processed.
1139 1660 of 2717 files processed.
1137 1600 of 2717 files processed.
1140 1570 of 2717 files processed.
1138 1610 of 2717 files processed.
1138 1620 of 2717 files processed.
1138 1630 of 2717 files processed.
1139 1670 of 2717 files processed.
1137 1610 of 2717 files processed.
1139 1680 of 2717 files processed.
1138 1640 of 2717 files processed.
1140 1580 of 2717 files processed.
1137 1620 of 2717 files processed.
1140 1590 of 2717 files processed.
1138 1650 of 2717 files processed.
1139 1690 of 2717 files processed.
1137 1630 of 2717 files processed.
1140 1600 of 2717 files processed.
1139 1700 of 2717 files processed.
1138 1660 of 2717 files processed.
1140 1610 of 2717 files processed.
1137 1640 of 2717 files processed.
1138 1670 of 2717 files processed.
1138 1680 of 2717 files processed.
1139 1710 of 2717 files processed.
1140 1620 of 2717 files processed.
1137 1650 of 2717 files processed.
1140 1630 of 2717 files processed.
1139 1720 of 2717 files processed.
1138 1690 of 2717 files processed.
1140 1640 of 2717 files processed.
11381137 17001600 of of 2717 2717files processed. files processed.

```

`ments/appleidai/microsoft malware/ASM Image Features.ipynb`

Documents/appleidai/microsoft malware/ASM Image Features.ipynb

1137 2530 of 2717 files processed.
1139 2610 of 2717 files processed.
1138 2540 of 2717 files processed.
1137 2540 of 2717 files processed.
1138 2550 of1140 27172590 offiles processed.
2717 files processed.
1139 2620 of 2717 files processed.
1137 2550 of 2717 files processed.
1138 2560 of 2717 files processed.
1137 2560 of 2717 files processed.
1139 2630 of 2717 files processed.
1138 2570 of 2717 files processed.
1137 2570 of 2717 files processed.
1140 2600 of 2717 files processed.
1138 2580 of 2717 files processed.
1137 2580 of 2717 files processed.
1139 2640 of 2717 files processed.
1140 2610 of 2717 files processed.
1138 2590 of 2717 files processed.
1139 2650 of 2717 files processed.
1140 2620 of 2717 files processed.
1139 2660 of 2717 files processed.
1137 2590 of 2717 files processed.
1138 2600 of 2717 files processed.
1140 2630 of 2717 files processed.
1139 2670 of 2717 files processed.
1140 2640 of 2717 files processed.
1139 2680 of 2717 files processed.
1138 2610 of 2717 files processed.
1137 2600 of 2717 files processed.
1139 2690 of 2717 files processed.
1138 2620 of 2717 files processed.
1140 2650 of 2717 files processed.
1138 2630 of 2717 files processed.
1140 2660 of 2717 files processed.
1139 2700 of 2717 files processed.
1138 2640 of 2717 files processed.
1137 2610 of 2717 files processed.
1139 2710 of 2717 files processed.
1137 2620 of 2717 files processed.
1137 2630 of 2717 files processed.
1138 2650 of 2717 files processed.
1140 2670 of 2717 files processed.
1137 2640 of 2717 files processed.
1138 2660 of 2717 files processed.
1140 2680 of 2717 files processed.
1137 2650 of 2717 files processed.
1140 2690 of 2717 files processed.
1138 2670 of 2717 files processed.
1137 2660 of 2717 files processed.
1137 2670 of 2717 files processed.
1138 2680 of 2717 files processed.
1140 2700 of 2717 files processed.
1137 2680 of 2717 files processed.
1138 2690 of 2717 files processed.
1137 2690 of 2717 files processed.
1140 2710 of 2717 files processed.
1138 2700 of 2717 files processed.
1138 2710 of 2717 files processed.
1137 2700 of 2717 files processed.
1137 2710 of 2717 files processed.
Elapsed time: 0.00 hours.

Merging the image dataset

```
In [6]: 1 file1 = pd.read_csv('1137-image-features-asm.csv')
2 file2 = pd.read_csv('1138-image-features-asm.csv')
3 file3 = pd.read_csv('1139-image-features-asm.csv')
4 file4 = pd.read_csv('1140-image-features-asm.csv')
5
6 image_data = np.concatenate([file1,file2,file3,file4],axis=0)
```

```
In [7]: 1 file1.columns

Out[7]: Index(['filename', 'ASM_0', 'ASM_1', 'ASM_2', 'ASM_3', 'ASM_4', 'ASM_5',
              'ASM_6', 'ASM_7', 'ASM_8',
              ...,
              'ASM_990', 'ASM_991', 'ASM_992', 'ASM_993', 'ASM_994', 'ASM_995',
              'ASM_996', 'ASM_997', 'ASM_998', 'ASM_999'],
              dtype='object', length=1001)
```

```
In [8]: 1 labels = pd.read_csv('trainLabels.csv')
```

```
In [9]: 1 labels.head(5)

Out[9]:
```

	Id	Class
0	01kcPWA9K2BOxQeS5Rju	1
1	04EjldbPV5e1XroFOpiN	1
2	05EeG39MTRi6VY21DPd	1
3	05rJTUWYAKNegBk2wE8X	1
4	0AnoOZDNbPXlr2MRBSCJ	1

```
In [10]: 1 image_data = pd.DataFrame(data=image_data,columns=file1.columns)
```

```
In [11]: 1 image_data.head()

Out[11]:
```

	filename	ASM_0	ASM_1	ASM_2	ASM_3	ASM_4	ASM_5	ASM_6	ASM_7	ASM_8	...	ASM_990	ASM_991	ASM_992	ASM_993	ASM_994	ASM_995	ASM_996	ASM_997	ASM_998	ASM_999
0	JKbYXxOv2Wa8eMNIq04	46	116	101	120	116	58	48	48	52	...	10	46	116	101	120	116	58	48	48	52
1	1ZYiDKrdm9yVLpeU5cSJ	72	69	65	68	69	82	58	48	48	...	59	32	70	111	114	109	97	116	9	32
2	HI3iXz1IhS69DCANrVT	72	69	65	68	69	82	58	48	48	...	116	101	120	116	58	48	48	52	48	49
3	HSyZm4zIbXMT5R37F1wK	46	116	101	120	116	58	48	48	52	...	10	46	116	101	120	116	58	48	48	52
4	ESoHvWwbZ4JgwzIFK6uA	72	69	65	68	69	82	58	48	48	...	59	32	70	111	114	109	97	116	9	32

5 rows × 1001 columns

```
In [13]: 1 image_data_sorted = image_data.sort_values(by='filename')
2 labels_sorted = labels.sort_values(by='Id')
3
4 ## removing the column file name from the dataset
5
6 X_ = image_data.iloc[:,1:]
7 Y_ = labels.iloc[:,1:]
```

```
In [14]: 1 print('Shape of the train data:',X_.shape,' shape of labels:',Y_.shape)

Shape of the train data: (10868, 1000) shape of labels: (10868, 1)
```

2.Feature Selection using Chi2- selecting top 50%

refer: <https://github.com/dchad/malware-detection/blob/master/mmcc/feature-reduction.ipynb> (<https://github.com/dchad/malware-detection/blob/master/mmcc/feature-reduction.ipynb>)

```
In [15]: 1 from sklearn.feature_selection import SelectKBest, SelectPercentile
2 from sklearn.feature_selection import chi2
3 from sklearn.ensemble import RandomForestClassifier, ExtraTreesClassifier
4 import warnings
5 warnings.filterwarnings("ignore", category=DeprecationWarning)
6 warnings.filterwarnings("ignore", category=FutureWarning)
```

```
In [17]: 1 # find the top 50 percent variance features, from 1006 -> 503 features
2 fsp = SelectPercentile(chi2, 50)
3 X_new_50 = fsp.fit_transform(X_,Y_)
4 X_new_50.shape
```

```
Out[17]: (10868, 500)
```



```
In [18]: 1 selected_names = fsp.get_support(indices=True)
2 selected_names = selected_names + 1
3 selected_names
```

Out[18]: array([15, 21, 22, 29, 30, 32, 33, 34, 35, 41, 42, 43, 44, 48, 50, 125, 126, 135, 136, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 151, 155, 156, 157, 158, 160, 161, 162, 163, 164, 165, 167, 169, 173, 174, 179, 186, 188, 190, 198, 201, 202, 205, 215, 216, 220, 221, 222, 223, 224, 226, 227, 236, 240, 242, 243, 244, 245, 246, 247, 248, 249, 252, 253, 260, 261, 262, 263, 264, 265, 266, 267, 268, 271, 272, 282, 287, 291, 292, 293, 294, 295, 296, 297, 307, 310, 311, 312, 313, 314, 315, 317, 318, 323, 326, 327, 328, 330, 334, 337, 338, 339, 340, 341, 343, 344, 345, 346, 349, 351, 352, 353, 354, 356, 357, 358, 359, 366, 370, 371, 372, 373, 374, 375, 376, 378, 379, 380, 381, 384, 399, 400, 401, 405, 408, 409, 410, 412, 413, 414, 415, 422, 423, 424, 425, 426, 427, 428, 436, 437, 439, 440, 441, 446, 447, 448, 449, 450, 451, 452, 453, 457, 460, 461, 464, 465, 466, 467, 477, 478, 479, 480, 481, 482, 538, 539, 555, 556, 557, 558, 559, 560, 561, 563, 564, 567, 568, 572, 573, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 597, 598, 600, 601, 602, 603, 606, 607, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 663, 664, 665, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 691, 692, 693, 694, 695, 696, 697, 701, 702, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 717, 718, 719, 720, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 744, 746, 747, 748, 749, 751, 752, 753, 754, 755, 756, 758, 759, 760, 761, 762, 763, 764, 765, 774, 775, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 788, 790, 791, 792, 793, 794, 795, 797, 798, 799, 800, 801, 802, 813, 814, 815, 816, 818, 819, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 841, 843, 844, 845, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 865, 866, 867, 868, 869, 870, 873, 874, 875, 876, 877, 878, 879, 882, 885, 887, 888, 889, 890, 892, 893, 894, 895, 896, 898, 900, 901, 902, 903, 904, 905, 906, 916, 917, 918, 919, 920, 923, 924, 930, 931, 932, 934, 935, 936, 937, 938, 943, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000])

```
In [21]: 1 data_trimmed = image_data_sorted.iloc[:,selected_names]
2 data_fnames = pd.DataFrame(image_data_sorted['filename'])
3 data_50 = data_fnames.join(data_trimmed)
4 data_50.head()
```

Out[21]:

	filename	ASM_14	ASM_20	ASM_21	ASM_28	ASM_29	ASM_31	ASM_32	ASM_33	ASM_34	...	ASM_990	ASM_991	ASM_992	ASM_993	ASM_994	ASM_995	ASM_996	ASM_997	ASM_998	ASM_999
3677	01IsoISMh5gxyDYT4CB	9	32	32	13	10	116	101	120	116	...	10	46	116	101	120	116	58	48	48	52
5749	01SuzwMJEIXsK7A8dQbl	48	9	9	68	69	58	48	48	52	...	116	101	120	116	58	48	48	52	48	49
3852	01azqd4InC7m9JpocGv5	48	9	9	68	69	58	48	48	52	...	116	101	120	116	58	48	48	52	48	49
9748	01jsnpXSAlgw6aPeDxrU	48	9	9	68	69	58	48	48	52	...	116	101	120	116	58	48	48	52	48	49
6133	01kcPWA9K2BOxQeS5Rju	48	9	9	68	69	58	49	48	48	...	71	77	69	78	84	32	72	69	65	68

5 rows × 501 columns

```
In [23]: 1 print('Shape of the data after data reduction with 50% variance features:',data_50.shape)

Shape of the data after data reduction with 50% variance features: (10868, 501)
```

```
In [34]: 1 ### save the file after performing feature selection
2 data_50.rename(columns={'filename':'ID'},inplace=True)
3 data_50.to_csv('asm_50%_variance_data.csv')
```

3. Implemation of different Feature Combination with ASM extracted image

3.1 Asm-uni grams + ASM-Extracted image features

```
In [127]: 1 labels =pd.read_csv('trainLabels.csv')
```

```
In [128]: 1 dfasm=pd.read_csv("asmoutputfile.csv")
2 labels.columns = ['ID', 'Class']
3 result_asm = pd.merge(dfasm, labels,on='ID', how='left')
4 result_asm.head()
```

Out[128]:

	ID	HEADER:	.text:	.Pav:	.ldata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	edx	esi	eax	ebx	ecx	edi	ebp	esp	ebp	Class
0	01kcPWA9K2BOxQeS5Rju	19	744	0	127	57	0	323	0	3	...	18	66	15	43	83	0	17	48	29	1
1	1E93CpP60RHFNIT5Qfn	17	838	0	103	49	0	0	0	3	...	18	29	48	82	12	0	14	0	20	1
2	3ekVowZajZHbTnBcsDFX	17	427	0	50	43	0	145	0	3	...	13	42	10	67	14	0	11	0	9	1
3	3X2nY7QaPBIWDrAZqJe	17	227	0	43	19	0	0	0	3	...	6	8	14	7	2	0	8	0	6	1
4	46OZdsSKDCFV8h7XWxf	17	402	0	59	170	0	0	0	3	...	12	9	18	29	5	0	11	0	11	1

5 rows × 53 columns

```
In [129]: 1 data_50 = pd.read_csv('asm_50%_variance_data.csv')
```

```
In [130]: 1 data_50.drop(columns=['Unnamed: 0'],inplace=True)
2 data_50.columns
```

Out[130]: Index(['ID', 'ASM_14', 'ASM_20', 'ASM_21', 'ASM_28', 'ASM_29', 'ASM_31', 'ASM_32', 'ASM_33', 'ASM_34', ..., 'ASM_990', 'ASM_991', 'ASM_992', 'ASM_993', 'ASM_994', 'ASM_995', 'ASM_996', 'ASM_997', 'ASM_998', 'ASM_999'], dtype='object', length=501)

```
In [131]: 1 X_merged = pd.merge(result_asm, data_50, on='ID', how='inner')
2 X_merged.head()
```

Out[131]:

	ID	HEADER:	.text:	.Pav:	.ldata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	ASM_990	ASM_991	ASM_992	ASM_993	ASM_994	ASM_995	ASM_996	ASM_997	ASM_998	ASM_999
0	01kcPWA9K2BOxQeS5Rju	19	744	0	127	57	0	323	0	3	...	71	77	69	78	84	32	72	69	65	68
1	1E93CpP60RHFNIT5Qfn	17	838	0	103	49	0	0	0	3	...	32	32	32	32	58	9	80	111	114	116
2	3ekVowZajZHbTnBcsDFX	17	427	0	50	43	0	145	0	3	...	59	32	70	111	114	109	97	116	9	32
3	3X2nY7QaPBIWDrAZqJe	17	227	0	43	19	0	0	0	3	...	32	32	32	32	58	9	80	111	114	116
4	46OZdsSKDCFV8h7XWxf	17	402	0	59	170	0	0	0	3	...	32	32	32	32	58	9	80	111	114	116

5 rows × 553 columns

```
In [132]: 1 Y_merged = X_merged['Class']
2 X_merged = X_merged.drop(['ID','rtn','.BSS:','.CODE','Class'], axis=1)
```

```
In [133]: 1 ### train test split
2 X_train, X_test, y_train, y_test = train_test_split(X_merged, Y_merged,stratify=Y_merged,test_size=0.20)
3 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,stratify=y_train,test_size=0.20)
```

```
In [134]: 1 X_train = X_train.apply(pd.to_numeric, errors='coerce')
2 X_test = X_test.apply(pd.to_numeric, errors='coerce')
3 X_cv = X_cv.apply(pd.to_numeric, errors='coerce')
```

- since random forest and XGB performed better on bytes/asm unigrams we will use both

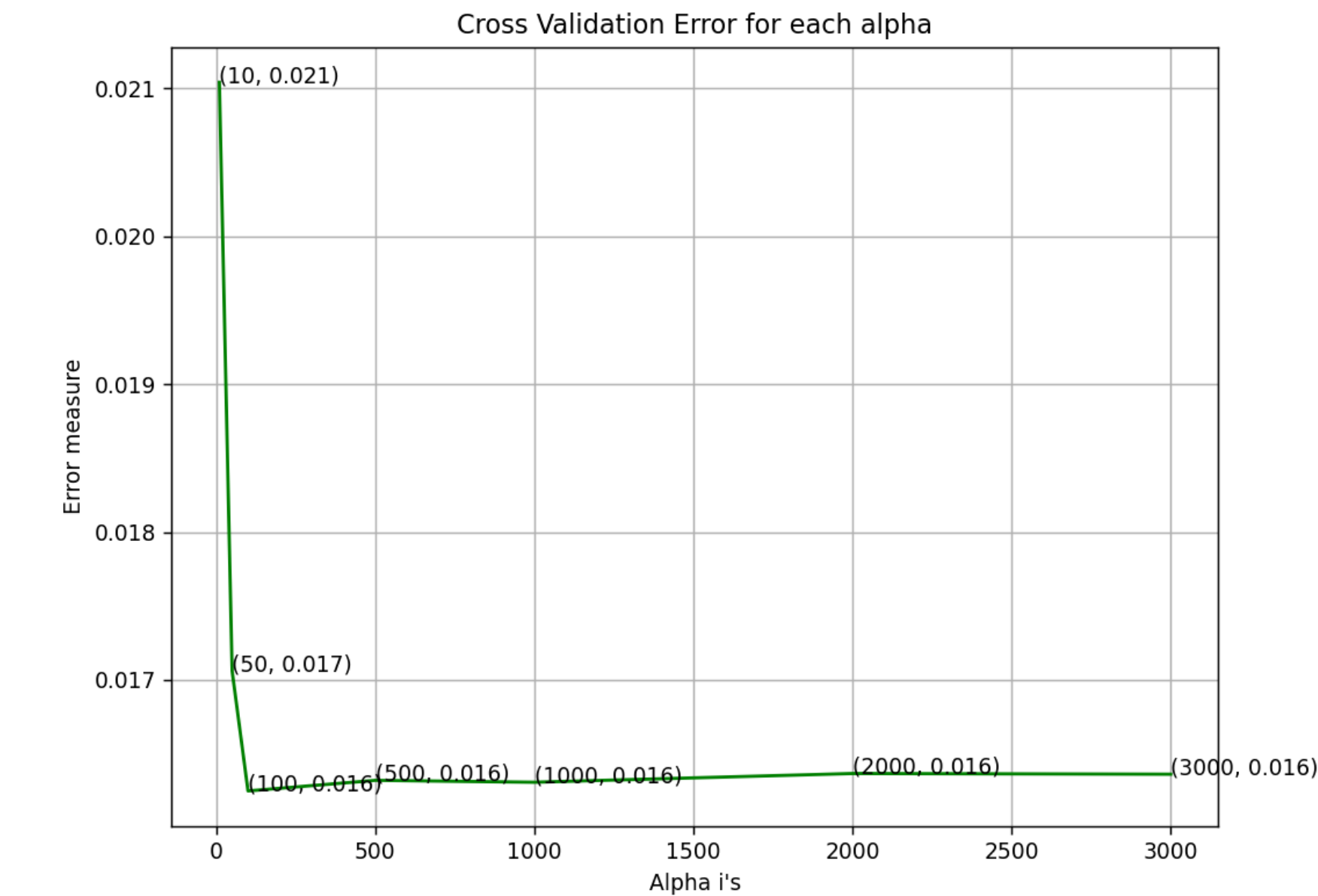
3.1.1 Random Forest Classifier with Asm-uni grams + ASM-Extracted image features

In [52]:

```
1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # predict_proba(X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
19 # -----
20
21 alpha=[10,50,100,500,1000,2000,3000]
22 cv_log_error_array=[]
23 train_log_error_array=[]
24 for i in alpha:
25     r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
26     r_cfl.fit(X_train,y_train)
27     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
28     sig_clf.fit(X_train, y_train)
29     predict_y = sig_clf.predict_proba(X_cv)
30     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
31     print('Parameter tuning for alpha:',i,' Loss:',log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
32 for i in range(len(cv_log_error_array)):
33     print ('log_loss for c = ',alpha[i],'is',cv_log_error_array[i])
34
35
36 best_alpha = np.argmin(cv_log_error_array)
37
38 fig, ax = plt.subplots()
39 ax.plot(alpha, cv_log_error_array,c='g')
40 for i, txt in enumerate(np.round(cv_log_error_array,3)):
41     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
42 plt.grid()
43 plt.title("Cross Validation Error for each alpha")
44 plt.xlabel("Alpha i's")
45 plt.ylabel("Error measure")
46 plt.show()
47
48 r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
49 r_cfl.fit(X_train,y_train)
50 sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
51 sig_clf.fit(X_train, y_train)
52
53
54 predict_y = sig_clf.predict_proba(X_train)
55 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
56 predict_y = sig_clf.predict_proba(X_cv)
57 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
58 predict_y = sig_clf.predict_proba(X_test)
59 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
60
```

Parameter tuning for alpha: 10 Loss: 0.021040358325913156
Parameter tuning for alpha: 50 Loss: 0.01706201515911518
Parameter tuning for alpha: 100 Loss: 0.01624964757759834
Parameter tuning for alpha: 500 Loss: 0.016321154903590304
Parameter tuning for alpha: 1000 Loss: 0.01630844779802613
Parameter tuning for alpha: 2000 Loss: 0.016367377506599643
Parameter tuning for alpha: 3000 Loss: 0.016362740697231582
log_loss for c = 10 is 0.021040358325913156
log_loss for c = 50 is 0.01706201515911518
log_loss for c = 100 is 0.01624964757759834
log_loss for c = 500 is 0.016321154903590304
log_loss for c = 1000 is 0.01630844779802613
log_loss for c = 2000 is 0.016367377506599643
log_loss for c = 3000 is 0.016362740697231582

<IPython.core.display.Javascript object>



For values of best alpha = 100 The train log loss is: 0.012043317011966856
For values of best alpha = 100 The cross validation log loss is: 0.01624964757759834
For values of best alpha = 100 The test log loss is: 0.02383102390750985

```
In [54]: 1 # predict_y is test prediction
        2 predicted_y = np.argmax(predict_y, axis=1)
        3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points 0.45998160073597055

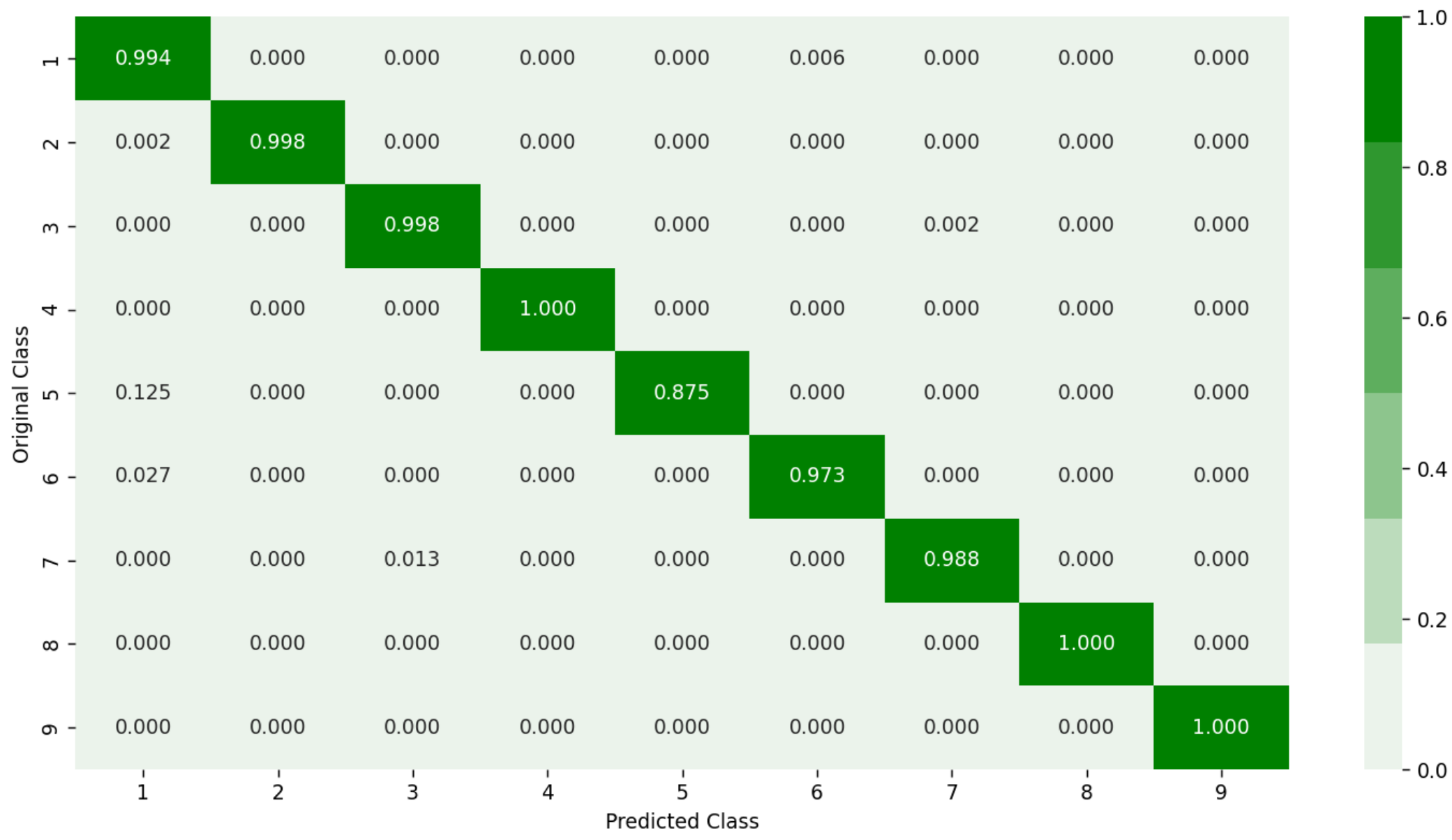
----- Confusion matrix -----
<IPython.core.display.Javascript object>



----- Precision matrix -----
<IPython.core.display.Javascript object>



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



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

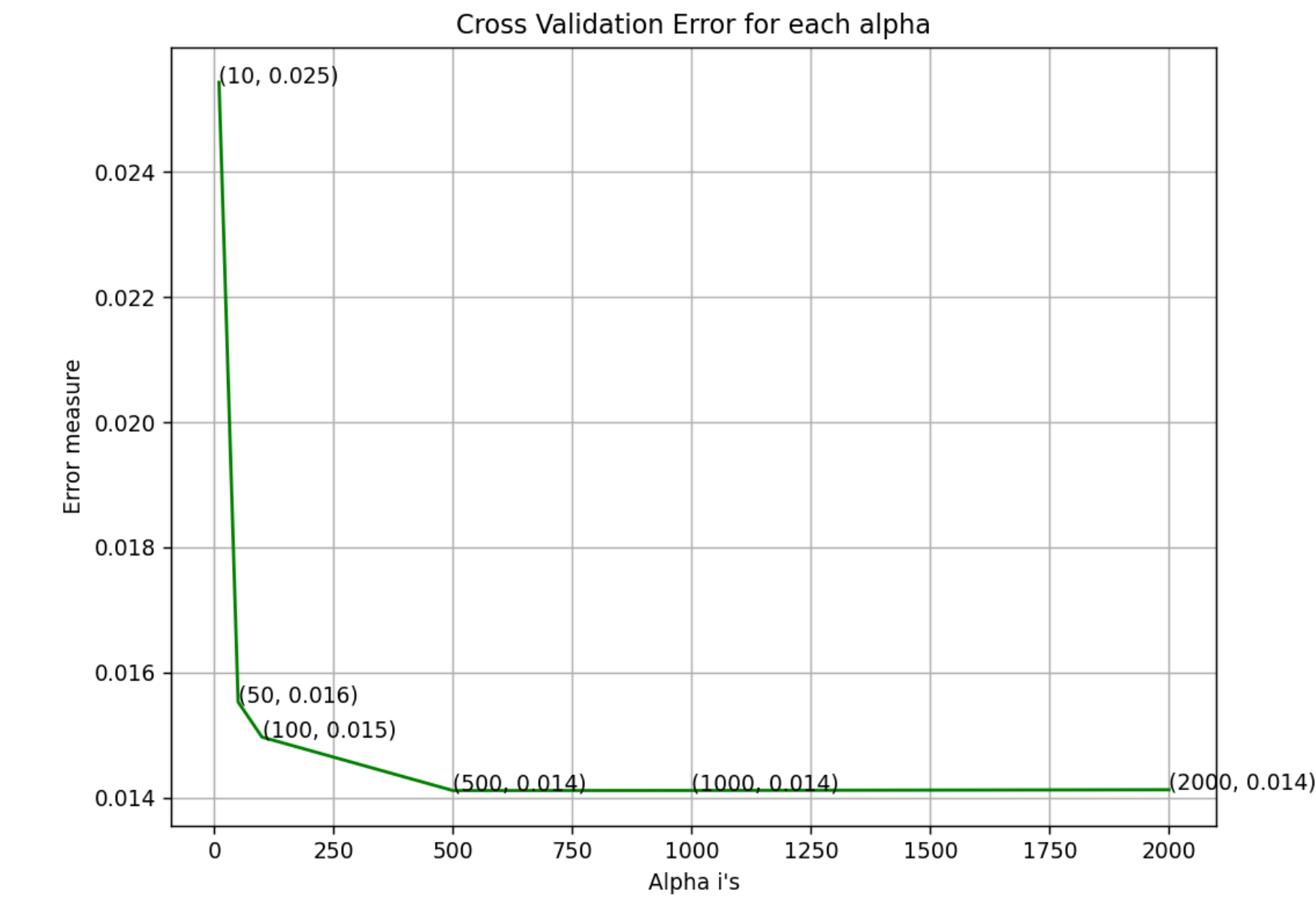
3.1.2 XGB Classifier with Asm-uni grams + ASM-Extracted image features

In [79]:

```
1 %%time
2 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
3
4 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
5 # -----
6 # default paramters
7 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
8 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1,
9 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1,
10 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
11
12 # some of methods of RandomForestRegressor()
13 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
14 # get_params([deep]) Get parameters for this estimator.
15 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
16 # get_score(importance_type='weight') -> get the feature importance
17 # -----
18 # video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-decision-trees-2/
19 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
20 # -----
21
22 alpha=[10,50,100,500,1000,2000]
23 cv_log_error_array=[]
24 for i in alpha:
25     x_cfl=XGBClassifier(n_estimators=i,nthread=-1,silent=True,verbosity=0)
26     x_cfl.fit(X_train,y_train)
27     sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
28     sig_clf.fit(X_train, y_train)
29     predict_y = sig_clf.predict_proba(X_cv)
30     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_, eps=1e-15))
31     print('Parameter tuning for alpha:',i,' Loss:',log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
32
33
34 # for i in range(len(cv_log_error_array)):
35 #     print ('Log_Loss for c = ',alpha[i],'is',cv_log_error_array[i])
36
37
38 best_alpha = np.argmin(cv_log_error_array)
39
40 fig, ax = plt.subplots()
41 ax.plot(alpha, cv_log_error_array,c='g')
42 for i, txt in enumerate(np.round(cv_log_error_array,3)):
43     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
44 plt.grid()
45 plt.title("Cross Validation Error for each alpha")
46 plt.xlabel("Alpha i's")
47 plt.ylabel("Error measure")
48 plt.show()
49
50 x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1,silent=True,verbosity=0)
51 x_cfl.fit(X_train,y_train)
52 sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
53 sig_clf.fit(X_train, y_train)
54
55 predict_y = sig_clf.predict_proba(X_train)
56 print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
57 predict_y = sig_clf.predict_proba(X_cv)
58 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
59 predict_y = sig_clf.predict_proba(X_test)
60 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
61
```

Parameter tuning for alpha: 10 Loss: 0.025430225353088646
Parameter tuning for alpha: 50 Loss: 0.01553217557088968
Parameter tuning for alpha: 100 Loss: 0.014970418779213806
Parameter tuning for alpha: 500 Loss: 0.014116494279968701
Parameter tuning for alpha: 1000 Loss: 0.014117920024732094
Parameter tuning for alpha: 2000 Loss: 0.014129280344415004
log_loss for c = 10 is 0.025430225353088646
log_loss for c = 50 is 0.01553217557088968
log_loss for c = 100 is 0.014970418779213806
log_loss for c = 500 is 0.014116494279968701
log_loss for c = 1000 is 0.014117920024732094
log_loss for c = 2000 is 0.014129280344415004

<IPython.core.display.Javascript object>



For values of best alpha = 500 The train log loss is: 0.010945769938987372
For values of best alpha = 500 The cross validation log loss is: 0.014116494279968701
For values of best alpha = 500 The test log loss is: 0.023801783541974996

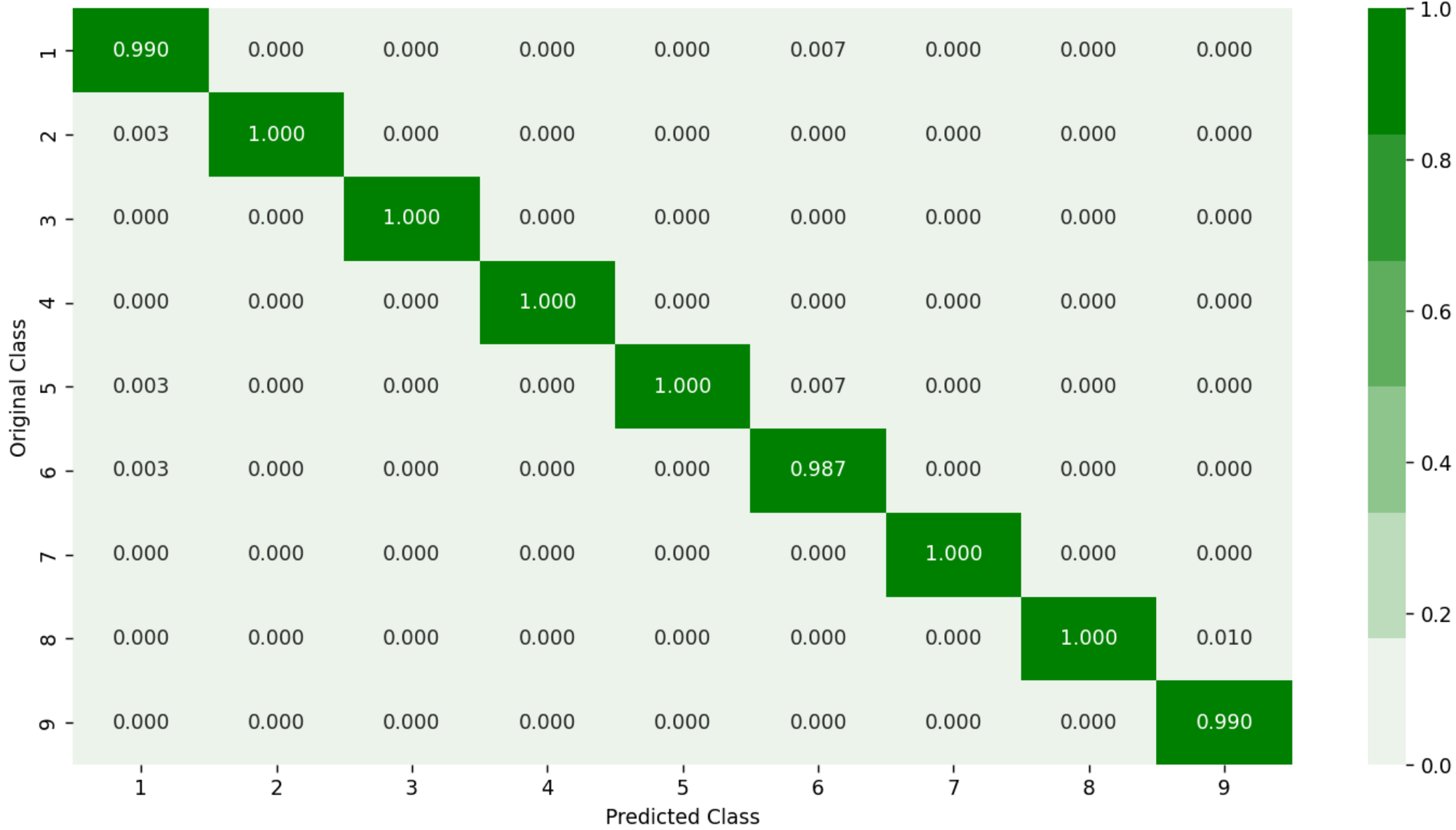
```
In [80]: 1 # predict_y is test prediction
        2 predicted_y = np.argmax(predict_y, axis=1)
        3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points 0.3219871205151794

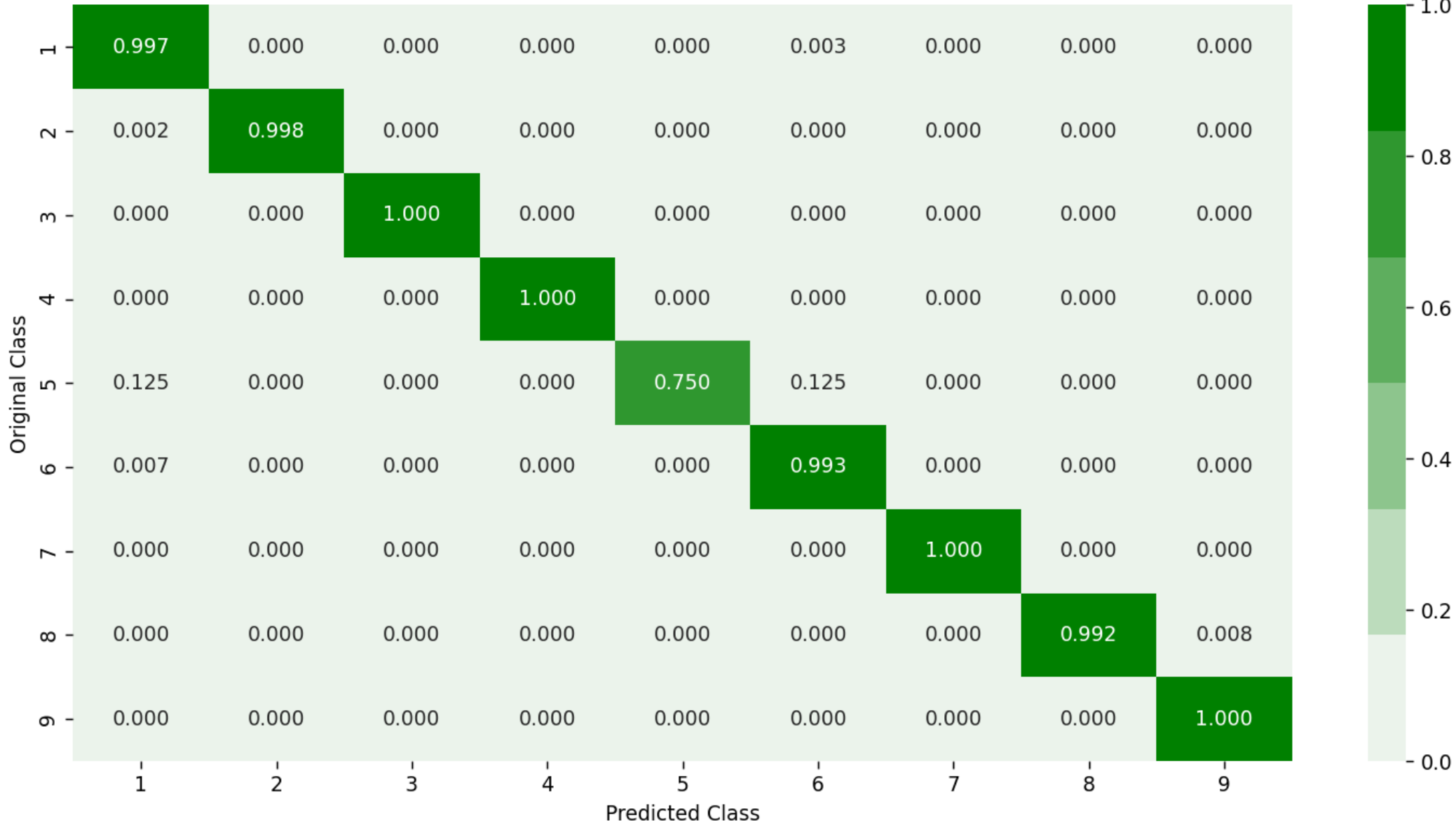
----- Confusion matrix -----
<IPython.core.display.Javascript object>



----- Precision matrix -----
<IPython.core.display.Javascript object>



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



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

3.1.3 XGB Classifier with Asm-uni grams + ASM-Extracted image features (with best hyper paramters)

<https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/> / <https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/>

In [16]:

```
1 %%time
2 x_cfl=XGBClassifier(verbosity=0,silent=True)
3 prams={
4     'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
5     'n_estimators':[100,200,500],
6     'max_depth':[3,5,10],
7     'colsample_bytree':[0.1,0.3,0.5,1],
8     'subsample':[0.1,0.3,0.5,1]
9 }
10 random_cfl1=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,cv=4)
11 random_cfl1.fit(X_train,y_train)
```

Fitting 4 folds for each of 10 candidates, totalling 40 fits

[CV 1/4; 1/10] START colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3

[CV 1/4; 1/10] END colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3, score=0.989 total time= 18.0s

[CV 2/4; 1/10] START colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3

[CV 2/4; 1/10] END colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3, score=0.992 total time= 32.2s

[CV 3/4; 1/10] START colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3

[CV 3/4; 1/10] END colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3, score=0.987 total time= 32.5s

[CV 4/4; 1/10] START colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3

[CV 4/4; 1/10] END colsample_bytree=0.3, learning_rate=0.03, max_depth=5, n_estimators=100, subsample=0.3, score=0.994 total time= 33.0s

[CV 1/4; 2/10] START colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3

[CV 1/4; 2/10] END colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3, score=0.997 total time= 43.6s

[CV 2/4; 2/10] START colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3

[CV 2/4; 2/10] END colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3, score=0.997 total time= 43.0s

[CV 3/4; 2/10] START colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3

[CV 3/4; 2/10] END colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3, score=0.995 total time= 42.8s

[CV 4/4; 2/10] START colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3

[CV 4/4; 2/10] END colsample_bytree=0.1, learning_rate=0.1, max_depth=10, n_estimators=500, subsample=0.3, score=0.997 total time= 43.4s

[CV 1/4; 3/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 1/4; 3/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.995 total time= 40.0s

[CV 2/4; 3/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 2/4; 3/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.992 total time= 39.3s

[CV 3/4; 3/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 3/4; 3/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.990 total time= 40.5s

[CV 4/4; 3/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 4/4; 3/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.996 total time= 41.0s

[CV 1/4; 4/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 1/4; 4/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.995 total time= 27.7s

[CV 2/4; 4/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 2/4; 4/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.993 total time= 27.0s

[CV 3/4; 4/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 3/4; 4/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.991 total time= 27.5s

[CV 4/4; 4/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1

[CV 4/4; 4/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=10, n_estimators=200, subsample=0.1, score=0.995 total time= 27.7s

[CV 1/4; 5/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3

[CV 1/4; 5/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3, score=0.995 total time= 1.6min

[CV 2/4; 5/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3

[CV 2/4; 5/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3, score=0.995 total time= 1.6min

[CV 3/4; 5/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3

[CV 3/4; 5/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3, score=0.993 total time= 1.6min

[CV 4/4; 5/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3

[CV 4/4; 5/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.3, score=0.997 total time= 1.6min

[CV 1/4; 6/10] START colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3

[CV 1/4; 6/10] END colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3, score=0.997 total time= 1.7min

[CV 2/4; 6/10] START colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3

[CV 2/4; 6/10] END colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3, score=0.995 total time= 1.7min

[CV 3/4; 6/10] START colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3

[CV 3/4; 6/10] END colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3, score=0.993 total time= 1.7min

[CV 4/4; 6/10] START colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3

[CV 4/4; 6/10] END colsample_bytree=1, learning_rate=0.2, max_depth=3, n_estimators=200, subsample=0.3, score=0.996 total time= 1.8min

[CV 1/4; 7/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3

[CV 1/4; 7/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3, score=0.995 total time= 24.8s

[CV 2/4; 7/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3

[CV 2/4; 7/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3, score=0.995 total time= 24.3s

[CV 3/4; 7/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3

[CV 3/4; 7/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3, score=0.993 total time= 24.4s

[CV 4/4; 7/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3

[CV 4/4; 7/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=100, subsample=0.3, score=0.996 total time= 24.3s

[CV 1/4; 8/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1

[CV 1/4; 8/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1, score=0.978 total time= 15.3s

[CV 2/4; 8/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1

[CV 2/4; 8/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1, score=0.979 total time= 15.5s

[CV 3/4; 8/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1

[CV 3/4; 8/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1, score=0.978 total time= 15.3s

[CV 4/4; 8/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1

[CV 4/4; 8/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=0.1, score=0.984 total time= 15.0s

[CV 1/4; 9/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1

[CV 1/4; 9/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1, score=0.995 total time= 1.2min

[CV 2/4; 9/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1

[CV 2/4; 9/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1, score=0.994 total time= 1.1min

[CV 3/4; 9/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1

[CV 3/4; 9/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1, score=0.992 total time= 1.2min

[CV 4/4; 9/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1

[CV 4/4; 9/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=10, n_estimators=500, subsample=0.1, score=0.997 total time= 1.2min

[CV 1/4; 10/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3

[CV 1/4; 10/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3, score=0.995 total time= 5.4min

[CV 2/4; 10/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3

[CV 2/4; 10/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3, score=0.996 total time= 5.4min

[CV 3/4; 10/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3

[CV 3/4; 10/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3, score=0.993 total time= 5.3min

[CV 4/4; 10/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3

[CV 4/4; 10/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=500, subsample=0.3, score=0.996 total time= 5.4min

CPU times: user 1h 43min 44s, sys: 5.84 s, total: 1h 43min 50s

Wall time: 52min 11s

```
Out[16]: RandomizedSearchCV(cv=4,
      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=nan,
                              monotone_constraints=None,
                              n_estimators=100,...
                              num_parallel_tree=None,
                              random_state=None, reg_alpha=None,
                              reg_lambda=None,
                              scale_pos_weight=None, silent=True,
                              subsample=None, tree_method=None,
                              validate_parameters=None,
                              verbosity=0),
      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],
                          'subsample': [0.1, 0.3, 0.5, 1]},
      verbose=10)
```

In [17]:

```
1 print (random_cfl1.best_params_)

{'subsample': 0.3, 'n_estimators': 500, 'max_depth': 10, 'learning_rate': 0.1, 'colsample_bytree': 0.1}
```

In [135]:

```
1 %%time
2 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
3
4 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
5 # -----
6 # default parameters
7 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
8 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1,
9 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1,
10 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
11
12 # some of methods of RandomForestRegressor()
13 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
14 # get_params([deep]) Get parameters for this estimator.
15 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
16 # get_score(importance_type='weight') -> get the feature importance
17 # -----
18 # video Link2: https://www.appliedaiacourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
19 # -----
20 x_cfl=XGBClassifier(n_estimators=500, subsample=0.3,learning_rate=0.1, colsample_bytree=0.1, max_depth=10)
21 x_cfl.fit(X_train,y_train)
22 c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
23 c_cfl.fit(X_train,y_train)
24
25 predict_y = c_cfl.predict_proba(X_train)
26 print ('train loss',log_loss(y_train, predict_y))
27 predict_y = c_cfl.predict_proba(X_cv)
28
29 print ('cv loss',log_loss(y_cv, predict_y))
30 predict_y = c_cfl.predict_proba(X_test)
31 print ('test loss',log_loss(y_test, predict_y))
```

train loss 0.011000827447175468

cv loss 0.021335735778038614

test loss 0.025312012599598657

CPU times: user 5min 22s, sys: 530 ms, total: 5min 22s

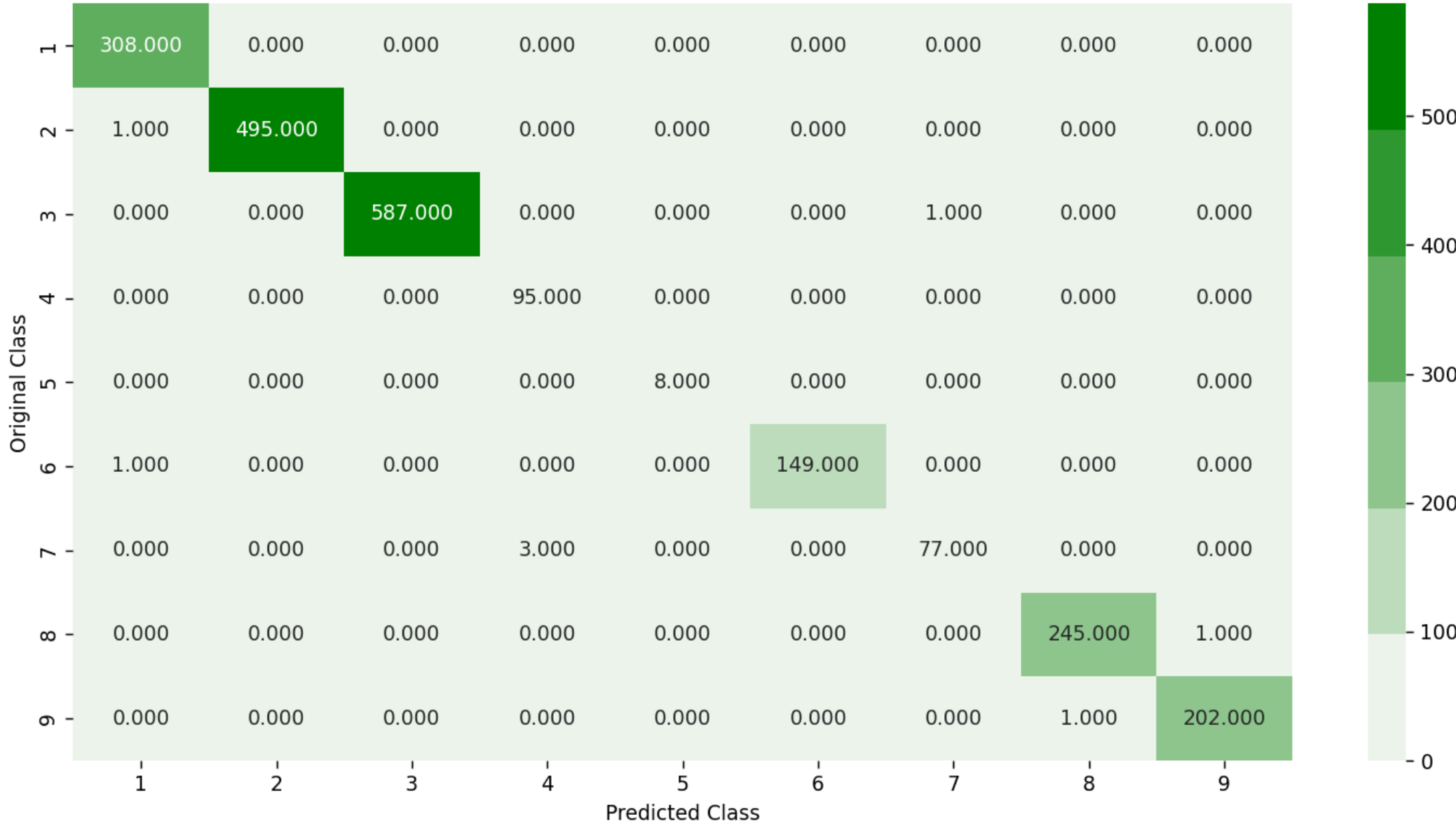
Wall time: 2min 43s

```
In [136]: 1 # predict_y is test prediction
          2 predicted_y = np.argmax(predict_y, axis=1)
          3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points 0.24798528058877645

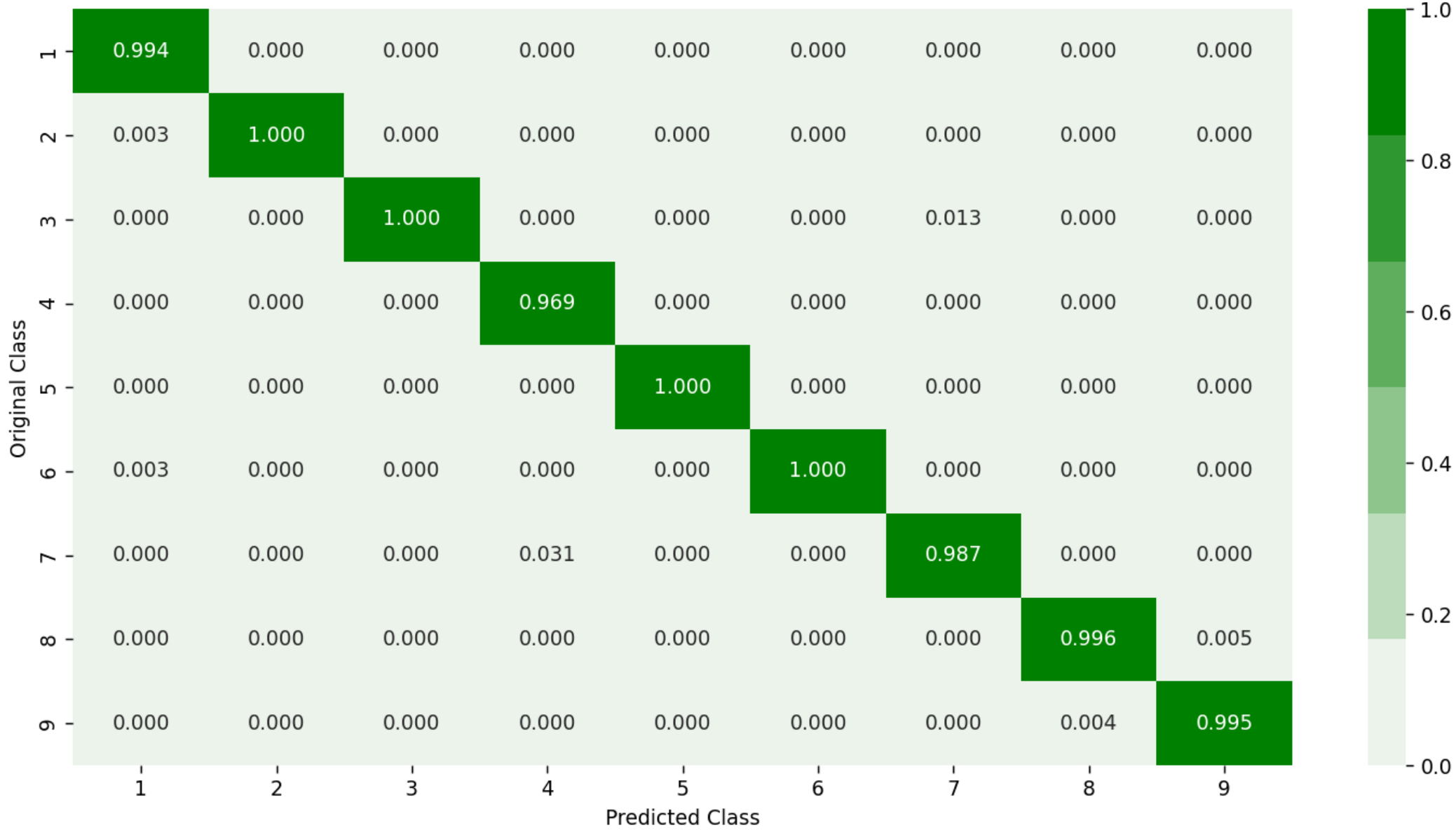
----- Confusion matrix -----

<IPython.core.display.Javascript object>



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

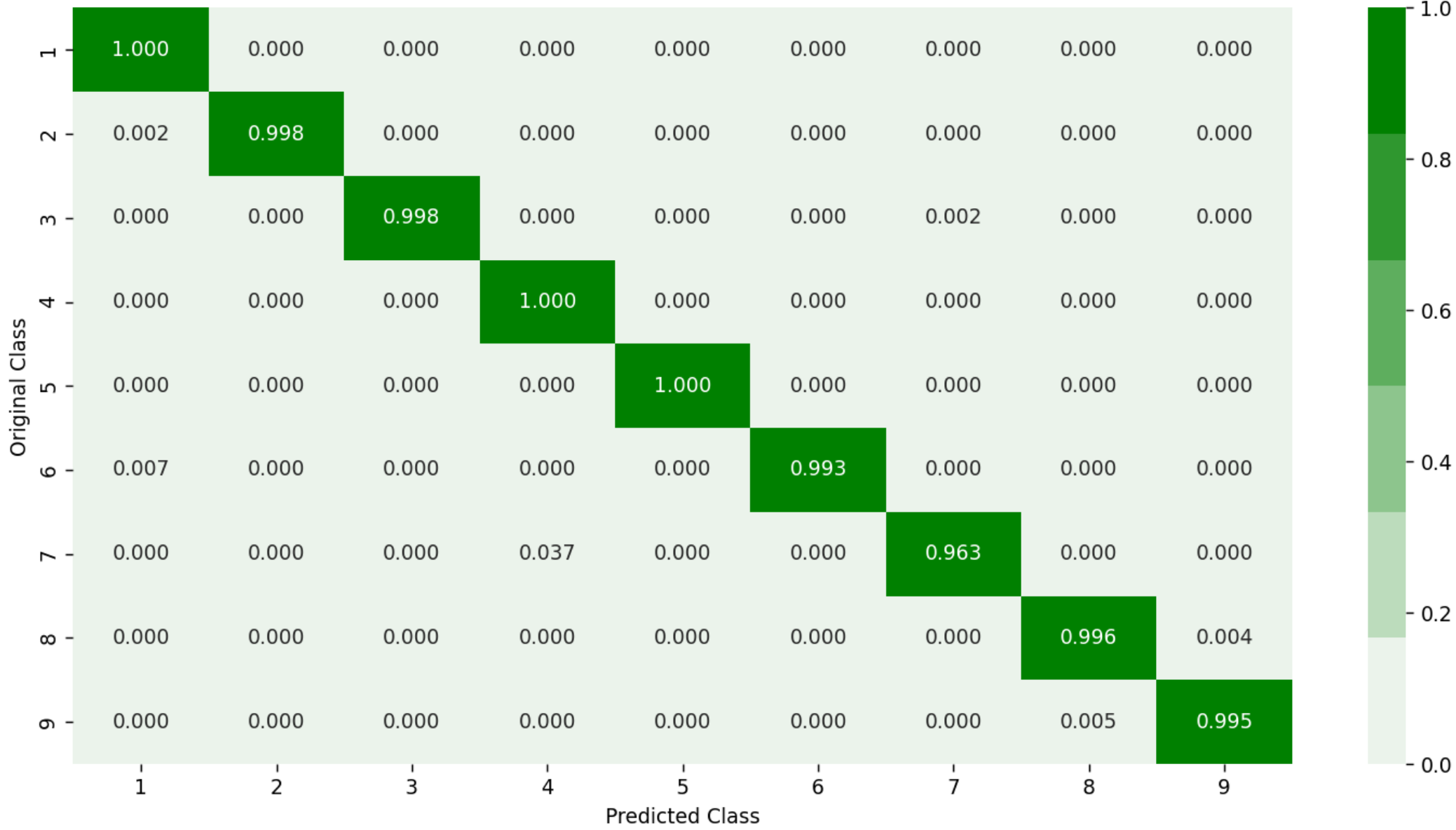
<IPython.core.display.Javascript object>



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

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

<IPython.core.display.Javascript object>



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

3.2 Random Forest Classifier , Asm-uni grams , ASM-Extracted image features , Bytes uni-grams

In [3]:

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

Out[3]:

	ID	0	1	2	3	4	5	6	7	8	...	f7	f8	f9	fa	fb	fc	fd	fe	ff	??
0	01azqd4InC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...	2804	3687	3101	3211	3097	2758	3099	2759	5753	1824
1	01IsoISMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...	451	6536	439	281	302	7639	518	17001	54902	8588

2 rows × 258 columns

In [4]:

```
1 byte_features_with_size = pd.read_csv("result_with_size.csv")
2 byte_features_with_size = byte_features_with_size.drop(columns='Unnamed: 0')
3 byte_features_with_size.head(5)
```

Out[4]:

	ID	0	1	2	3	4	5	6	7	8	...	f9	fa	fb	fc	fd	fe	ff	??	size	Class
0	01azqd4InC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	2965	...	3101	3211	3097	2758	3099	2759	5753	1824	4.148438	9
1	01IsoISMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	9291	...	439	281	302	7639	518	17001	54902	8588	5.425781	2
2	01jsnpXSAIgw6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	9107	...	2242	2885	2863	2471	2786	2680	49144	468	3.808594	9
3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	1078	...	485	462	516	1133	471	761	7998	13940	0.562500	1
4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	422	...	350	209	239	653	221	242	2199	9008	0.363281	8

5 rows × 260 columns

In [5]:

```
1 # https://stackoverflow.com/a/29651514
2 def normalize(df):
3     result1 = df.copy()
4     for feature_name in df.columns:
5         if (str(feature_name) != str('ID') and str(feature_name)!=str('Class')):
6             max_value = df[feature_name].max()
7             min_value = df[feature_name].min()
8             result1[feature_name] = (df[feature_name] - min_value) / (max_value - min_value)
9     return result1
10 result = normalize(byte_features_with_size)
11
```

In [7]:

```
1 labels = pd.read_csv('trainLabels.csv')
2 data_50 = pd.read_csv('asm_50%_variance_data.csv')
```

In [8]:

```
1 result_asm = pd.read_csv("asmoutputfile.csv")
2
3 result_asm = normalize(result_asm)
4
5 labels.columns = ['ID', 'Class']
6 result_asm = pd.merge(result_asm, labels,on='ID', how='inner')
7 result_asm.head()
```

Out[8]:

	ID	HEADER:	.text:	.Pav:	.ldata:	.data:	.bss:	.rdata:	.edata:	.rsrc:	...	edx	esi	eax	ebx	ecx	edi	ebp	esp	eip	Class	
0	01kcPWA9K2BOxQeS5Rju		0.107345	0.001092	0.0	0.000761	0.000023	0.0	0.000084	0.0	0.000072	...	0.000343	0.000746	0.000301	0.000360	0.001057	0.0	0.030797	0.001468	0.003173	1
1	1E93CpP60RHFNITSQfn		0.096045	0.001230	0.0	0.000617	0.000019	0.0	0.000000	0.0	0.000072	...	0.000343	0.000328	0.000965	0.000686	0.000153	0.0	0.025362	0.000000	0.002188	1
2	3ekVow2ajZhTnBcsDX		0.096045	0.000627	0.0	0.000300	0.000017	0.0	0.000038	0.0	0.000072	...	0.000248	0.000475	0.000201	0.000560	0.000178	0.0	0.019928	0.000000	0.000985	1
3	3X2nY7QaPBIWDrAZqJe		0.096045	0.000333	0.0	0.000258	0.000008	0.0	0.000000	0.0	0.000072	...	0.000114	0.000090	0.000281	0.000059	0.000025	0.0	0.014493	0.000000	0.000657	1
4	46OZzdsSKDCFV8h7XWxf		0.096045	0.000590	0.0	0.000353	0.000068	0.0	0.000000	0.0	0.000072	...	0.000229	0.000102	0.000362	0.000243	0.000064	0.0	0.019928	0.000000	0.001204	1

5 rows × 53 columns

In [9]:

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

Out[9]:

	ID	0	1	2	3	4	5	6	7	8	...	f9	fa	fb	fc	fd	fe	ff	??	size	Class
0	01azqd4InC7m9JpocGv5	0.262806	0.005498	0.001567	0.002067	0.002048	0.001835	0.002058	0.002946	0.002638	...	0.013560	0.013107	0.013634	0.031724	0.014549	0.014348	0.007843	0.000129	0.092219	9
1	01IsoISMh5gxyDYTI4CB	0.017358	0.011737	0.004033	0.003876	0.005303	0.003873	0.004747	0.006984	0.008267	...	0.001920	0.001147	0.001329	0.087867	0.002432	0.088411	0.074851	0.000606	0.121236	2
2	01jsnpXSAIgw6aPeDxrU	0.040827	0.013434	0.001429	0.001315	0.005464	0.005280	0.005078	0.002155	0.008104	...	0.009804	0.011777	0.012604	0.028423	0.013080	0.013937	0.067001	0.000033	0.084499	9
3	01kcPWA9K2BOxQeS5Rju	0.009209	0.001708	0.000404	0.000441	0.000770	0.000354	0.000310	0.000481	0.000959	...	0.002121	0.001886	0.002272	0.013032	0.002211	0.003957	0.010904	0.000984	0.010759	1
4	01SuzwMJEIXsK7A8dQbl	0.008629	0.001000	0.000168	0.000234	0.000342	0.000232	0.000148	0.000229	0.000376	...	0.001530	0.000853	0.001052	0.007511	0.001038	0.001258	0.002998	0.000636	0.006233	8

5 rows × 260 columns

In [10]:

```
1 print('total features (asm unigram + bytes-unigram + asm-img):',result.shape[1]+result_asm.shape[1]+data_50.shape[1])
```

total features (asm unigram + bytes-unigram + asm-img): 815

In [11]:

```
1 X_merged = pd.merge(result_asm,result.drop(columns=['size','Class']),on='ID',how='inner')
2 X_merged = pd.merge(X_merged, data_50, on='ID', how='inner')
```

In [12]:

```
1 X_merged.shape
```

Out[12]: (10868, 811)

In [13]:

```
1 Y_merged = X_merged['Class']
2 X_merged = X_merged.drop(['ID','rtn','.BSS:','.CODE','Class'], axis=1)
3 print('shape of x',X_merged.shape,' y:',Y_merged.shape)
```

shape of x (10868, 806) y: (10868,)

In [14]:

```
1 ### train test split
2 X_train, X_test, y_train, y_test = train_test_split(X_merged, Y_merged,stratify=Y_merged,test_size=0.20)
3 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,stratify=y_train,test_size=0.20)
```

In [15]:

```
1 X_train = X_train.astype(float)
2 X_test = X_test.astype(float)
3 X_cv = X_cv.astype(float)
```

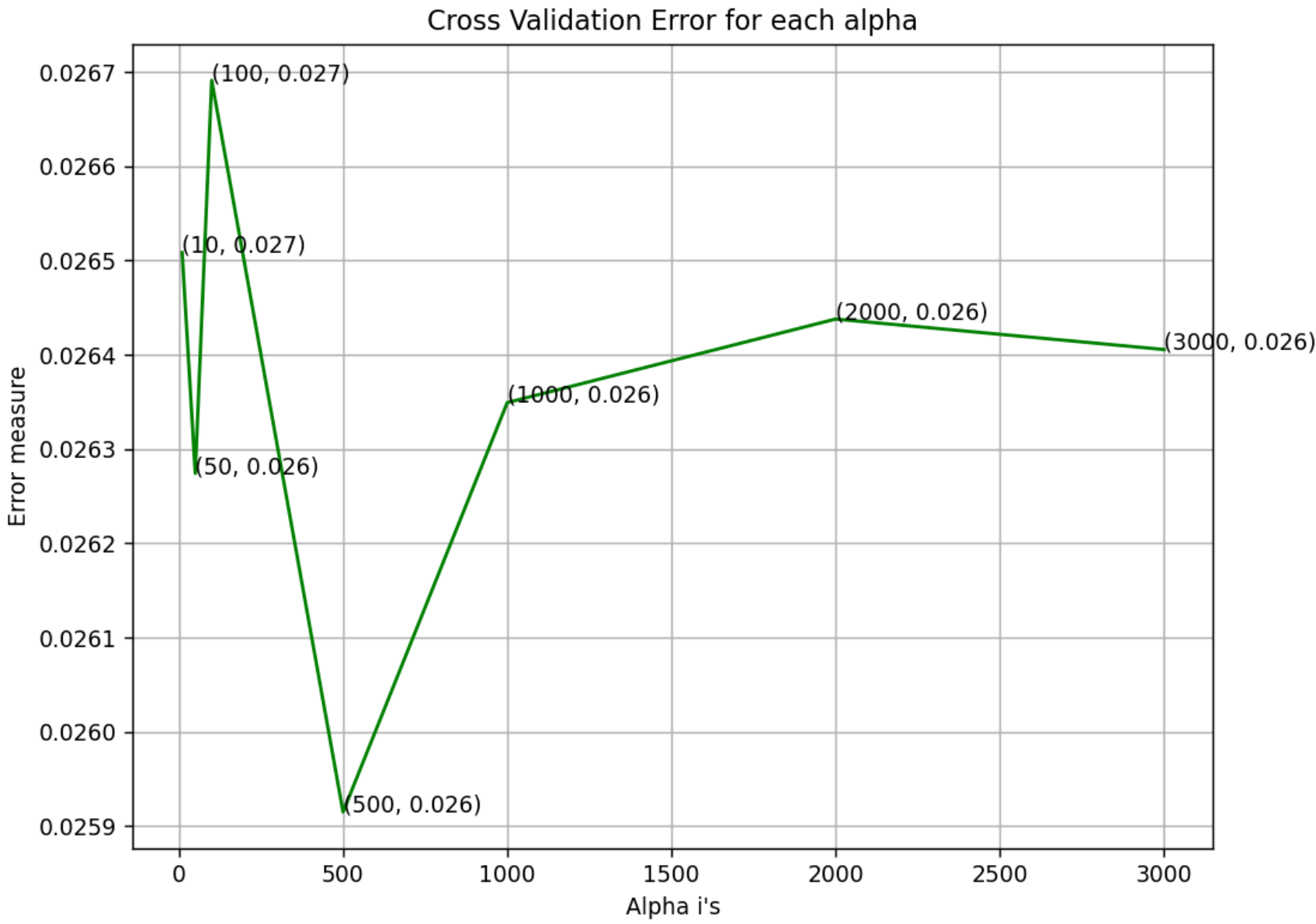
3.2.1 Random Forest Classifier with Asm-uni grams + ASM-Extracted image features + Byte uni gram

In [138]:

```
1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # predict_proba(X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
19 # -----
20
21 alpha=[10,50,100,500,1000,2000,3000]
22 cv_log_error_array=[]
23 train_log_error_array=[]
24 for i in alpha:
25     r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
26     r_cfl.fit(X_train,y_train)
27     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
28     sig_clf.fit(X_train, y_train)
29     predict_y = sig_clf.predict_proba(X_cv)
30     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
31     print('Parameter tuning for alpha:',i,' Loss:',log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
32
33 # for i in range(len(cv_log_error_array)):
34 #     print ('Log_Loss for c = ',alpha[i],'is',cv_log_error_array[i])
35
36 best_alpha = np.argmin(cv_log_error_array)
37
38 fig, ax = plt.subplots()
39 ax.plot(alpha, cv_log_error_array,c='g')
40
41 for i, txt in enumerate(np.round(cv_log_error_array,3)):
42     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
43
44 plt.grid()
45 plt.title("Cross Validation Error for each alpha")
46 plt.xlabel("Alpha i's")
47 plt.ylabel("Error measure")
48 plt.show()
49
50 r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
51 r_cfl.fit(X_train,y_train)
52 sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
53 sig_clf.fit(X_train, y_train)
54
55 predict_y = sig_clf.predict_proba(X_train)
56 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
57 predict_y = sig_clf.predict_proba(X_cv)
58 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
59 predict_y = sig_clf.predict_proba(X_test)
60 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
61
62
```

Parameter tuning for alpha: 10 Loss: 0.026508445703490754
Parameter tuning for alpha: 50 Loss: 0.026274183912199087
Parameter tuning for alpha: 100 Loss: 0.026691217424787736
Parameter tuning for alpha: 500 Loss: 0.025914948184497147
Parameter tuning for alpha: 1000 Loss: 0.02634931812198413
Parameter tuning for alpha: 2000 Loss: 0.02643794647709983
Parameter tuning for alpha: 3000 Loss: 0.026405623759426712

<IPython.core.display.Javascript object>



For values of best alpha = 500 The train log loss is: 0.010860945011193368
For values of best alpha = 500 The cross validation log loss is: 0.025914948184497147
For values of best alpha = 500 The test log loss is: 0.024136455428315266

```
In [139]: 1 # predict_y is test prediction
          2 predicted_y = np.argmax(predict_y, axis=1)
          3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points 0.41398344066237347

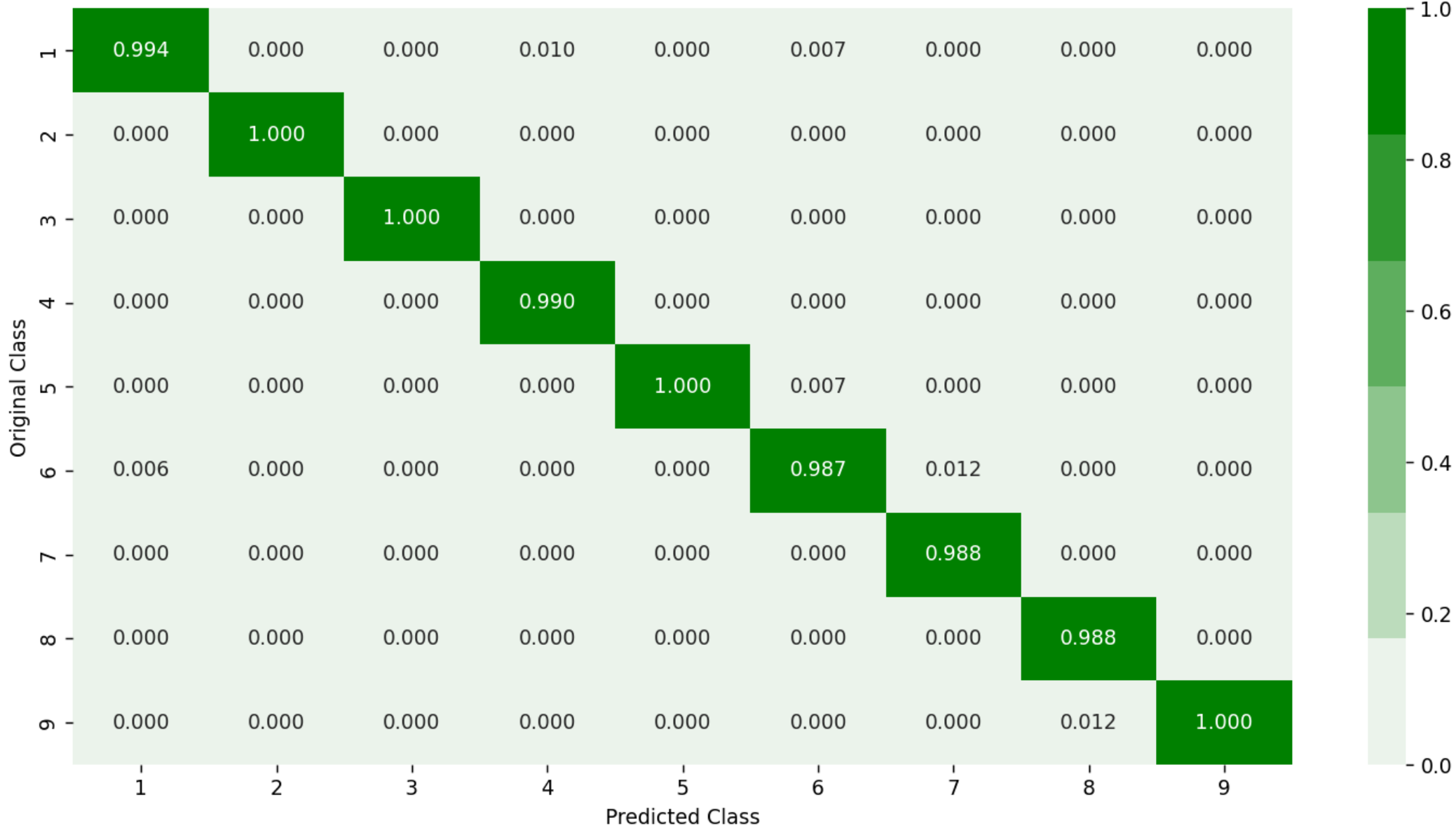
----- Confusion matrix -----

<IPython.core.display.Javascript object>



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

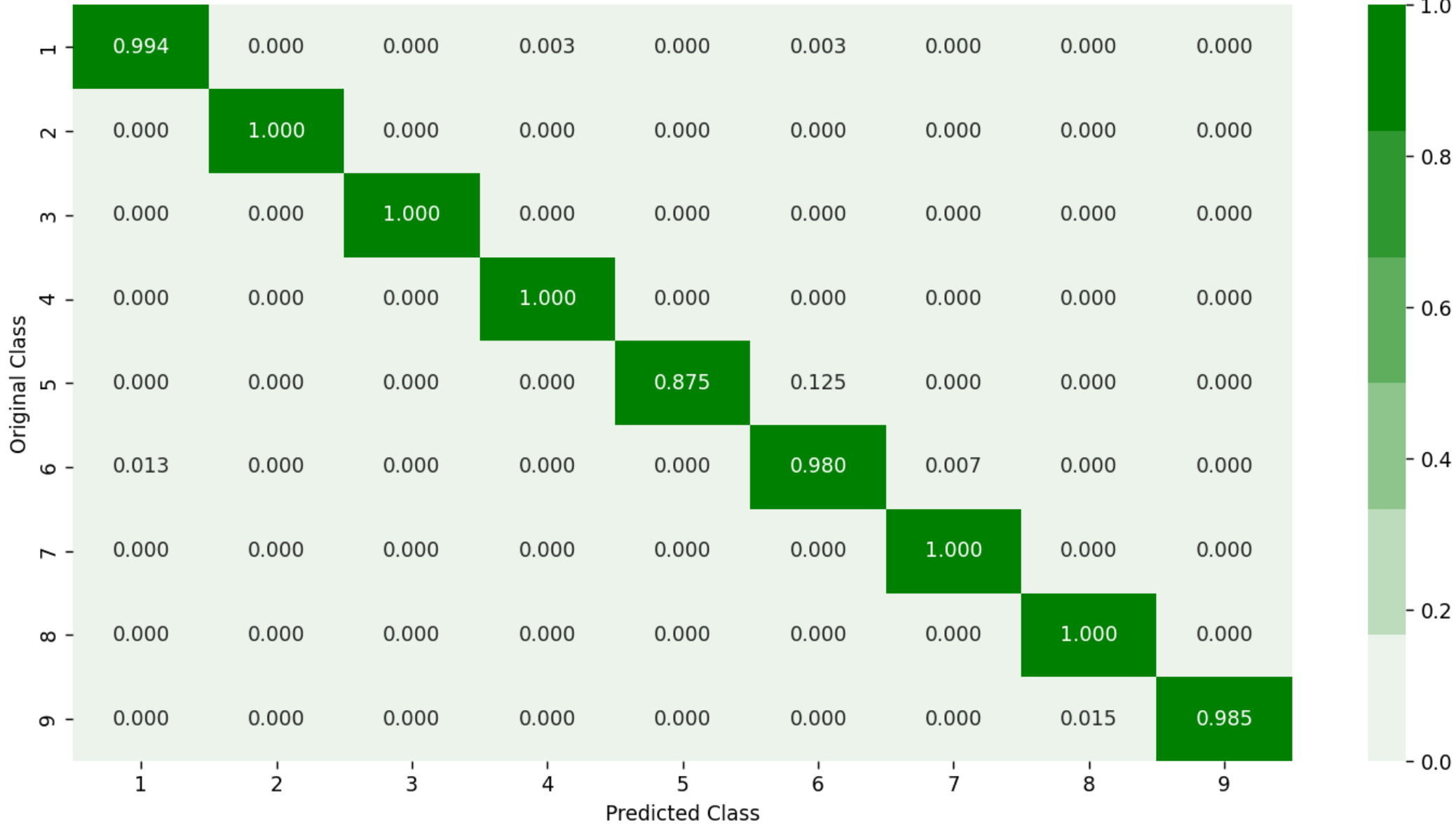
<IPython.core.display.Javascript object>



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

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

<IPython.core.display.Javascript object>



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

3.2.2 Xgboost Classifier with best hyperparameters


```
In [25]: 1 %%time
2 x_cfl=XGBClassifier(silent=True,verbosity=0)
3
4 prams={
5     'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
6     'n_estimators':[100,200,500,1000,2000],
7     'max_depth':[3,5,10,12],
8     'colsample_bytree':[0.1,0.3,0.5,1],
9     'subsample':[0.1,0.3,0.5,1]
10 }
11 random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,cv=4)
12 random_cfl.fit(X_train,y_train)
```

Fitting 4 folds for each of 10 candidates, totalling 40 fits

[CV 1/4; 1/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1

[CV 1/4; 1/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1; score=0.996 total time= 6.6min

[CV 2/4; 1/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1

[CV 2/4; 1/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1; score=0.997 total time= 6.5min

[CV 3/4; 1/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1

[CV 3/4; 1/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1; score=0.993 total time= 6.5min

[CV 4/4; 1/10] START colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1

[CV 4/4; 1/10] END colsample_bytree=0.5, learning_rate=0.1, max_depth=10, n_estimators=1000, subsample=0.1; score=0.997 total time= 6.6min

[CV 1/4; 2/10] START colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1

[CV 1/4; 2/10] END colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1; score=0.997 total time=11.4min

[CV 2/4; 2/10] START colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1

[CV 2/4; 2/10] END colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1; score=0.997 total time=11.3min

[CV 3/4; 2/10] START colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1

[CV 3/4; 2/10] END colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1; score=0.994 total time=10.9min

[CV 4/4; 2/10] START colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1

[CV 4/4; 2/10] END colsample_bytree=0.3, learning_rate=0.05, max_depth=12, n_estimators=1000, subsample=1; score=0.995 total time=11.3min

[CV 1/4; 3/10] START colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3

[CV 1/4; 3/10] END colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3; score=0.996 total time= 5.8min

[CV 2/4; 3/10] START colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3

[CV 2/4; 3/10] END colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3; score=0.995 total time= 5.6min

[CV 3/4; 3/10] START colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3

[CV 3/4; 3/10] END colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3; score=0.994 total time= 5.6min

[CV 4/4; 3/10] START colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3

[CV 4/4; 3/10] END colsample_bytree=0.5, learning_rate=0.03, max_depth=5, n_estimators=200, subsample=0.3; score=0.995 total time= 5.7min

[CV 1/4; 4/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1

[CV 1/4; 4/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1; score=0.997 total time=95.4min

[CV 2/4; 4/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1

[CV 2/4; 4/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1; score=0.995 total time=61.4min

[CV 3/4; 4/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1

[CV 3/4; 4/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1; score=0.996 total time=55.4min

[CV 4/4; 4/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1

[CV 4/4; 4/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=2000, subsample=1; score=0.994 total time=56.0min

[CV 1/4; 5/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1

[CV 1/4; 5/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1; score=0.991 total time= 2.3min

[CV 2/4; 5/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1

[CV 2/4; 5/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1; score=0.989 total time= 2.3min

[CV 3/4; 5/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1

[CV 3/4; 5/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1; score=0.991 total time= 2.3min

[CV 4/4; 5/10] START colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1

[CV 4/4; 5/10] END colsample_bytree=1, learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.1; score=0.994 total time= 2.3min

[CV 1/4; 6/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5

[CV 1/4; 6/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5; score=0.997 total time= 3.7min

[CV 2/4; 6/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5

[CV 2/4; 6/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5; score=0.997 total time= 3.7min

[CV 3/4; 6/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5

[CV 3/4; 6/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5; score=0.995 total time= 3.7min

[CV 4/4; 6/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5

[CV 4/4; 6/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=10, n_estimators=200, subsample=0.5; score=0.995 total time= 3.7min

[CV 1/4; 7/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1

[CV 1/4; 7/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1; score=0.997 total time= 8.6min

[CV 2/4; 7/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1

[CV 2/4; 7/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1; score=0.997 total time= 8.5min

[CV 3/4; 7/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1

[CV 3/4; 7/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1; score=0.994 total time= 8.2min

[CV 4/4; 7/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1

[CV 4/4; 7/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=5, n_estimators=1000, subsample=1; score=0.995 total time= 8.7min

[CV 1/4; 8/10] START colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1

[CV 1/4; 8/10] END colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1; score=0.991 total time=11.5min

[CV 2/4; 8/10] START colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1

[CV 2/4; 8/10] END colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1; score=0.990 total time=11.6min

[CV 3/4; 8/10] START colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1

[CV 3/4; 8/10] END colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1; score=0.992 total time=12.0min

[CV 4/4; 8/10] START colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1

[CV 4/4; 8/10] END colsample_bytree=1, learning_rate=0.01, max_depth=5, n_estimators=1000, subsample=0.1; score=0.993 total time=11.9min

[CV 1/4; 9/10] START colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3

[CV 1/4; 9/10] END colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3; score=0.991 total time= 2.4min

[CV 2/4; 9/10] START colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3

[CV 2/4; 9/10] END colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3; score=0.991 total time= 2.4min

[CV 3/4; 9/10] START colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3

[CV 3/4; 9/10] END colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3; score=0.991 total time= 2.4min

[CV 4/4; 9/10] START colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3

[CV 4/4; 9/10] END colsample_bytree=1, learning_rate=0.03, max_depth=3, n_estimators=100, subsample=0.3; score=0.992 total time= 2.4min

[CV 1/4; 10/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1

[CV 1/4; 10/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1; score=0.988 total time= 1.1min

[CV 2/4; 10/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1

[CV 2/4; 10/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1; score=0.990 total time= 1.1min

[CV 3/4; 10/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1

[CV 3/4; 10/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1; score=0.987 total time= 1.1min

[CV 4/4; 10/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1

[CV 4/4; 10/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=3, n_estimators=100, subsample=1; score=0.991 total time= 1.1min

CPU times: user 14h 3min 2s, sys: 55.2 s, total: 14h 3min 58s

Wall time: 8h 6min 31s

```
Out[25]: RandomizedSearchCV(cv=4,
      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=nan,
      monotone_constraints=None,
      n_estimators=100,...
      random_state=None, reg_alpha=None,
      reg_lambda=None,
      scale_pos_weight=None, silent=True,
      subsample=None, tree_method=None,
      validate_parameters=None,
      verbosity=0),
      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, 12],
      'n_estimators': [100, 200, 500, 1000,
      2000],
      'subsample': [0.1, 0.3, 0.5, 1]},
      verbose=10)
```

```
In [61]: 1 # files = pd.DataFrame(data=os.listdir('byteFiles'),columns=['ID'])
2 # files = pd.DataFrame(data = files['ID'].str.replace('.txt',''),columns=['ID'])
3 # labels = result[['ID','Class']]
```

```
In [65]: 1 # pd.merge(files,labels,on='ID',how='inner').to_csv('bi_gram_file_label_order.csv')
```

```
In [67]: 1 print (random_cfl.best_params_)

{'subsample': 0.5, 'n_estimators': 200, 'max_depth': 10, 'learning_rate': 0.05, 'colsample_bytree': 0.5}
```


In [71]:

```
1 %%time
2 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
3
4 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
5 # -----
6 # default paramters
7 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
8 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1,
9 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1,
10 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
11
12 # some of methods of RandomForestRegressor()
13 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
14 # get_params([deep]) Get parameters for this estimator.
15 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
16 # get_score(importance_type='weight') -> get the feature importance
17 # -----
18 # video link2: https://www.applidaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
19 # -----
20 x_cfl=XGBClassifier(n_estimators=200, subsample=0.5,learning_rate=0.05, colsample_bytree=0.5, max_depth=10)
21 x_cfl.fit(X_train,y_train)
22 c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
23 c_cfl.fit(X_train,y_train)
24
25 predict_y = c_cfl.predict_proba(X_train)
26 print ('train loss',log_loss(y_train, predict_y))
27 predict_y = c_cfl.predict_proba(X_cv)
28
29 print ('cv loss',log_loss(y_cv, predict_y))
30 predict_y = c_cfl.predict_proba(X_test)
31 print ('test loss',log_loss(y_test, predict_y))
```

train loss 0.011415465124124244
cv loss 0.016410075115980253)
test loss 0.01475890827104327
CPU times: user 45min 36s, sys: 714 ms, total: 45min 37s
Wall time: 22min 55s

```
In [101]: 1 # predict_y is test prediction
          2 predicted_y = np.argmax(predict_y, axis=1)
          3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points 0.20198712051517942

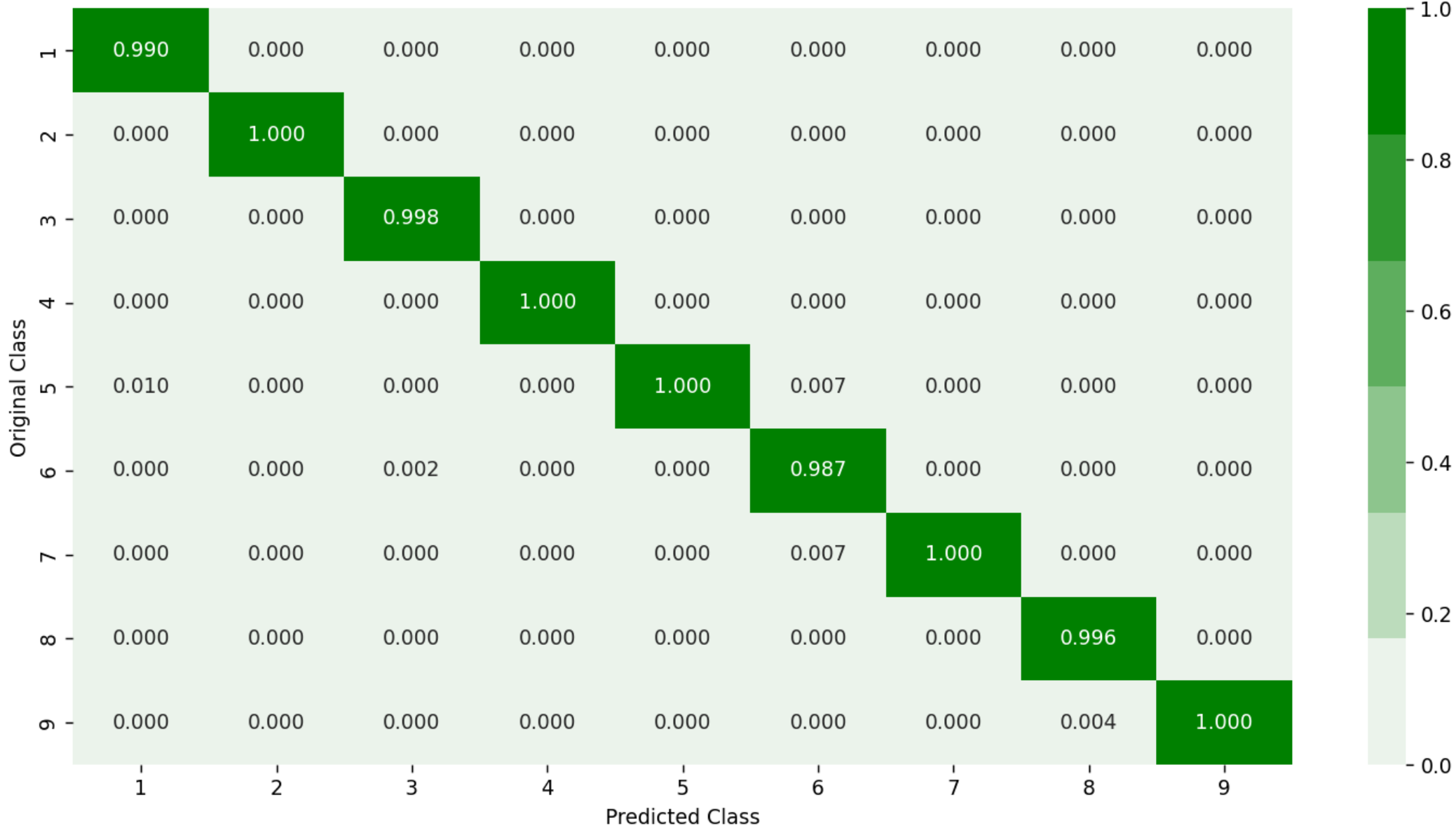
Confusion matrix

<IPython.core.display.Javascript object>



Precision matrix

<IPython.core.display.Javascript object>



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

Recall matrix

<IPython.core.display.Javascript object>



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

3.3 Byte- uni + ASM - Image Features

```
In [85]: 1 X_merged = pd.merge(result.drop(columns=['size']),data_50,on='ID',how='inner')
          2 Y_merged = X_merged['Class']
          3 X_merged = X_merged.drop(columns=['ID'])
          4 print('shape of x',X_merged.shape,' y:',Y_merged.shape)
```

shape of x (10868, 759) y: (10868,)

```
In [86]: 1 X_merged.columns
```

Out[86]: Index(['0', '1', '2', '3', '4', '5', '6', '7', '8', '9',
...
'ASM_990', 'ASM_991', 'ASM_992', 'ASM_993', 'ASM_994', 'ASM_995',
'ASM_996', 'ASM_997', 'ASM_998', 'ASM_999'],
dtype='object', length=759)

```
In [87]: 1 ### train test split
          2 X_train, X_test, y_train, y_test = train_test_split(X_merged, Y_merged,stratify=Y_merged,test_size=0.20)
          3 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,stratify=y_train,test_size=0.20)
```

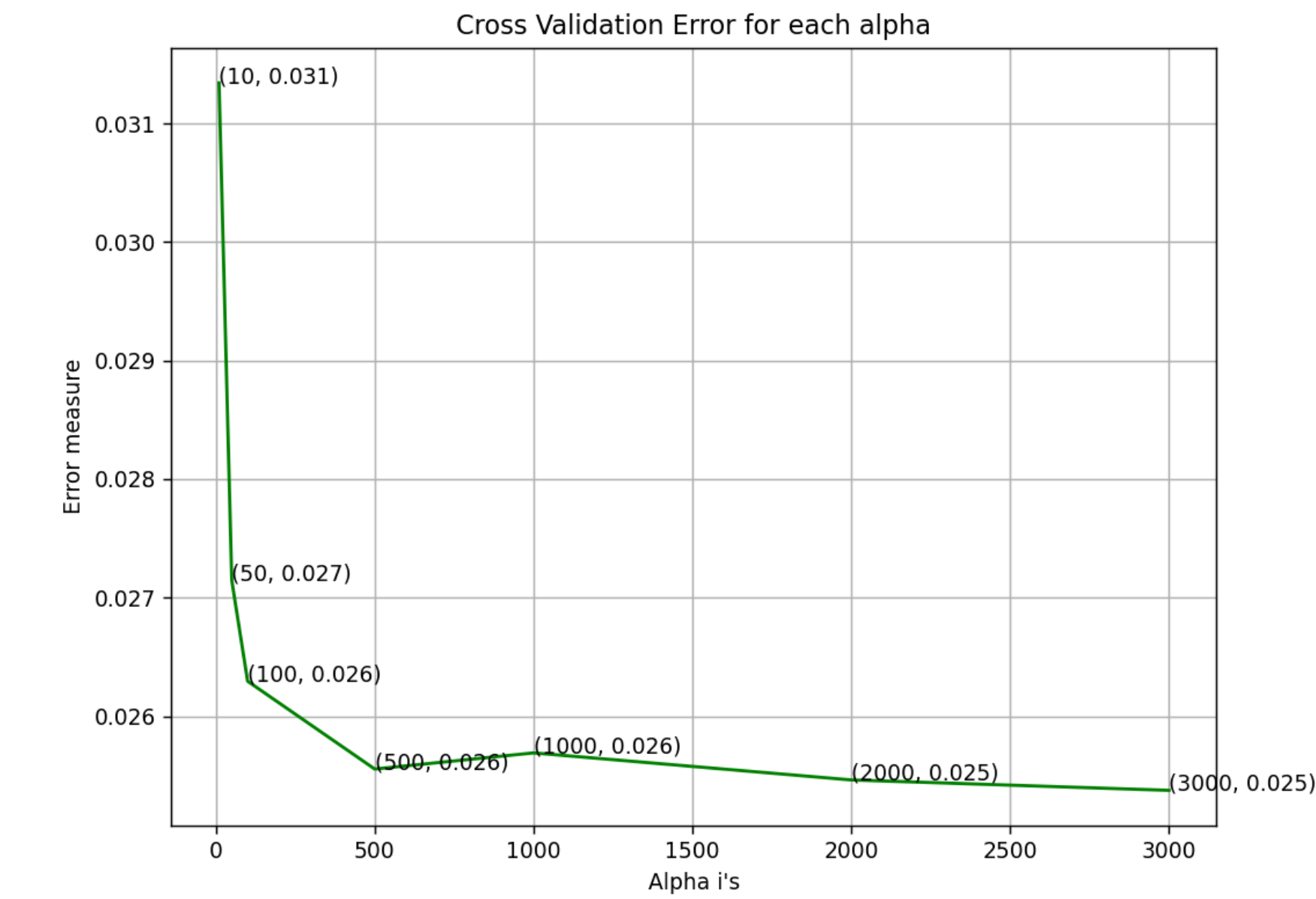
```
In [88]: 1 X_train = X_train.astype(float)
          2 X_test = X_test.astype(float)
          3 X_cv = X_cv.astype(float)
```

3.3.1 Random Forest Classifier for ASM-Extracted image features + Byte uni gram

```
In [89]: 1 # -----
2 # default parameters
3 # sklearn.ensemble.RandomForestClassifier(n_estimators=10, criterion='gini', max_depth=None, min_samples_split=2,
4 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
5 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False,
6 # class_weight=None)
7
8 # Some of methods of RandomForestClassifier()
9 # fit(X, y, [sample_weight]) Fit the SVM model according to the given training data.
10 # predict(X) Perform classification on samples in X.
11 # predict_proba(X) Perform classification on samples in X.
12
13 # some of attributes of RandomForestClassifier()
14 # feature_importances_ : array of shape = [n_features]
15 # The feature importances (the higher, the more important the feature).
16
17 # -----
18 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/random-forest-and-their-construction-2/
19 # -----
20
21 alpha=[10,50,100,500,1000,2000,3000]
22 cv_log_error_array=[]
23 train_log_error_array=[]
24 for i in alpha:
25     r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
26     r_cfl.fit(X_train,y_train)
27     sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
28     sig_clf.fit(X_train, y_train)
29     predict_y = sig_clf.predict_proba(X_cv)
30     cv_log_error_array.append(log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
31     print('Parameter tuning for alpha:',i,' Loss:',log_loss(y_cv, predict_y, labels=r_cfl.classes_, eps=1e-15))
32
33 # for i in range(len(cv_log_error_array)):
34 #     print ('Log_Loss for c = ',alpha[i],'is',cv_log_error_array[i])
35
36
37 best_alpha = np.argmin(cv_log_error_array)
38
39 fig, ax = plt.subplots()
40 ax.plot(alpha, cv_log_error_array,c='g')
41
42 for i, txt in enumerate(np.round(cv_log_error_array,3)):
43     ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
44
45 plt.grid()
46 plt.title("Cross Validation Error for each alpha")
47 plt.xlabel("Alpha i's")
48 plt.ylabel("Error measure")
49 plt.show()
50
51
52 r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=-1)
53 r_cfl.fit(X_train,y_train)
54 sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
55 sig_clf.fit(X_train, y_train)
56
57 predict_y = sig_clf.predict_proba(X_train)
58 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y))
59 predict_y = sig_clf.predict_proba(X_cv)
60 print('For values of best alpha = ', alpha[best_alpha], "The cross validation log loss is:",log_loss(y_cv, predict_y))
61 predict_y = sig_clf.predict_proba(X_test)
62 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y))
63
```

Parameter tuning for alpha: 10 Loss: 0.0313426542602312
Parameter tuning for alpha: 50 Loss: 0.02714685745594213
Parameter tuning for alpha: 100 Loss: 0.026295166349999472
Parameter tuning for alpha: 500 Loss: 0.025554102636726757
Parameter tuning for alpha: 1000 Loss: 0.025690196753958717
Parameter tuning for alpha: 2000 Loss: 0.025462073418468168
Parameter tuning for alpha: 3000 Loss: 0.025374208756105795

<IPython.core.display.Javascript object>



For values of best alpha = 3000 The train log loss is: 0.010686512051691533
For values of best alpha = 3000 The cross validation log loss is: 0.025374208756105795
For values of best alpha = 3000 The test log loss is: 0.03933044530569612

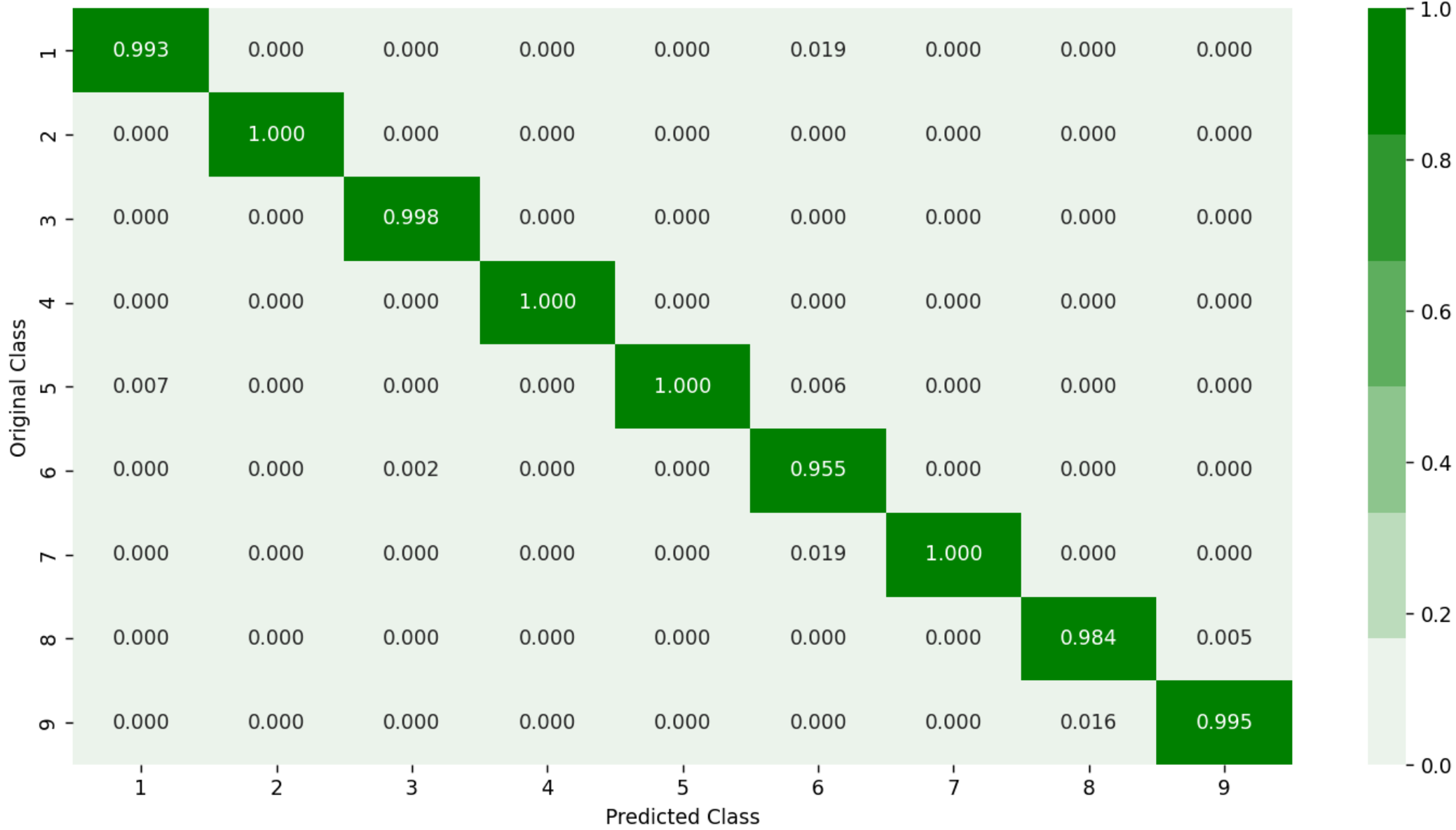

```
In [90]: 1 # predict_y is test prediction
        2 predicted_y = np.argmax(predict_y, axis=1)
        3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points 0.6899724011039559

----- Confusion matrix -----
<IPython.core.display.Javascript object>



----- Precision matrix -----
<IPython.core.display.Javascript object>



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



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

3.3.2 XgBoost Classifier with best params for ASM-Extracted image features + Byte uni gram

In [103]:

```
1 %%time
2 x_cfl=XGBClassifier(silent=True,verbosity=0)
3
4 prams={
5     'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
6     'n_estimators':[100,200,500,1000,2000],
7     'max_depth':[3,5,10,12],
8     'colsample_bytree':[0.1,0.3,0.5,1],
9     'subsample':[0.1,0.3,0.5,1]
10 }
11 random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=10,cv=4)
12 random_cfl.fit(X_train,y_train)
```

Fitting 4 folds for each of 10 candidates, totalling 40 fits

[CV 1/4; 1/10] START colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1

[CV 1/4; 1/10] END colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1, score=0.999 total time= 39.7s

[CV 2/4; 1/10] START colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1

[CV 2/4; 1/10] END colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1, score=0.999 total time= 40.1s

[CV 3/4; 1/10] START colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1

[CV 3/4; 1/10] END colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1, score=0.998 total time= 39.5s

[CV 4/4; 1/10] START colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1

[CV 4/4; 1/10] END colsample_bytree=0.5, learning_rate=0.2, max_depth=5, n_estimators=500, subsample=0.1, score=1.000 total time= 1.3min

[CV 1/4; 2/10] START colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5

[CV 1/4; 2/10] END colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5, score=0.998 total time= 2.8min

[CV 2/4; 2/10] START colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5

[CV 2/4; 2/10] END colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5, score=1.000 total time= 2.9min

[CV 3/4; 2/10] START colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5

[CV 3/4; 2/10] END colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5, score=0.999 total time= 2.8min

[CV 4/4; 2/10] START colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5

[CV 4/4; 2/10] END colsample_bytree=1, learning_rate=0.15, max_depth=10, n_estimators=200, subsample=0.5, score=1.000 total time= 2.8min

[CV 1/4; 3/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3

[CV 1/4; 3/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3, score=0.999 total time=17.1min

[CV 2/4; 3/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3

[CV 2/4; 3/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3, score=1.000 total time=17.2min

[CV 3/4; 3/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3

[CV 3/4; 3/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3, score=0.999 total time=17.1min

[CV 4/4; 3/10] START colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3

[CV 4/4; 3/10] END colsample_bytree=1, learning_rate=0.01, max_depth=12, n_estimators=1000, subsample=0.3, score=1.000 total time=17.1min

[CV 1/4; 4/10] START colsample_bytree=0.1, learning_rate=0.03, max_depth=5, n_estimators=500, subsample=0.5

[CV 1/4; 4/10] END colsample_bytree=0.1, learning_rate=0.03, max_depth=5, n_estimators=500, subsample=0.5, score=0.997 total time= 1.7min

[CV 2/4; 4/10] START colsample_bytree=0.1, learning_rate=0.03, max_depth=5, n_estimators=500, subsample=0.5

[CV 2/4; 4/10] END colsample_bytree=0.1, learning_rate=0.03, max_depth=5, n_estimators=500, subsample=0.5, score=0.998 total time= 1.7min

[CV 3/4; 4/10] START colsample_bytree=0.1, learning_rate=0.03, max_depth=5, n_estimators=500, subsample=0.5, score=0.996 total time= 1.7min

[CV 4/4; 4/10] START colsample_bytree=0.1, learning_rate=0.03, max_depth=5, n_estimators=500, subsample=0.5, score=0.997 total time= 1.7min

[CV 1/4; 5/10] START colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1

[CV 1/4; 5/10] END colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1, score=0.999 total time= 29.3s

[CV 2/4; 5/10] START colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1

[CV 2/4; 5/10] END colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1, score=0.997 total time= 29.7s

[CV 3/4; 5/10] START colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1

[CV 3/4; 5/10] END colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1, score=0.997 total time= 29.4s

[CV 4/4; 5/10] START colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1

[CV 4/4; 5/10] END colsample_bytree=0.3, learning_rate=0.15, max_depth=12, n_estimators=200, subsample=0.1, score=0.998 total time= 29.6s

[CV 1/4; 6/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1

[CV 1/4; 6/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1, score=0.999 total time= 2.1min

[CV 2/4; 6/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1

[CV 2/4; 6/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1, score=1.000 total time= 2.1min

[CV 3/4; 6/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1

[CV 3/4; 6/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1, score=0.999 total time= 2.1min

[CV 4/4; 6/10] START colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1

[CV 4/4; 6/10] END colsample_bytree=0.3, learning_rate=0.1, max_depth=3, n_estimators=500, subsample=1, score=0.999 total time= 2.1min

[CV 1/4; 7/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1

[CV 1/4; 7/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1, score=0.995 total time= 1.4min

[CV 2/4; 7/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1

[CV 2/4; 7/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1, score=0.997 total time= 1.4min

[CV 3/4; 7/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1

[CV 3/4; 7/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1, score=0.994 total time= 1.4min

[CV 4/4; 7/10] START colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1

[CV 4/4; 7/10] END colsample_bytree=0.3, learning_rate=0.01, max_depth=5, n_estimators=100, subsample=1, score=0.998 total time= 1.4min

[CV 1/4; 8/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1

[CV 1/4; 8/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1, score=0.998 total time= 1.5min

[CV 2/4; 8/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1

[CV 2/4; 8/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1, score=1.000 total time= 1.5min

[CV 3/4; 8/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1

[CV 3/4; 8/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1, score=0.999 total time= 1.5min

[CV 4/4; 8/10] START colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1

[CV 4/4; 8/10] END colsample_bytree=0.5, learning_rate=0.15, max_depth=3, n_estimators=200, subsample=1, score=1.000 total time= 1.5min

[CV 1/4; 9/10] START colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3

[CV 1/4; 9/10] END colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3, score=0.994 total time= 1.6min

[CV 2/4; 9/10] START colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3

[CV 2/4; 9/10] END colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3, score=0.992 total time= 1.6min

[CV 3/4; 9/10] START colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3

[CV 3/4; 9/10] END colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3, score=0.989 total time= 1.6min

[CV 4/4; 9/10] START colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3

[CV 4/4; 9/10] END colsample_bytree=0.1, learning_rate=0.01, max_depth=5, n_estimators=500, subsample=0.3, score=0.996 total time= 1.6min

[CV 1/4; 10/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1

[CV 1/4; 10/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1, score=0.997 total time= 34.2s

[CV 2/4; 10/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1

[CV 2/4; 10/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1, score=0.993 total time= 34.0s

[CV 3/4; 10/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1

[CV 3/4; 10/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1, score=0.993 total time= 33.8s

[CV 4/4; 10/10] START colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1

[CV 4/4; 10/10] END colsample_bytree=0.5, learning_rate=0.05, max_depth=3, n_estimators=100, subsample=0.1, score=0.996 total time= 34.0s

CPU times: user 4h 45min 26s, sys: 15.3 s, total: 4h 45min 41s

Wall time: 2h 23min 39s

```
Out[103]: RandomizedSearchCV(cv=4,
                             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=nan,
                                                       monotone_constraints=None,
                                                       n_estimators=100,...
                                                       random_state=None, reg_alpha=None,
                                                       reg_lambda=None,
                                                       scale_pos_weight=None, silent=True,
                                                       subsample=None, tree_method=None,
                                                       validate_parameters=None,
                                                       verbosity=0),
                             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, 12],
                                                  'n_estimators': [100, 200, 500, 1000,
                                                                  2000],
                                                  'subsample': [0.1, 0.3, 0.5, 1]},
                             verbose=10)
```

In [106]:

```
1 print (random_cfl.best_params_)

{'subsample': 0.3, 'n_estimators': 1000, 'max_depth': 12, 'learning_rate': 0.01, 'colsample_bytree': 1}
```

In [125]:

```
1 %%time
2 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
3
4 # find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#xgboost.XGBClassifier
5 # -----
6 # default paramters
7 # class xgboost.XGBClassifier(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
8 # objective='binary:logistic', booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1,
9 # max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1,
10 # scale_pos_weight=1, base_score=0.5, random_state=0, seed=None, missing=None, **kwargs)
11
12 # some of methods of RandomForestRegressor()
13 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
14 # get_params([deep]) Get parameters for this estimator.
15 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
16 # get_score(importance_types='weight') -> get the feature importance
17 # -----
18 # video Link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
19 # -----
20 x_cfl=XGBClassifier(n_estimators=1000, subsample=0.3,learning_rate=0.01, colsample_bytree=1, max_depth=12)
21 x_cfl.fit(X_train,y_train)
22 c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
23 c_cfl.fit(X_train,y_train)
24
25 predict_y = random_cfl.best_estimator_.predict_proba(X_train)
26 print ('train loss',log_loss(y_train, predict_y))
27 predict_y = random_cfl.best_estimator_.predict_proba(X_cv)
28
29 print ('cv loss',log_loss(y_cv, predict_y))
30 predict_y = random_cfl.best_estimator_.predict_proba(X_test)
31 print ('test loss',log_loss(y_test, predict_y))
```

train loss 0.002346527399823316

cv loss 0.004251100181973078

test loss 0.007054805073482494

CPU times: user 1.58 s, sys: 5.76 ms, total: 1.59 s

Wall time: 802 ms

```
In [126]: 1 # predict_y is test prediction
          2 predicted_y = np.argmax(predict_y, axis=1)
          3 plot_confusion_matrix(y_test, predicted_y+1)
```

Number of misclassified points -0.028003679852805885

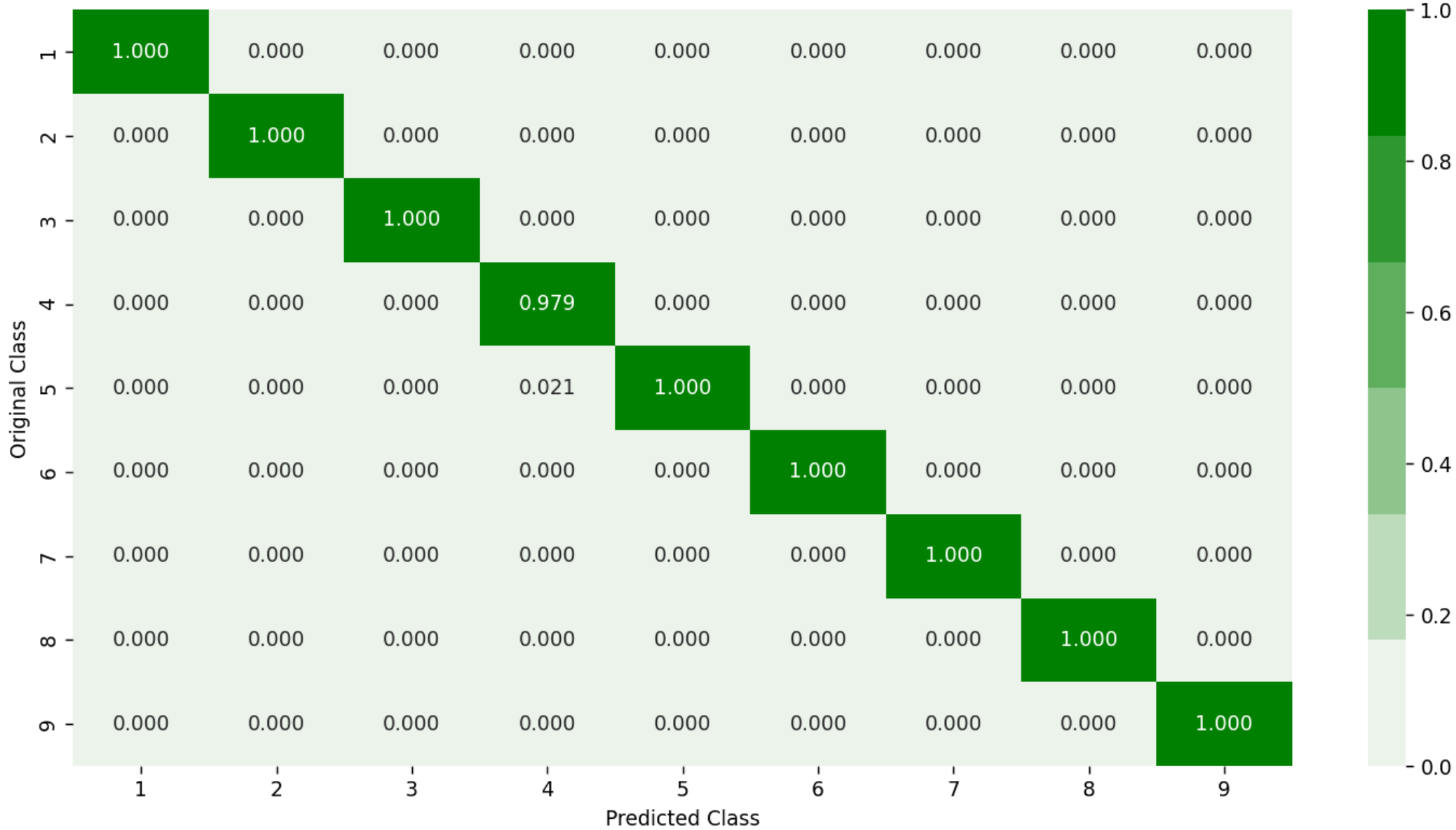
----- Confusion matrix -----

<IPython.core.display.Javascript object>



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

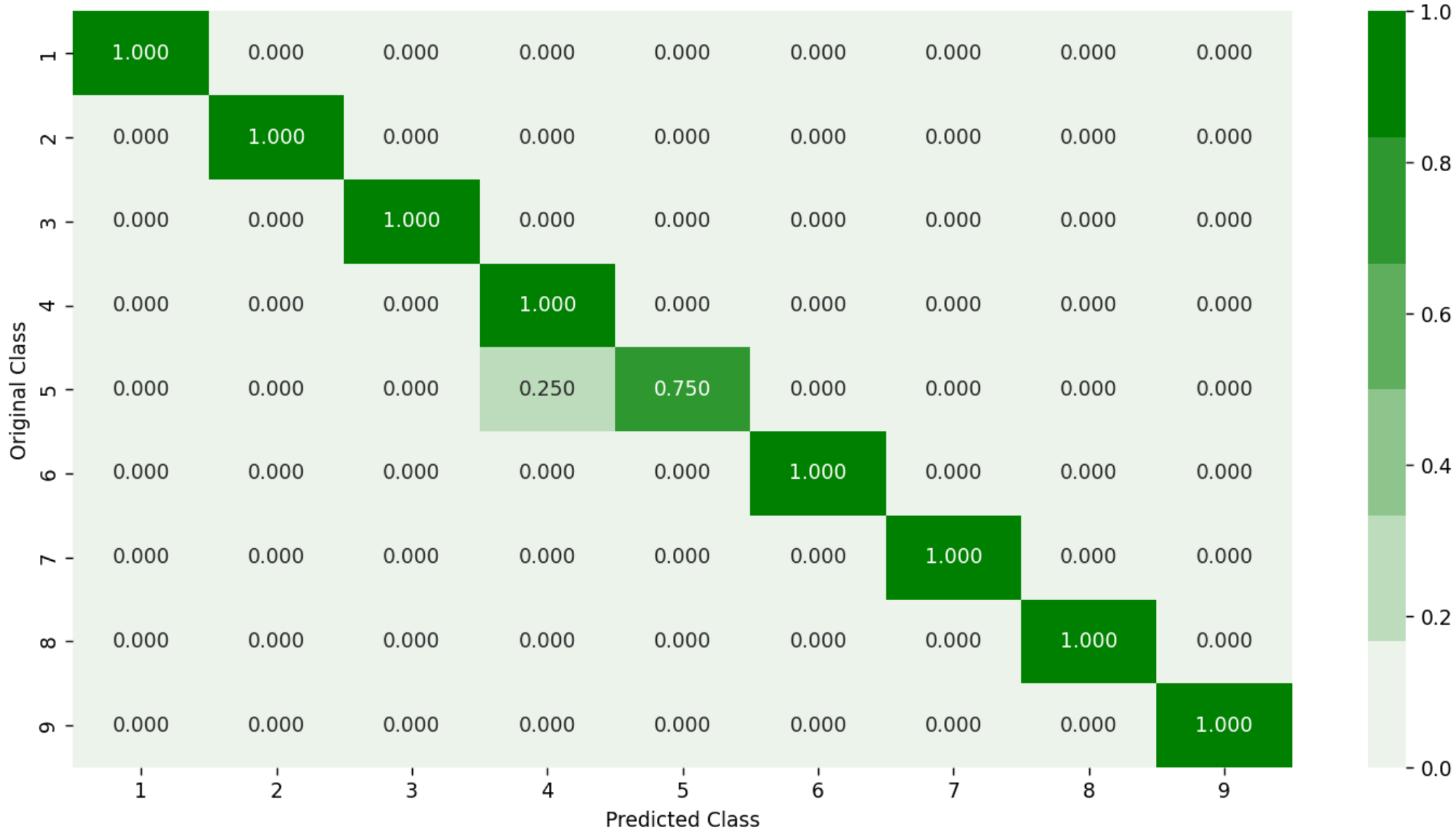
<IPython.core.display.Javascript object>



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

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

<IPython.core.display.Javascript object>



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

Final Summary


```
In [139]: 1 from prettytable import PrettyTable
2 table = PrettyTable()
3 table.field_names = ['Model', 'Features', 'Best HyperParameter', 'Test Log Loss']
4
5 table.add_row(['Random Forest Classifier', 'Byte-Bigrams', 500, 0.0538])
6 table.add_row(['Random Forest Classifier', 'Byte-Bigrams + ASM Image Features', 50, 0.0206])
7 table.add_row(['Random Forest Classifier', 'ASM-unigrams + ASM Image Features', 100, 0.0238])
8 table.add_row(['Random Forest Classifier', 'Byte-unigrams + ASM Image Features', 3000, 0.0393])
9 table.add_row(['Random Forest Classifier', 'Byte-unigrams + ASM Image Features + ASM-unigrams', 500, 0.0241])
10 table.add_row(['', '', '', ''])
11 table.add_row(['XGB(best_params) Classifier', 'Byte-Bigrams', '500,max_depth=3', 0.034])
12 table.add_row(['XGB(best_params) Classifier', 'Byte-Bigrams + ASM Image Features', 500, 0.034])
13 table.add_row(['XGB(best_params) Classifier', 'ASM-unigrams + ASM Image Features', '500,max_depth=10', 0.0253])
14 table.add_row(['XGB(best_params) Classifier', 'Byte-unigrams + ASM Image Features', '1000,0.007'])
15 table.add_row(['XGB(best_params) Classifier', 'Byte-unigrams + ASM Image Features + ASM-unigrams', '200,max_depth=10', 0.0147])
16
17 table.add_row(['XGB Classifier', 'ASM-unigrams + ASM image features', 500, 0.0238])
18
19 print(table)
```

Model	Features	Best HyperParameter	Test Log Loss
Random Forest Classifier	Byte-Bigrams	500	0.0538
Random Forest Classifier	Byte-Bigrams + ASM Image Features	50	0.0206
Random Forest Classifier	ASM-unigrams + ASM Image Features	100	0.0238
Random Forest Classifier	Byte-unigrams + ASM Image Features	3000	0.0393
Random Forest Classifier	Byte-unigrams + ASM Image Features + ASM-unigrams	500	0.0241
XGB(best_params) Classifier	Byte-Bigrams	500,max_depth=3	0.034
XGB(best_params) Classifier	Byte-Bigrams + ASM Image Features	500	0.05
XGB(best_params) Classifier	ASM-unigrams + ASM Image Features	500,max_depth=10	0.0253
XGB(best_params) Classifier	Byte-unigrams + ASM Image Features	1000	0.007
XGB(best_params) Classifier	Byte-unigrams + ASM Image Features + ASM-unigrams	200,max_depth=10	0.0147
XGB Classifier	ASM-unigrams + ASM image features	500	0.0238

Conclusion :

Carried out the following note-book in the following sequence--

- Created byte-bigrams from the byte files and stored the generated 66000 unique keys in pickle file.
- created bi-gram datamatrix with(10868 rows,~66000 columns).
- Using Randomforest classifier performed feature selection and took top 200 bibram features.
- Created ASM image features then chi square test ,selected optimal features which stores 50% variance.
- Trained multiple models like (RandomForest,XGB) models combination of features like
 - Byte-Bigrams
 - Byte-Bigrams + ASM Image Features
 - ASM-unigrams + ASM Image Features
 - Byte-unigrams + ASM Image Features
 - Byte-unigrams + ASM Image Features + ASM-unigrams

XGB classifier with best parameters featureized using Byte-unigrams + ASM Image Features performed best on test data reducing loss to 0.007