

MICROSOFT MALWARE DETECTION ASSIGNMENT'

IN THIS PROBLEM STATEMENT WE HAVE THE RAW DATA WITH 9 TYPES OF MALWARE WITH TWO TYPES OF FILES THE FILES ARE

- .ASM FILES
- .BYTE FILES

WE ARE JUST GIVEN WITH THESE DATA WE HAVE TO FIND THE MALWARE EACH PARTICULAR FILE IT BELONGS TO

MALWARES WE HAVE ARE

- RAMMIT
- LOLLIPOP
- KELIHOS
- VUNDO
- SIMDA
- TRACUR
- KELIHOS VERSION1
- OBFUSCVATOR
- GATAK

WE HAVE THE 10868 BYTE FILES AND .ASM FILES WE HAVE TO PERFORM THE EXPLORATORY DATA ANALYSIS, FEATURE EXTRACTION, FEATURE ENGINEERING, MODEL SELECTION, HYPERPARAMETER TUNING THE MODELS AND SHOULD TAKE THE STEPS TO

```
In [0]: import pandas as pd
```

```
In [0]: import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use('nbAgg')
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process# this is used for multithreading
import multiprocessing
import codecs# this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\ensemble\weight_boosting.py:
29: DeprecationWarning: numpy.core.umath_tests is an internal NumPy module and
should not be imported. It will be removed in a future NumPy release.
    from numpy.core.umath_tests import inner1d
```

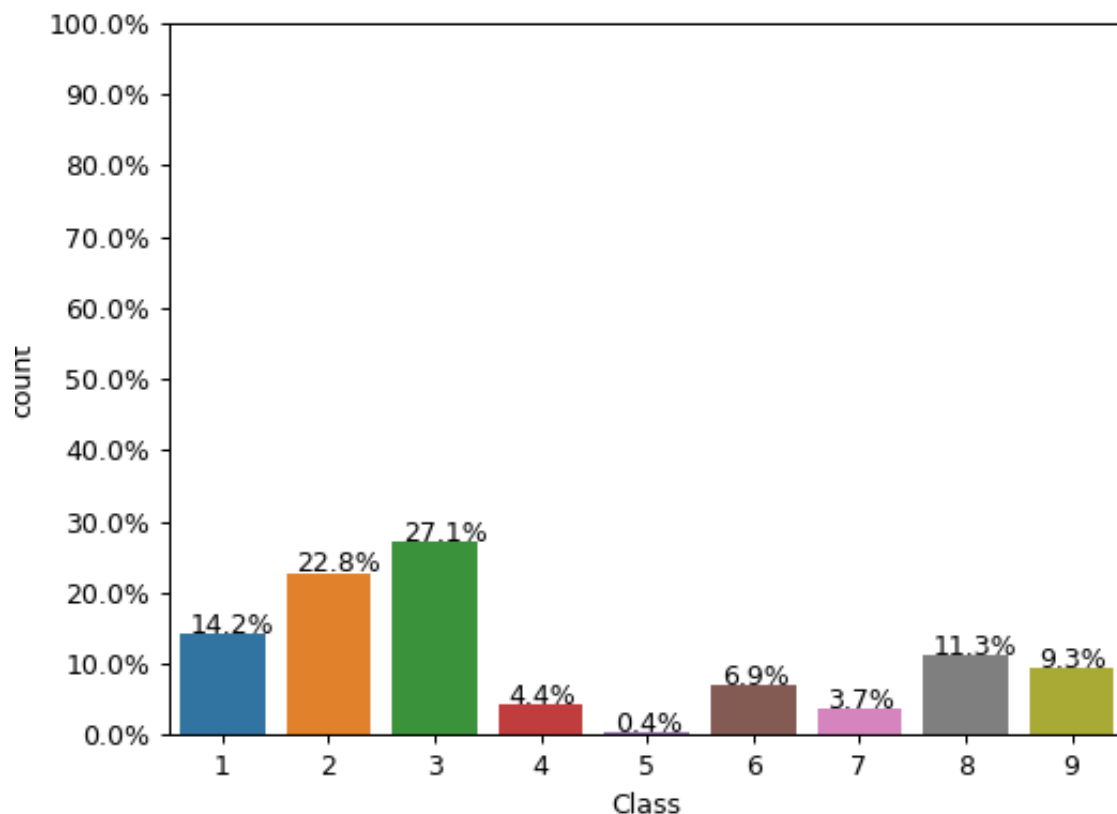
THIS IS THE PROCESS OF EXPLORATORY DATA ANALYSIS CHECKING THE IMBALANCE IN THE DATA WE CAN APPLY STARATEGY TO HANDLE THE DATA IMBALANACE BY TAKING MEASURES OF UNDERSAMPLING AND OVER SAMPLING OR TAKING CARE WE DIVIDE THE DATA INTO TRAIN TEST SPLITTING WE CAN APPLLY STRATEGY TO SPLIT THE DATA BASED ON THE Y LABEL

```
In [0]: Y=pd.read_csv("trainLabels.csv")
total = len(Y)*1.
ax=sns.countplot(x="Class", data=Y)
for p in ax.patches:
    ax.annotate('{:.1f}%'.format(100*p.get_height()/total), (p.get_x()+0.1, p.get_y()+0.1))

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

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

<IPython.core.display.Javascript object>



THIS IS THE IMPORTANT STEP OF FEATURE EXTRACTION WHERE WE SEE THE SIZE OF THE FILE WHICH IS THE IMPORTANT THING IN THE PREDICTION OF THE MALWARE BECAUSE THE MALWARE CAN BE BELONG TO MANY SIZES LATER WE CAN APPEND THIS AS THE FEATURE.

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

```
files=os.listdir('byteFiles')
filenames=Y['Id'].tolist()
class_y=Y['Class'].tolist()
class_bytes=[]
sizebytes=[]
fnames=[]
for file in files:
    # print(os.stat('byteFiles/0A32eTdBKayjCWhZqDOQ.txt'))
    # os.stat_result(st_mode=33206, st_ino=1125899906874507, st_dev=3561571700, s
    # st_size=3680109, st_atime=1519638522, st_mtime=1519638522, st_ctime=1519638
    # read more about os.stat: here https://www.tutorialspoint.com/python/os_stat
    statinfo=os.stat('byteFiles/'+file)
    # split the file name at '.' and take the first part of it i.e the file name
    file=file.split('.')[0]
    if any(file == filename for filename in filenames):
        i=filenames.index(file)
        class_bytes.append(class_y[i])
        # converting into Mb's
        sizebytes.append(statinfo.st_size/(1024.0*1024.0))
        fnames.append(file)
data_size_byte=pd.DataFrame({'ID':fnames,'size':sizebytes,'Class':class_bytes})
print (data_size_byte.head())
```

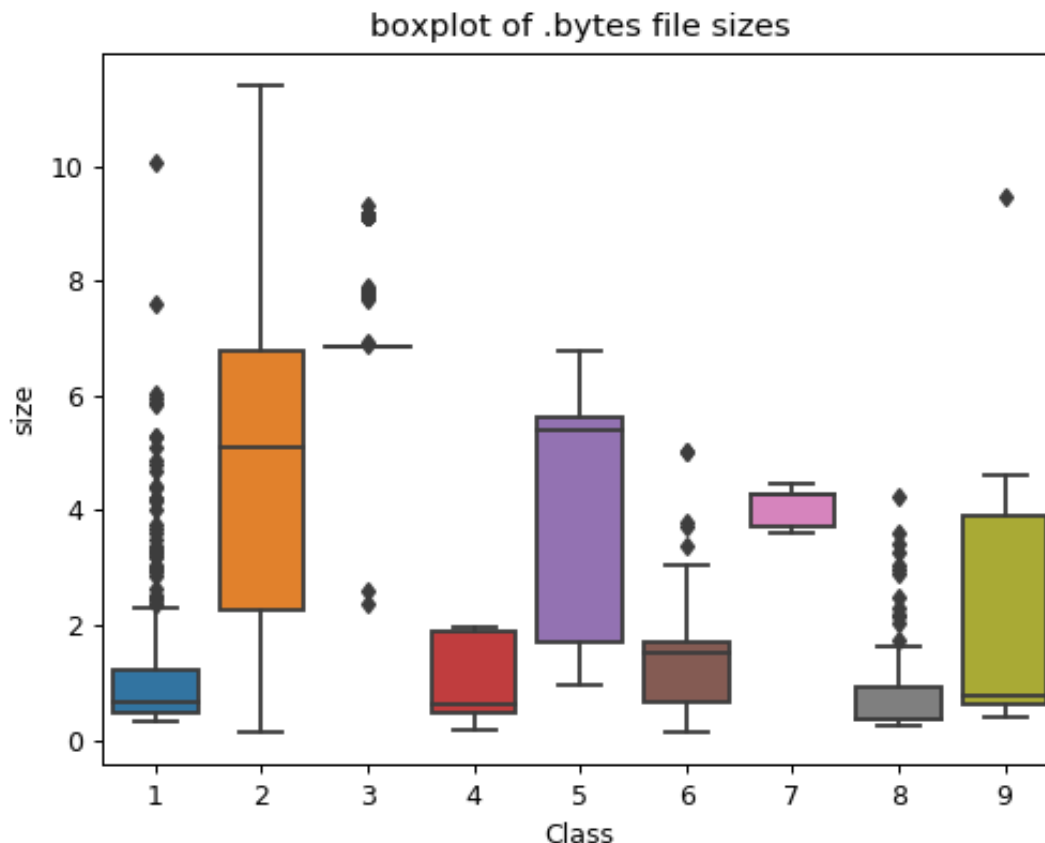
	ID	size	Class
0	05aiMRw13bYWqZ80Hvj1	5.945557	2
1	05EeG39MTRrI6VY21DPd	0.430664	1
2	05IXcWGxvnkto4sq17zZ	1.303955	2
3	05Kps4iFw8mOLJZQrb1H	6.220703	2
4	05LHG8fR3iPn6agIo9z7	2.380615	6

In [0]: data_size_byte.to_csv('classsize.csv')

In [0]: data_size_byte=pd.read_csv('classsize.csv')

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

<IPython.core.display.Javascript object>



ASSIGNMENT 1 STARTS FROM HERE WE WILL TAKE THE BIGRAMS OF THE DATA IN FILES AND VISUALISE THE PERFORMANCE USING THE MODELS.

THE CELL BELOW SHOWS THE GENERATION OF UNIGRAMS FORM THE BYTE FILES. AS THIS IS ALREADY DONE AS A PART OF ASSIGNMENT WE GONNA CONSIDER THE BIGRAMS GENERATED FROM THE FILES NAD OBTAIN THE PERFORMANCE.

```

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

files = os.listdir('byteFiles')
filenames=[]
array=[]
for file in files:
    if(file.endswith("bytes")):
        file=file.split('.')[0]
        text_file = open('byteFiles/'+file+".txt", 'w+')
        with open('byteFiles/'+file+".bytes", "r") as fp:
            lines=""

            for line in fp:
                a=line.rstrip().split(" ")[1:]
                b=' '.join(a)
                b=b+"\n"
                text_file.write(b)
            fp.close()
            os.remove('byteFiles/'+file+".bytes")
        text_file.close()

```

```

In [0]: byte_features1=byte_features.iloc[:, :258]
byte_features1=byte_features1['ID'].map(lambda x: x.replace('.txt', ''))

```


GENERATION OF BIGRAMS AND FEATURISE THEM WITH THE BASED ON THE HEXADECIMALDATA PRESENT INSIDE THE EACH FILE AND APPENDING THEM WITH THE SIZE AND CLASS LABEL BEFORE APPLYING THE MODELS.

```
In [0]: print(256*256)
```

65536

```
In [0]: files = os.listdir('byteFiles')
filenames2=[]
feature_matrix = np.zeros((len(files),65537),dtype=int)
```

```
In [0]: k=0
byte_feature_file=open('resultfin.csv','w+')
for file in files:
    filenames2.append(file)
    byte_feature_file.write(file+",")
    if(file.endswith(".txt")):
        with open('byteFiles/'+file,"r") as byte_flie:
            for lines in byte_flie:
                line=line.rstrip().split(" ")
                for i in range(len(line)):
                    if i<(len(line)-1):
                        a=line[i]
                        b=line[i+1]
                        if (a!='??') and (b!='??'):
                            score=arr[int(a,16)][int(b,16)]
                            feature_matrix[k][score]+=1

            byte_flie.close()
        for i in feature_matrix[k]:

            byte_feature_file.write(str(i)+",")
            byte_feature_file.write("\n")

        k += 1

byte_feature_file.close()
```



```
In [0]: import pandas as pd
das=pd.read_csv('resultfin.csv',nrows=10)
print(das)
print(das.shape)
```

```
05aiMRw13bYWqZ80Hvj1.txt 0 25141 782 140 77 211 21 25 27 \
0 05EeG39MTRrI6VY21DPd.txt 0 10352 178 74 122 77 41 13 30
1 05IXcWGxvnkto4sq17zZ.txt 0 8195 144 63 32 34 7 8 8
2 05Kps4iFw8mOLJZQrb1H.txt 0 16493 1185 560 2149 439 285 278 282
3 05LHG8fR3iPn6agIo9z7.txt 0 93134 325 256 403 251 237 273 311
4 05rJTUWYAKNegBk2wE8X.txt 0 37601 533 317 138 233 55 59 693
5 065EZhxgbLRSHsB87uIF.txt 0 20862 711 63 38 171 21 21 14
6 06aLOj8EUXMByS423sum.txt 0 26357 1014 457 491 518 371 317 193
7 06arUi9q3wHS2C8RZxeB.txt 0 5729 44 28 83 13 47 81 31
8 06KfrF7ltESna2ZHPVp5.txt 0 6122 33 23 70 13 56 16 38
9 06osXqPUVM1HbvBGNncT.txt 0 4051 21 16 44 14 29 10 26
```

```
... 220.1 133.2 8.1020 18.91 20.79 5.998 24.46 22.74 \
0 ... 2 9 1 3 2 8 3 5
1 ... 7 2 1 0 0 1 2 6
2 ... 22 12 0 8 18 7 14 13
3 ... 118 131 109 113 103 112 118 122
4 ... 39 24 3 5 2 14 13 17
5 ... 37 62 49 87 227 40 38 124
6 ... 43 29 33 53 86 73 108 111
7 ... 9 11 12 12 10 8 13 10
8 ... 15 7 8 13 14 9 20 16
9 ... 12 14 11 16 9 12 15 10
```

```
6836 Unnamed: 65538
0 512 NaN
1 2405 NaN
2 3445 NaN
3 29213 NaN
4 14747 NaN
5 8864 NaN
6 7107 NaN
7 74 NaN
8 177 NaN
9 86 NaN
```

```
[10 rows x 65539 columns]
(10, 65539)
```

**THE FILES ARE UPLOADED INTO THE DRIVE
AFTER GENERATING THE IMAGE FEATURES
AND MODELS ARE APPLIED ON THEM**

we have obtained the bigrams of the byte files using the hexadecimal data .

we got the features and named the features we got finally the 65538 features

```
In [0]: # Code to read csv file into Colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

100% |████████████████████████████████████████████████████████████████████████████████| 993kB 20.2MB/s ta 0:00:01
Building wheel for PyDrive (setup.py) ... done
```

```
In [0]: link = 'https://drive.google.com/open?id=1cphD9gkKoD80jw8No_K9Wv2nyp1kNUSs' # The
```

```
In [0]: fluff, id = link.split('=')
print(id) # Verify that you have everything after '='

1cphD9gkKoD80jw8No_K9Wv2nyp1kNUSs
```

generating the column names for byte files ranging to 65539

```
In [0]: l=[i for i in range(65539)]
l[0]='ID'
print(l)

['ID', 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20,
21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 4
0, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 5
9, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 7
8, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 9
7, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 1
13, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 12
8, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143,
144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 15
9, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174,
175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 19
0, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205,
206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 22
1, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236,
237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 25
2, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267,
268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 28
3, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298,
299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 31
4, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329,
330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344,
345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359,
360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374,
375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389,
390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404,
405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419,
420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434,
435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449,
450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464,
465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479,
480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494,
495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509,
510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524,
525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539,
540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554,
555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569,
570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584,
585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599,
600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614,
615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629,
630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644,
645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659,
660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674,
675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689,
690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704,
705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719,
720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734,
735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749,
750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764,
765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779,
780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794,
795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809,
810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824,
825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839,
840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854,
855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869,
870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884,
885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899,
900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914,
915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929,
930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944,
945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959,
960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974,
975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989,
990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000, 1001, 1002, 1003,
1004, 1005, 1006, 1007, 1008, 1009, 1010, 1011, 1012, 1013, 1014, 1015, 1016,
1017, 1018, 1019, 1020, 1021, 1022, 1023, 1024, 1025, 1026, 1027, 1028, 1029,
1030, 1031, 1032, 1033, 1034, 1035, 1036, 1037, 1038, 1039, 1040, 1041, 1042,
1043, 1044, 1045, 1046, 1047, 1048, 1049, 1050, 1051, 1052, 1053, 1054, 1055,
1056, 1057, 1058, 1059, 1060, 1061, 1062, 1063, 1064, 1065, 1066, 1067, 1068,
1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080, 1081,
1082, 1083, 1084, 1085, 1086, 1087, 1088, 1089, 1090, 1091, 1092, 1093, 1094,
1095, 1096, 1097, 1098, 1099, 1100, 1101, 1102, 1103, 1104, 1105, 1106, 1107,
1108, 1109, 1110, 1111, 1112, 1113, 1114, 1115, 1116, 1117, 1118, 1119, 1120,
1121, 1122, 1123, 1124, 1125, 1126, 1127, 1128, 1129, 1130, 1131, 1132, 1133,
1134, 1135, 1136, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144, 1145, 1146,
1147, 1148, 1149, 1150, 1151, 1152, 1153, 1154, 1155, 1156, 1157, 1158, 1159,
1160, 1161, 1162, 1163, 1164, 1165, 1166, 1167, 1168, 1169, 1170, 1171, 1172,
1173, 1174, 1175, 1176, 1177, 1178, 1179, 1180, 1181, 1182, 1183, 1184, 1185,
1186, 1187, 1188, 1189, 1190, 1191, 1192, 1193, 1194, 1195, 1196, 1197, 1198,
1199, 1200, 1201, 1202, 1203, 1204, 1205, 1206, 1207, 1208, 1209, 1210, 1211,
1212, 1213, 1214, 1215, 1216, 1217, 1218, 1219, 1220, 1221, 1222, 1223, 1224,
1225, 1226, 1227, 1228, 1229, 1230, 1231, 1232, 1233, 1234, 1235, 1236, 1237,
1238, 1239, 1240, 1241, 1242, 1243, 1244, 1245, 1246, 1247, 1248, 1249, 1250,
1251, 1252, 1253, 1254, 1255, 1256, 1257, 1258, 1259, 1260, 1261, 1262, 1263,
1264, 1265, 1266, 1267, 1268, 1269, 1270, 1271, 1272, 1273, 1274, 1275, 1276,
1277, 1278, 1279, 1280, 1281, 1282, 1283, 1284, 1285, 1286, 1287, 1288, 1289,
1290, 1291, 1292, 1293, 1294, 1295, 1296, 1297, 1298, 1299, 1300, 1301, 1302,
1303, 1304, 1305, 1306, 1307, 1308, 1309, 1310, 1311, 1312, 1313, 1314, 1315,
1316, 1317, 1318, 1319, 1320, 1321, 1322, 1323, 1324, 1325, 1326, 1327, 1328,
1329, 1330, 1331, 1332, 1333, 1334, 1335, 1336, 1337, 1338, 1339, 1340, 1341,
1342, 1343, 1344, 1345, 1346, 1347, 1348, 1349, 1350, 1351, 1352, 1353, 1354,
1355, 1356, 1357, 1358, 1359, 1360, 1361, 1362, 1363, 1364, 1365, 1366, 1367,
1368, 1369, 1370, 1371, 1372, 1373, 1374, 1375, 1376, 1377, 1378, 1379, 1380,
1381, 1382, 1383, 1384, 1385, 1386, 1387, 1388, 1389, 1390, 1391, 1392, 1393,
1394, 1395, 1396, 1397, 1398, 1399, 1400, 1401, 1402, 1403, 1404, 1405, 1406,
1407, 1408, 1409, 1410, 1411, 1412, 1413, 1414, 1415, 1416, 1417, 1418, 1419,
1420, 1421, 1422, 1423, 1424, 1425, 1426, 1427, 1428, 1429, 1430, 1431, 1432,
1433, 1434, 1435, 1436, 1437, 1438, 1439, 1440, 1441, 1442, 1443, 1444, 1445,
1446, 1447, 1448, 1449, 1450, 1451, 1452, 1453, 1454, 1455, 1456, 1457, 1458,
1459, 1460, 1461, 1462, 1463, 1464, 1465, 1466, 1467, 1468, 1469, 1470, 1471,
1472, 1473, 1474, 1475, 1476, 1477, 1478, 1479, 1480, 1481, 1482, 1483, 1484,
1485, 1486, 1487, 1488, 1489, 1490, 1491, 1492, 1493, 1494, 1495, 1496, 1497,
1498, 1499, 1500, 1501, 1502, 1503, 1504, 1505, 1506, 1507, 1508, 1509, 1510,
1511, 1512, 1513, 1514, 1515, 1516, 1517, 1518, 1519, 1520, 1521, 1522, 1523,
1524, 1525, 1526, 1527, 1528, 1529, 1530, 1531, 1532, 1533, 1534, 1535, 1536,
1537, 1538, 1539, 1540, 1541, 1542, 1543, 1544, 1545, 1546, 1547, 1548, 1549,
1550, 1551, 1552, 1553, 1554, 1555, 1556, 1557, 1558, 1559, 1560, 1561, 1562,
1563, 1564, 1565, 1566, 1567, 1568, 1569, 1570, 1571, 1572, 1573, 1574, 1575,
1576, 1577, 1578, 1579, 1580, 1581, 1582, 1583, 1584, 1585, 1586, 1587, 1588,
1589, 1590, 1591, 1592, 1593, 1594, 1595, 1596, 1597, 1598, 1599, 1600, 1601,
1602, 1603, 1604, 1605, 1606, 1607, 1608, 1609, 1610, 1611, 1612, 1613, 1614,
1615, 1616, 1617, 1618, 1619, 1620, 1621, 1622, 1623, 1624, 1625, 1626, 1627,
1628, 1629, 1630, 1631, 1632, 1633, 1634, 1635, 1636, 1637, 1638, 1639, 1640,
1641, 1642, 1643, 1644, 1645, 1646, 1647, 1648, 1649, 1650, 1651, 1652, 1653,
1654, 1655, 1656, 1657, 1658, 1659, 1660, 1661, 1662, 1663, 1664, 1665, 1666,
1667, 1668, 1669, 1670, 1671, 1672, 1673, 1674, 1675, 1676, 1677, 1678, 1679,
1680, 1681, 1682, 1683, 1684, 1685, 1686, 1687, 1688, 1689, 1690, 1691, 1692,
1693, 1694, 1695, 1696, 1697, 1698, 1699, 1700, 1701, 1702, 1703, 1704, 1705,
1706, 1707, 1708, 1709, 1710, 1711, 1712, 1713, 1714, 1715, 1716, 1717, 1718,
1719, 1720, 1721, 1722, 1723, 1724, 1725, 1726, 1727, 1728, 1729, 1730, 1731,
1732, 1733, 1734, 1735, 1736, 1737, 1738, 1739, 1740, 1741, 1742, 1743, 1744,
1745, 1746, 1747, 1748, 1749, 1750, 1751, 1752, 1753, 1754, 1755, 1756, 1757,
1758, 1759, 1760, 1761, 1762, 1763, 1764, 1765, 1766, 1767, 1768, 1769, 1770,
1771, 1772, 1773, 1774, 1775, 1776, 1777, 1778, 1779, 1780, 1781, 1782, 1783,
1784, 1785, 1786, 1787, 1788, 1789, 1790, 1791, 1792, 1793, 1794, 1795, 1796,
1797, 1798, 1799, 1800, 1801, 1802, 1803, 1804, 1805, 1806, 1807, 1808, 1809,
1810, 1811, 1812, 1813, 1814, 1815, 1816, 1817, 1818, 1819, 1820, 1821, 1822,
1823, 1824, 1825, 1826, 1827, 1828, 1829, 1830, 1831, 1832, 1833, 1834, 1835,
1836, 1837, 1838, 1839, 1840, 1841, 1842, 1843, 1844, 1845, 1846, 1847, 1848,
1849, 1850, 1851, 1852, 1853, 1854, 1855, 1856, 1857, 1858, 1859, 1860, 1861,
1862, 1863, 1864, 1865, 1866, 1867, 1868, 1869, 1870, 1871, 1872, 1873, 1874,
1875, 1876, 1877, 1878, 1879, 1880, 1881, 1882, 1883, 1884, 1885, 1886, 1887,
1888, 1889, 1890, 1891, 1892, 1893, 1894, 1895, 1896, 1897, 1898, 1899, 1900,
1901, 1902, 1903, 1904, 1905, 1906, 1907, 1908, 1909, 1910, 1911, 1912, 1913,
1914, 1915, 1916, 1917, 1918, 1919, 1920, 1921, 1922, 1923, 1924, 1925, 1926,
1927, 1928, 1929, 1930, 1931, 1932, 1933, 1934, 1935, 1936, 1937, 1938, 1939,
1940, 1941, 1942, 1943, 1944, 1945, 1946, 1947, 1948, 1949, 1950, 1951, 1952,
1953, 1954, 1955, 1956, 1957, 1958, 1959, 1960, 1961, 1962, 1963, 1964, 1965,
1966, 1967, 1968, 1969, 1970, 1971, 1972, 1973, 1974, 1975, 1976, 1977, 1978,
1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991,
1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004,
2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017,
2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030,
2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043,
2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056,
2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069,
2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082,
2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095,
2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108,
2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121,
2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134,
2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 
```

```
In [0]: import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('resultfin.csv')
df3 = pd.read_csv('resultfin.csv',names=1)
```

```
In [0]: print(df3.shape)
```

```
In [0]: df3.head(5)
```

```
In [0]: df4=df3['ID'].map(lambda x: x.replace('.txt',''))
df4.head(5)
df3['ID']=df4
df3.head(5)
df3=df3.drop(columns=[65538])
```

```
In [0]: import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use('nbAgg')
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process# this is used for multithreading
import multiprocessing
import codecs# this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

```
In [0]: print(df3.columns)
```

```
In [0]: # Code to read csv file into Colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)
```

```
In [0]: link = 'https://drive.google.com/open?id=1EeuLtqkPWV875y0YUwTN-jeOzWCrS8Ac' # The
```

```
In [0]: fluff, id = link.split('=')
print(id) # Verify that you have everything after '='

1EeuLtqkPWV875y0YUwTN-jeOzWCrS8Ac
```

```
In [0]: import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('classsize.csv')
dataframe = pd.read_csv('classsize.csv')
```

```
In [0]: dataframe.head(5)
print(dataframe.columns)
dataframe=dataframe.drop(columns=['Unnamed: 0'])
dataframe.head(5)

Index(['Unnamed: 0', 'ID', 'size', 'Class'], dtype='object')
```

```
Out[6]:
```

	ID	size	Class
0	05aiMRw13bYWqZ8OHvjI	5.945557	2
1	05EeG39MTRrI6VY21DPd	0.430664	1
2	05IXcWGxvnkto4sq17zZ	1.303955	2
3	05Kps4iFw8mOLJZQrb1H	6.220703	2
4	05LHG8fR3iPn6aglo9z7	2.380615	6

```
In [0]: df3.head(5)
```

```
Out[16]:
```

	ID	1	2	3	4	5	6	7	8	9	...	65528	65529	65530
0	05aiMRw13bYWqZ8OHvjI	0	25141	782	140	77	211	21	25	27	...	49	220	1
1	05EeG39MTRrI6VY21DPd	0	10352	178	74	122	77	41	13	30	...	1	2	
2	05IXcWGxvnkto4sq17zZ	0	8195	144	63	32	34	7	8	8	...	10	7	
3	05Kps4iFw8mOLJZQrb1H	0	16493	1185	560	2149	439	285	278	282	...	135	22	
4	05LHG8fR3iPn6aglo9z7	0	93134	325	256	403	251	237	273	311	...	120	118	1

5 rows × 65538 columns

we have obtained the features and we will add one more feature which is the size of the malware and we append the class label to the data. later we will use the models to train and test

```
In [0]: result=pd.merge(df3,dataframe,on='ID',how='left')
result.head(5)
```

```
In [0]: result.to_csv('resultingdataframe.csv')
```

```
In [0]: !pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

# Authenticate and create the PyDrive client.
# This only needs to be done once in a notebook.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

# Create & upload a file.
uploaded = drive.CreateFile({'title': 'resultingdataframe.csv'})
uploaded.SetContentFile('resultingdataframe.csv')
uploaded.Upload()
print('Uploaded file with ID {}'.format(uploaded.get('id')))
```

we have uploaded th resultant file into google drive. start executing the file fro the below.otherwise it rewrites the existing file.

```
In [0]: #resulting dataframe csv file link
https://drive.google.com/open?id=1buRmJX5bRd1y3S_t3RV1hkYNvz7tOx96
```

```
In [0]: # Code to read csv file into Colaboratory:  
!pip install -U -q PyDrive  
from pydrive.auth import GoogleAuth  
from pydrive.drive import GoogleDrive  
from google.colab import auth  
from oauth2client.client import GoogleCredentials  
# Authenticate and create the PyDrive client.  
auth.authenticate_user()  
gauth = GoogleAuth()  
gauth.credentials = GoogleCredentials.get_application_default()  
drive = GoogleDrive(gauth)  
  
100% | ████████████████████████████████████████ | 993kB 23.1MB/s ta 0:00:01  
Building wheel for PyDrive (setup.py) ... done
```

```
In [0]: link = 'https://drive.google.com/open?id=1buRmJX5bRd1y3S_t3RV1hkYNvz7t0x96' # The
```

```
In [0]: fluff, id = link.split('=')
print (id) # Verify that you have everything after '='
```

```
1buRmJX5bRd1y3S_t3RV1hkYNvz7t0x96
```

```
In [0]: import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('resultingdataframe.csv')
```

```
In [0]: result = pd.read_csv('resultingdataframe.csv')
```

```
In [0]: result=result.set_index('Unnamed: 0')
```

```
In [0]: print(result.head(5))
```

		ID	1	2	3	4	5	6	7	8	\
Unnamed: 0											
0	05aiMRw13bYWqZ80Hvj1	0	25141	782	140	77	211	21	25		
1	05EeG39MTRrI6VY21DPd	0	10352	178	74	122	77	41	13		
2	05IXcWGxvnkto4sq17zZ	0	8195	144	63	32	34	7	8		
3	05Kps4iFw8mOLJZQrb1H	0	16493	1185	560	2149	439	285	278		
4	05LHG8fR3iPn6agIo9z7	0	93134	325	256	403	251	237	273		
	9 ...	65530	65531	65532	65533	65534	65535	65536	\		
Unnamed: 0	...										
0	27 ...	133	8	18	20	5	24	22			
1	30 ...	9	1	3	2	8	3	5			
2	8 ...	2	1	0	0	1	2	6			
3	282 ...	12	0	8	18	7	14	13			
4	311 ...	131	109	113	103	112	118	122			
	65537	size	Class								
Unnamed: 0											
0	6836	5.945557	2								
1	512	0.430664	1								
2	2405	1.303955	2								
3	3445	6.220703	2								
4	29213	2.380615	6								

```
[5 rows x 65540 columns]
```

APPENDING THE SIZE AND CLASS FOR THE DATA SPLITTING THE DATA BASED ON THE Y LABEL INTO TRAIN CROSS VALIDATAION AND TEST DATA .

```
In [0]: data_y = result['Class']
# split the data into test and train by maintaining same distribution of output variable

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

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

```
Number of data points in train data: 1931
Number of data points in test data: 604
Number of data points in cross validation data: 483
(1931, 65538)
(1931,)
(604, 65538)
(604,)
(483, 65538)
(483,)
```

now we have the data of bigrams we train the models based on the data and predict the y labels

In [0]:

```

# This function plots the confusion matrices given y_i, y_i_hat.
%matplotlib inline
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    print("Number of misclassified points ", (len(test_y)-np.trace(C))/len(test_y))
    # C = 9,9 matrix, each cell (i,j) represents number of points of class i are j

    A = (((C.T)/(C.sum(axis=1))).T)

    B = (C/C.sum(axis=0))
    #divid each element of the confusion matrix with the sum of elements in that
    # C = [[1, 2],
    #      [3, 4]]
    # C.sum(axis= 0)  axis=0 corresponds to columns and axis=1 corresponds to rows
    # C.sum(axis=0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                      [3/4, 4/6]]
    plt.figure(figsize=(15,8))
    labels = [1,2,3,4,5,6,7,8,9]
    cmap=sns.light_palette("green")
    # representing A in heatmap format
    print("-"*50, "Confusion matrix", "-"*50)
    plt.figure(figsize=(10,5))
    sns.heatmap(C, annot=True, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()

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

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

```



```

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

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

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

print('-'*80)

test_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in test data')
plt.grid()
plt.show()

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

print('-'*80)

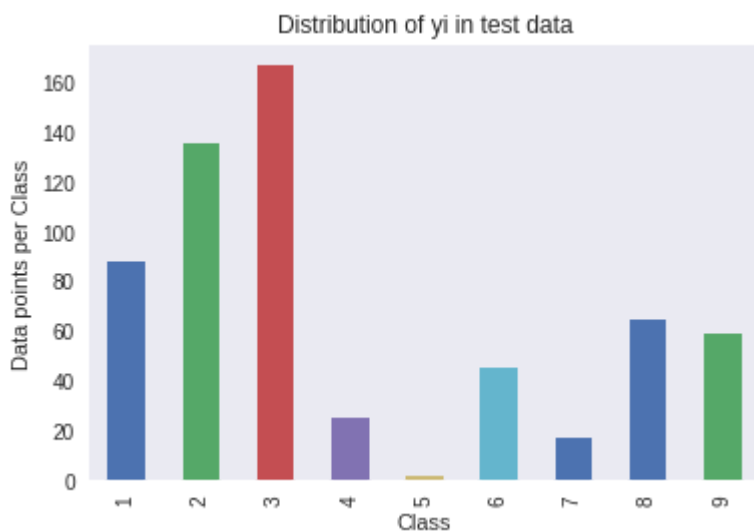
cv_class_distribution.plot(kind='bar')
plt.xlabel('Class')
plt.ylabel('Data points per Class')
plt.title('Distribution of yi in cross validation data')
plt.grid()
plt.show()

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

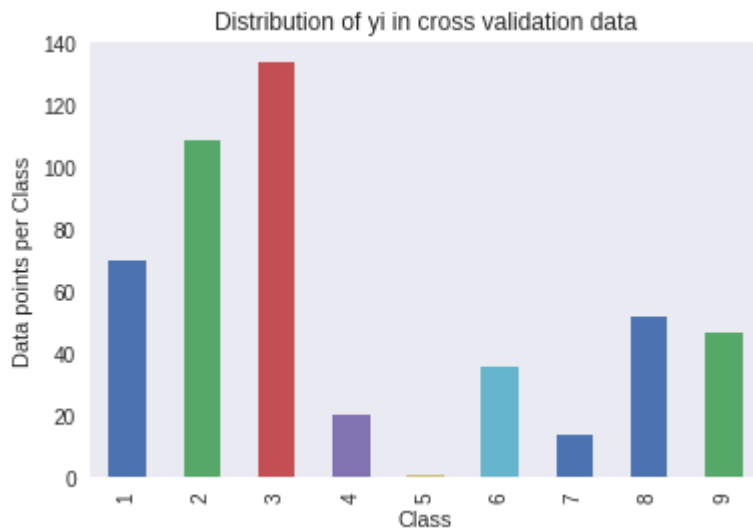
```



Number of data points in class 3 : 534 (27.654 %)
 Number of data points in class 2 : 433 (22.424 %)
 Number of data points in class 1 : 280 (14.5 %)
 Number of data points in class 8 : 210 (10.875 %)
 Number of data points in class 9 : 190 (9.839 %)
 Number of data points in class 6 : 144 (7.457 %)
 Number of data points in class 4 : 78 (4.039 %)
 Number of data points in class 7 : 56 (2.9 %)
 Number of data points in class 5 : 6 (0.311 %)



Number of data points in class 3 : 167 (27.649 %)
 Number of data points in class 2 : 136 (22.517 %)
 Number of data points in class 1 : 88 (14.57 %)
 Number of data points in class 8 : 65 (10.762 %)
 Number of data points in class 9 : 59 (9.768 %)
 Number of data points in class 6 : 45 (7.45 %)
 Number of data points in class 4 : 25 (4.139 %)
 Number of data points in class 7 : 17 (2.815 %)
 Number of data points in class 5 : 2 (0.331 %)



Number of data points in class 3 : 134 (27.743 %)
Number of data points in class 2 : 109 (22.567 %)
Number of data points in class 1 : 70 (14.493 %)
Number of data points in class 8 : 52 (10.766 %)
Number of data points in class 9 : 47 (9.731 %)
Number of data points in class 6 : 36 (7.453 %)
Number of data points in class 4 : 20 (4.141 %)
Number of data points in class 7 : 14 (2.899 %)
Number of data points in class 5 : 1 (0.207 %)

#APPLYING THE MACHINE LEARNING MODELS ON TOP OF DATA GENERAITED AFTER BIGRAMS

4. Machine Learning Models

4.1. Machine Leaning Models on bytes files

4.1.1. Random Model

```

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

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

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

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

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

plot_confusion_matrix(y_test, predicted_y+1)

```

Log loss on Cross Validation Data using Random Model 2.4623613440107714

Log loss on Test Data using Random Model 2.5097805028149645

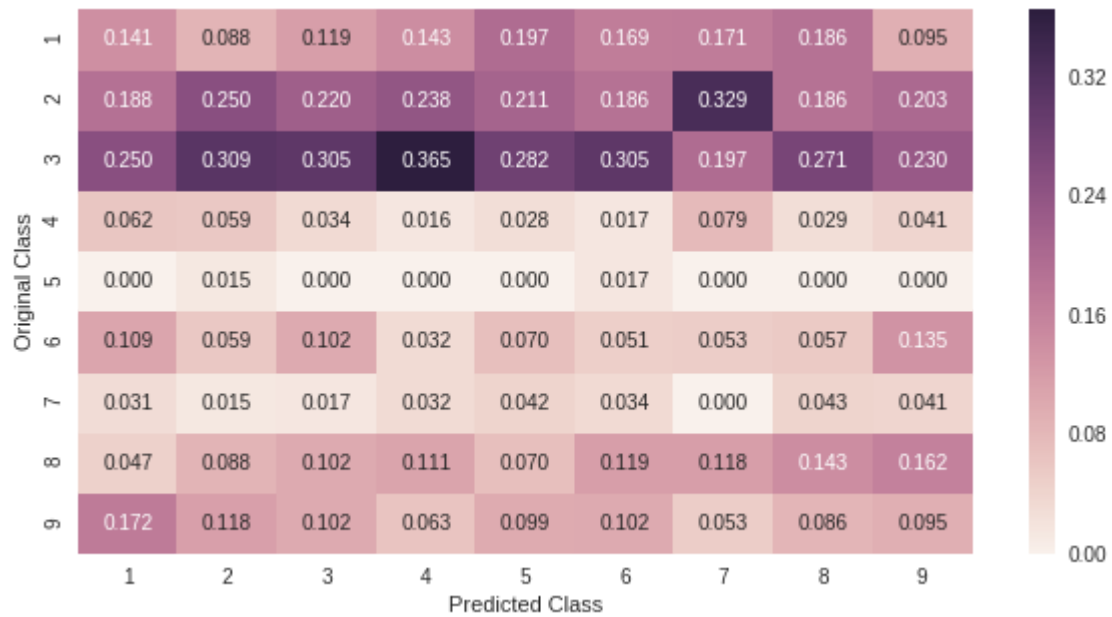
Number of misclassified points 89.23841059602648

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

<Figure size 1080x576 with 0 Axes>

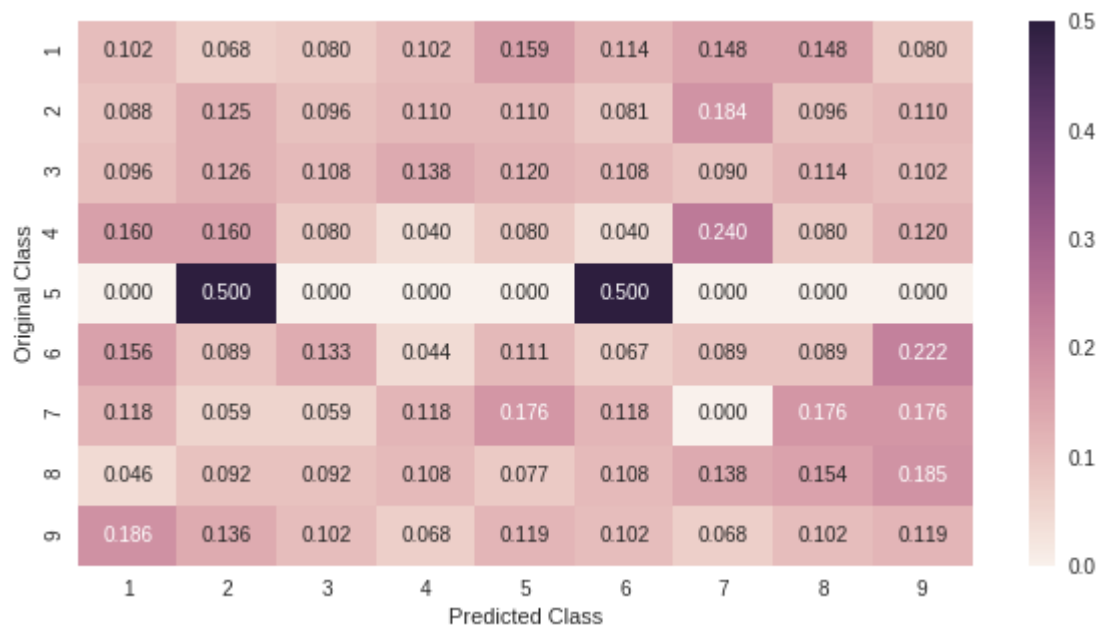


----- Precision matrix -----



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

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



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

4.1.2. K Nearest Neighbour Classification

In [0]:

```

alpha = [x for x in range(1, 10, 2)]
cv_log_error_array=[]
for i in alpha:
    print(i)
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cfl.classes_, ep

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

best_alpha = np.argmin(cv_log_error_array)

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

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

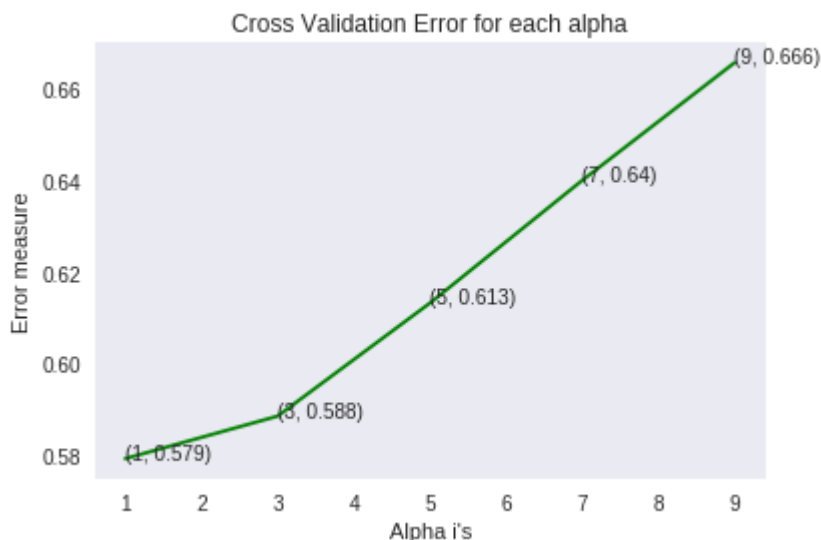
predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test)
print('For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

```

```

1
3
5
7
9
log_loss for k = 1 is 0.5791606031704255
log_loss for k = 3 is 0.5884541997326155
log_loss for k = 5 is 0.6131761414736544
log_loss for k = 7 is 0.6400133868683122
log_loss for k = 9 is 0.6657502168011453

```



For values of best alpha = 1 The train log loss is: 0.2901329071244141

For values of best alpha = 1 The cross validation log loss is: 0.5791606031704255

For values of best alpha = 1 The test log loss is: 0.7038217108964496

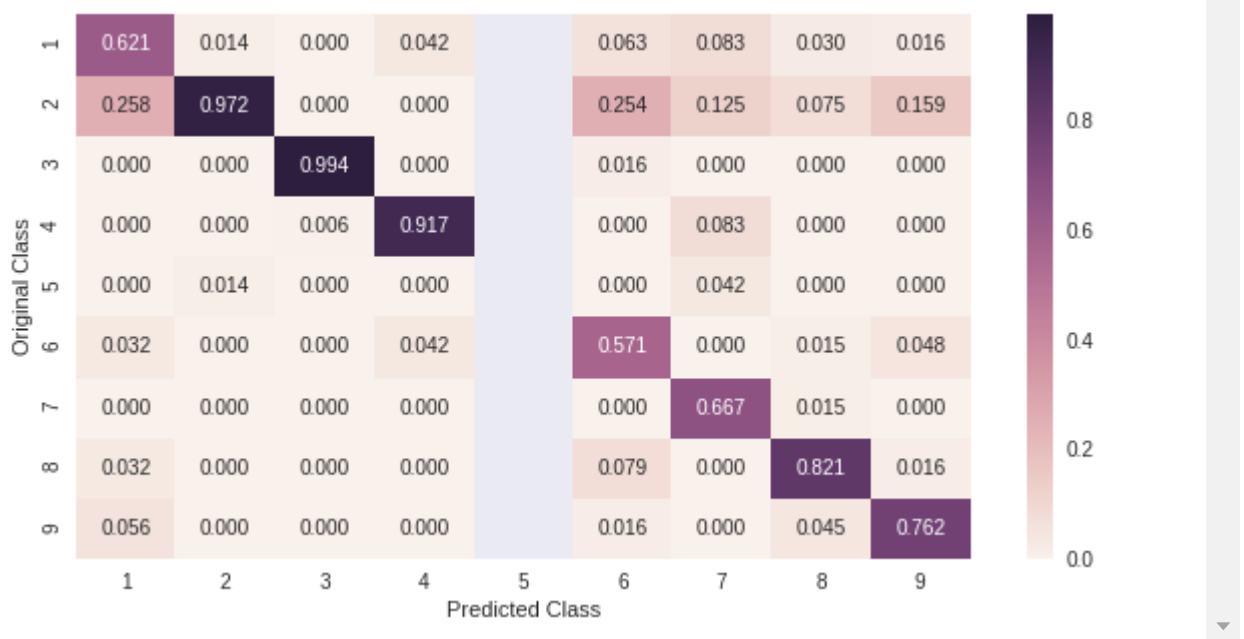
Number of misclassified points 18.874172185430464

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

<Figure size 1080x576 with 0 Axes>

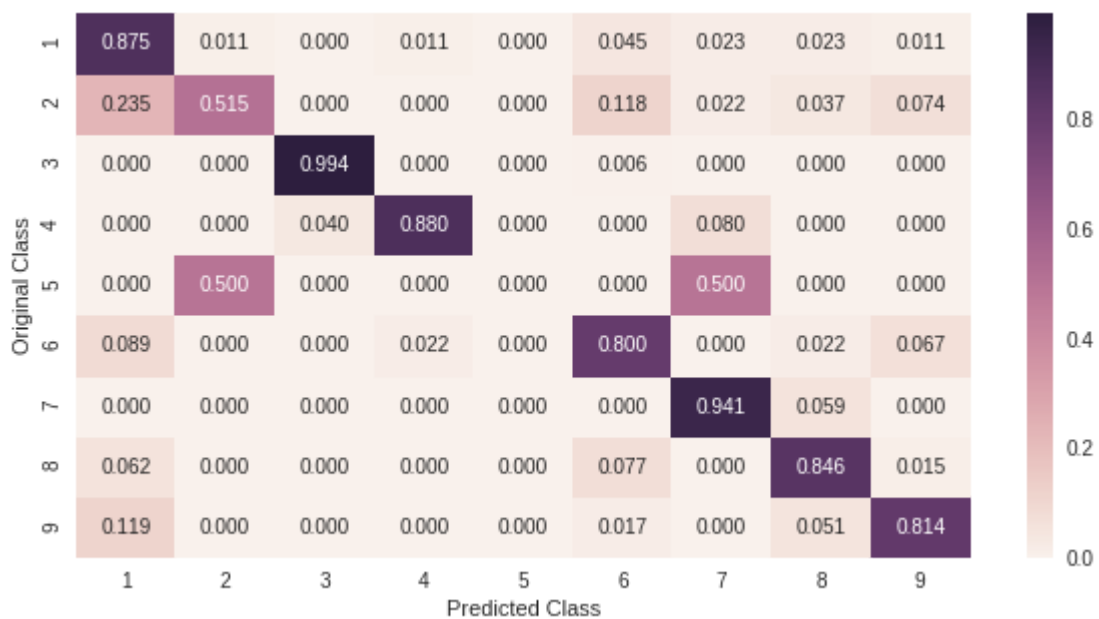


----- Precision matrix -----



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

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



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

4.1.3. Logistic Regression

In [0]:

```

alpha = [10 ** x for x in range(-2, 2)]
cv_log_error_array=[]
for i in alpha:
    print('current alpha value is',i)
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=logisticR.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

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

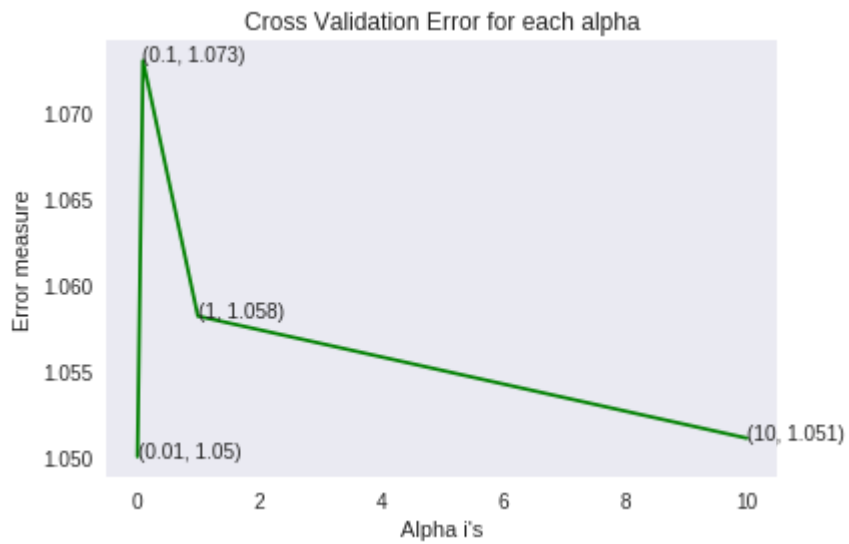
predict_y = sig_clf.predict_proba(X_train)
print ('log loss for train data',log_loss(y_train, predict_y, labels=logisticR.cl
predict_y = sig_clf.predict_proba(X_cv)
print ('log loss for cv data',log_loss(y_cv, predict_y, labels=logisticR.classes_
predict_y = sig_clf.predict_proba(X_test)
print ('log loss for test data',log_loss(y_test, predict_y, labels=logisticR.clas
plot_confusion_matrix(y_test, sig_clf.predict(X_test))

```

```

current alpha value is 0.01
current alpha value is 0.1
current alpha value is 1
current alpha value is 10
log_loss for c = 0.01 is 1.0499902467589572
log_loss for c = 0.1 is 1.0729569508613177
log_loss for c = 1 is 1.0581383830487323
log_loss for c = 10 is 1.051061617048442

```



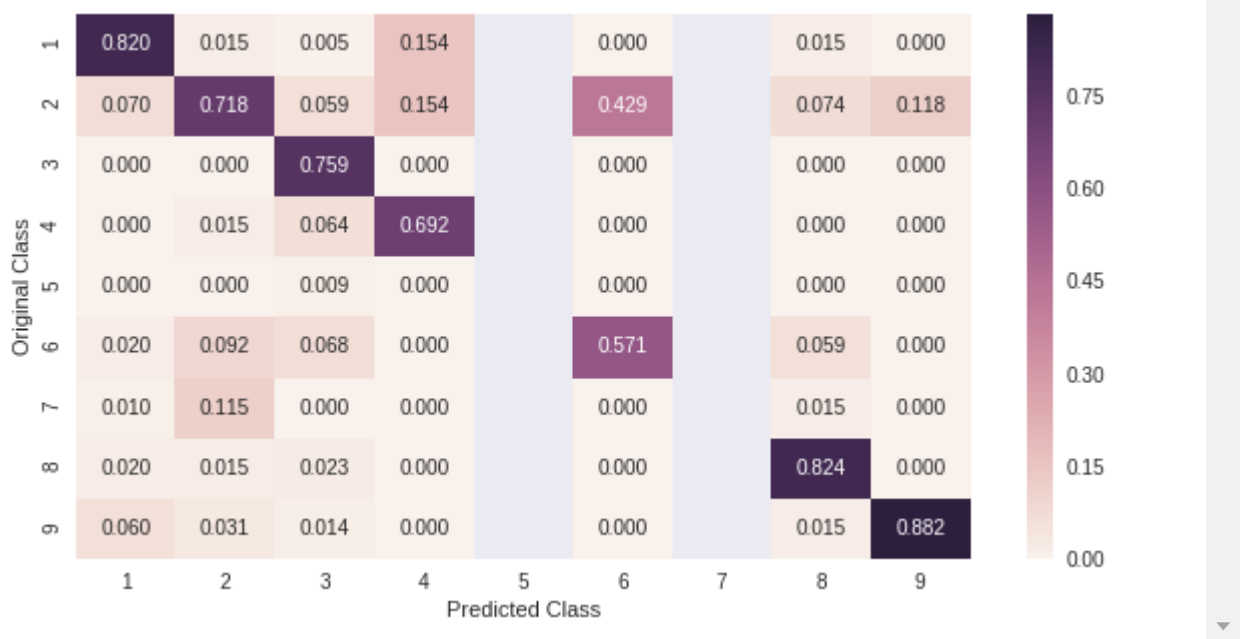
log loss for train data 0.9075279622711779
log loss for cv data 1.0499902467589572
log loss for test data 1.0407109799229066
Number of misclassified points 23.013245033112582

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

<Figure size 1080x576 with 0 Axes>

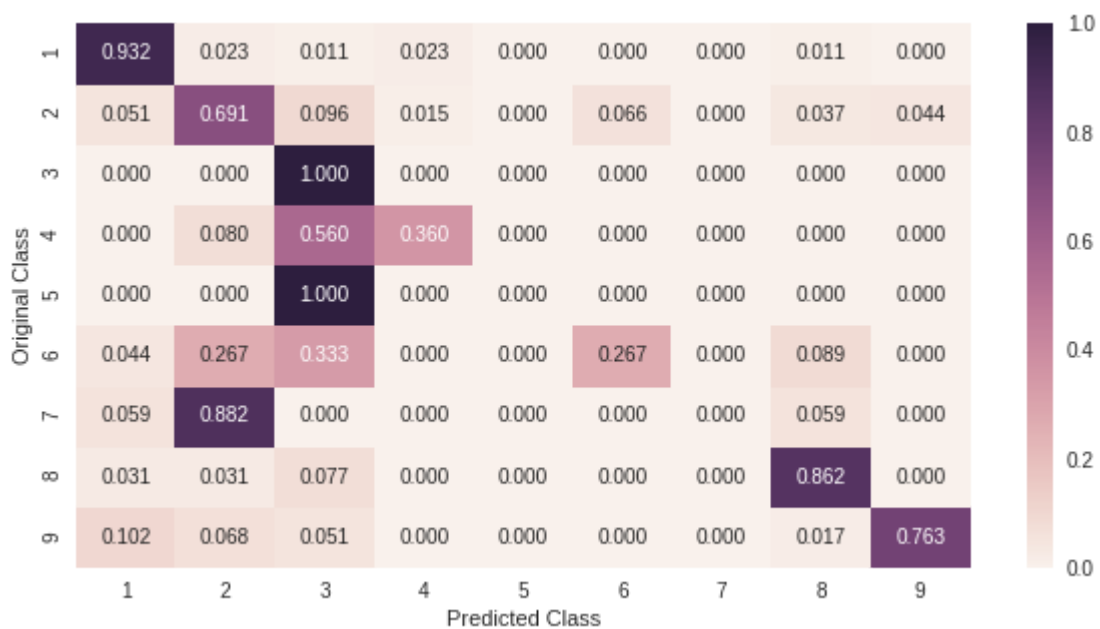


----- Precision matrix -----



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

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



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

4.1.4. Random Forest Classifier

In [0]:

```

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

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

best_alpha = np.argmin(cv_log_error_array)

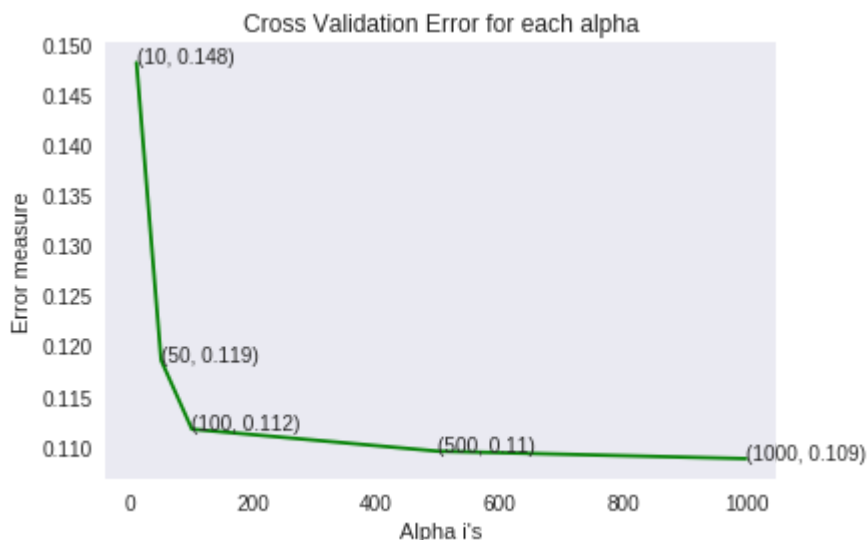
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

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

current alpha values is 10
current alpha values is 50
current alpha values is 100
current alpha values is 500
current alpha values is 1000
log_loss for c = 10 is 0.1481877887567859
log_loss for c = 50 is 0.11856598065687102
log_loss for c = 100 is 0.11170462724068957
log_loss for c = 500 is 0.1095091190751437
log_loss for c = 1000 is 0.1087707479310827

```



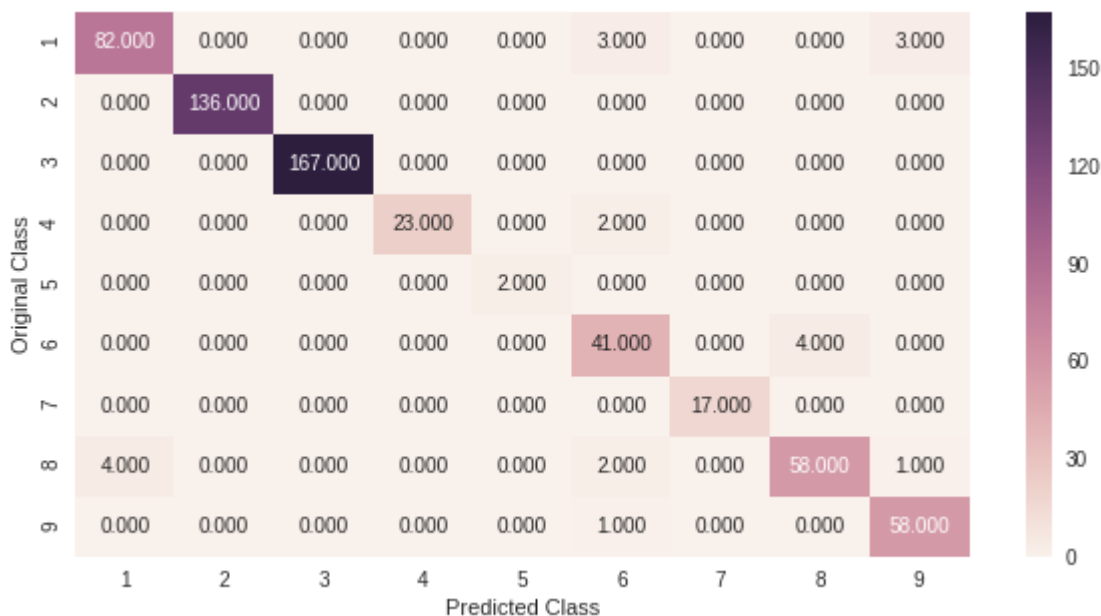
For values of best alpha = 1000 The train log loss is: 0.045470831665197164
 For values of best alpha = 1000 The cross validation log loss is: 0.1087707479310827

For values of best alpha = 1000 The test log loss is: 0.15214755949951503

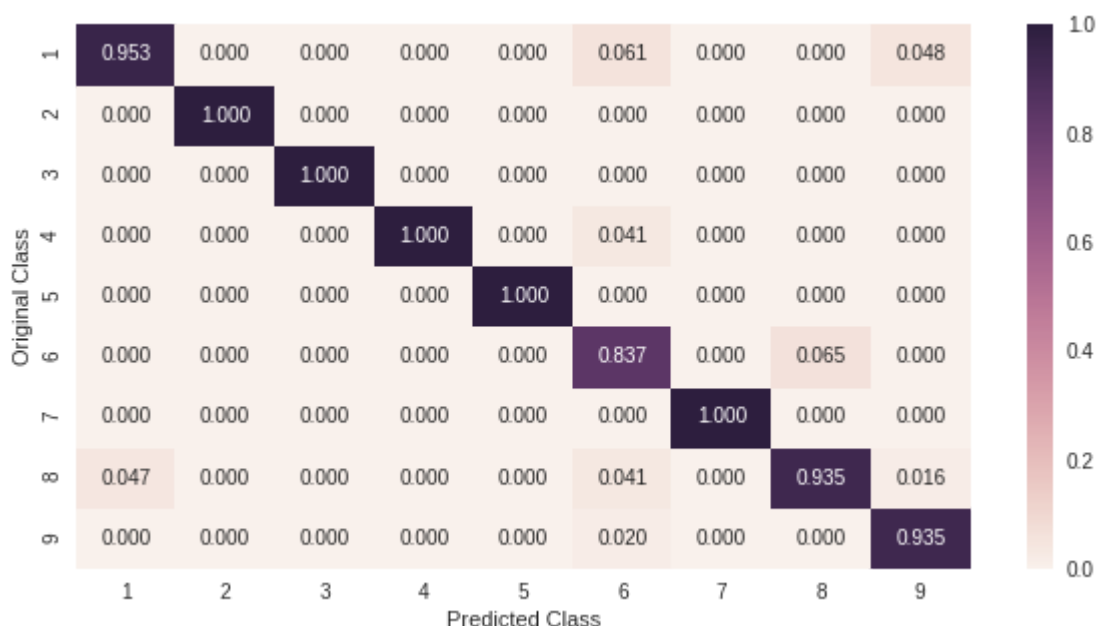
Number of misclassified points 3.3112582781456954

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

<Figure size 1080x576 with 0 Axes>



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



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

Recall matrix



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

SEESION IS GETTING CRASHED DUE TO THE LARGE DATA REDUCING THEDATA TO 20 PERCENT

```
In [0]: data_y = result['Class']
zz=result.drop(['ID','Class'], axis=1)
zz=zz.sample(frac=0.2, random_state=99,axis=1)
```

In [0]:

```
# split the data into test and train by maintaining same distribution of output variable
X_train, X_test, y_train, y_test = train_test_split(zz, data_y, stratify=data_y, test_size=0.2)
# split the train data into train and cross validation by maintaining same distribution
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify=y_train, test_size=0.2)
```

In [0]:

```
print(X_train.shape)
print(X_test.shape)
print(X_cv.shape)
print(y_train.shape)
print(y_test.shape)
print(y_cv.shape)
```

```
(1931, 13108)
(604, 13108)
(483, 13108)
(1931,)
(604,)
(483,)
```

+

4.1.5. XgBoost Classification

In [0]:

```

alpha=[10,500,1000]
cv_log_error_array=[]
for i in alpha:
    print('current alpha value is',i)
    x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
    x_cfl.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cfl.classes_, ep

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

best_alpha = np.argmin(cv_log_error_array)

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

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

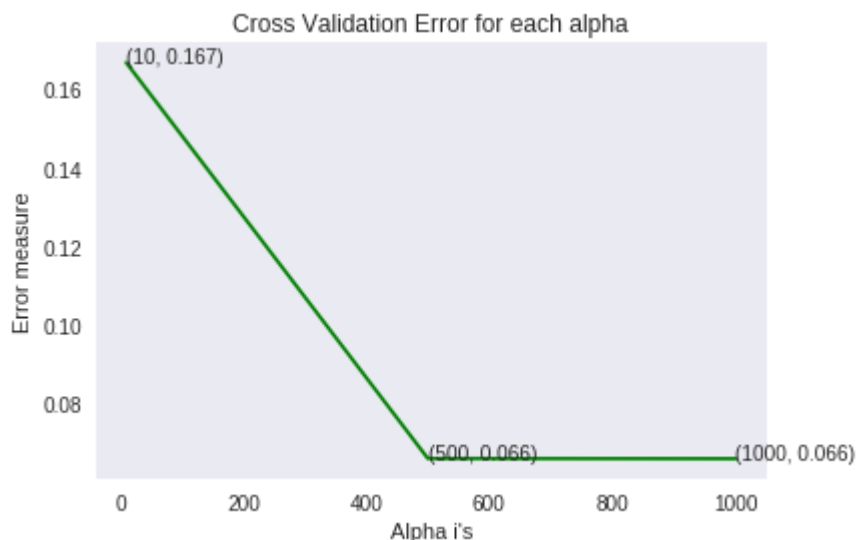
predict_y = sig_clf.predict_proba(X_train)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test)
print('For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

```

```

current alpha value is 10
current alpha value is 500
current alpha value is 1000
log_loss for c = 10 is 0.1668037527416192
log_loss for c = 500 is 0.06614198935543888
log_loss for c = 1000 is 0.06609387321850578

```

For values of best alpha = 1000 The train log loss is: 0.04081008058831673
 For values of best alpha = 1000 The cross validation log loss is: 0.06609387321850578

For values of best alpha = 1000 The test log loss is: 0.08880656976892347

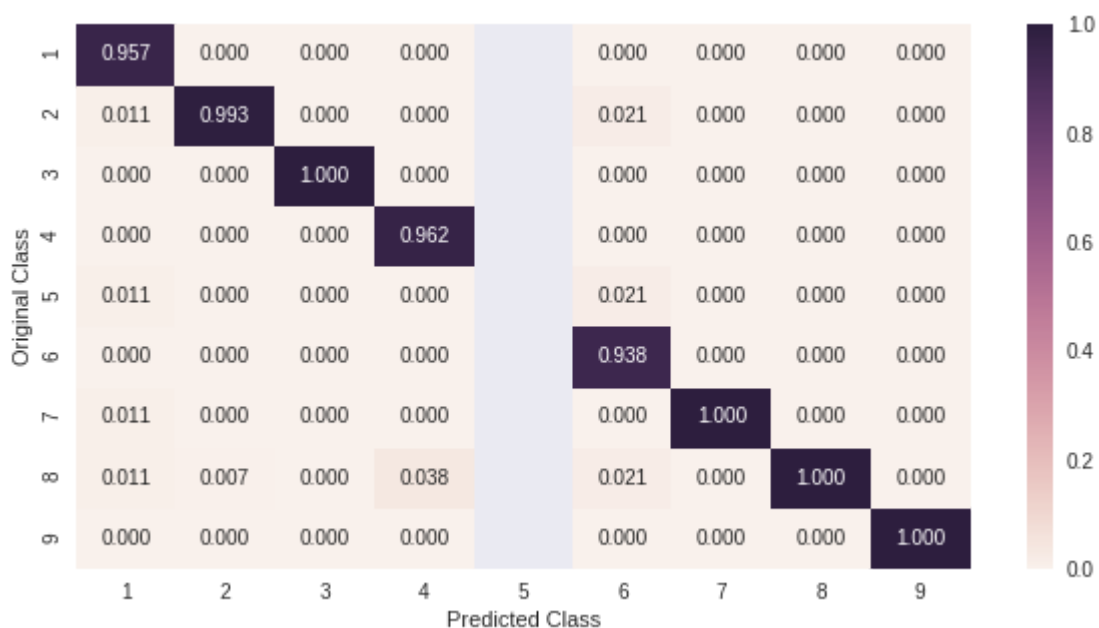
Number of misclassified points 1.490066225165563

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

<Figure size 1080x576 with 0 Axes>

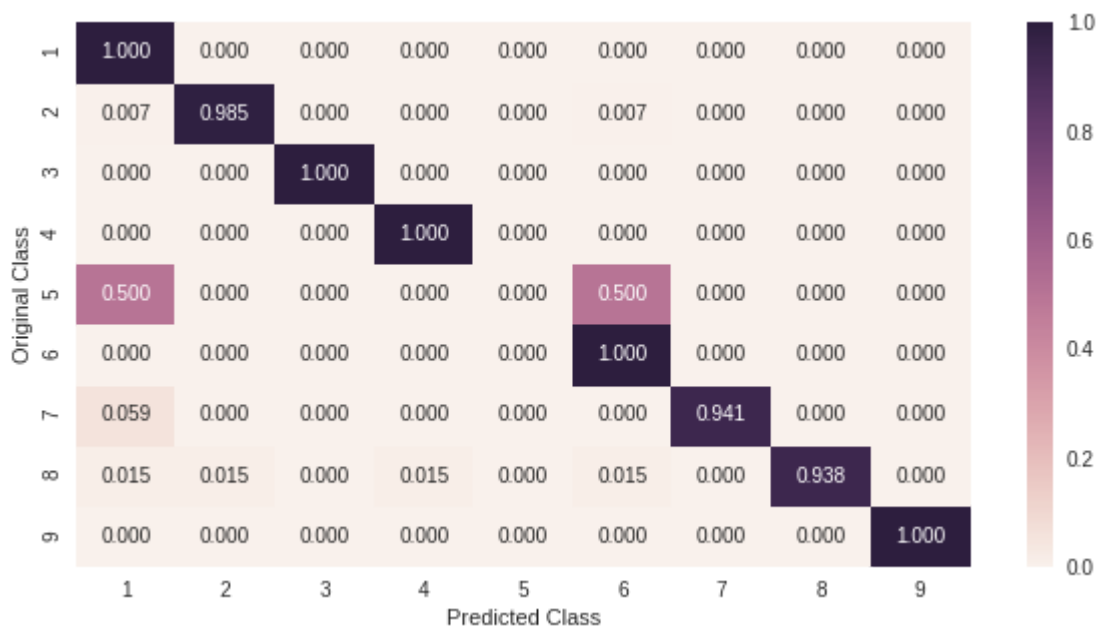


----- Precision matrix -----



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

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



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

4.1.5. XgBoost Classification with best hyper parameters using RandomSearch

```
In [0]: # https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xg
x_cfl=XGBClassifier()

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

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

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done   1 tasks      | elapsed: 39.0min
[Parallel(n_jobs=-1)]: Done   4 tasks      | elapsed: 93.5min
[Parallel(n_jobs=-1)]: Done   9 tasks      | elapsed: 147.8min
[Parallel(n_jobs=-1)]: Done  14 tasks      | elapsed: 208.6min
[Parallel(n_jobs=-1)]: Done  21 tasks      | elapsed: 426.0min
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed: 511.2min finished
```

```
Out[17]: RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
    estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_b
ylevel=1,
    colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
    max_depth=3, min_child_weight=1, missing=None, n_estimators=100,
    n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
    reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
    silent=True, subsample=1),
    fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
    param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15,
0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'co
lsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
    pre_dispatch='2*n_jobs', random_state=None, refit=True,
    return_train_score='warn', scoring=None, verbose=10)
```

```
In [0]: print (random_cfl1.best_params_)

{'subsample': 0.5, 'n_estimators': 2000, 'max_depth': 5, 'learning_rate': 0.2,
'colsample_bytree': 0.5}
```

```
In [0]: x_cfl=XGBClassifier(n_estimators=2000, learning_rate=0.2, colsample_bytree=0.5, m
x_cfl.fit(X_train,y_train)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train,y_train)

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

```
train loss 0.03551925222127126
cv loss 0.09008565321978028
test loss 0.09088136809655797
```

DOCUMENTATION, CONCLUSIONS AND KEYTAKEAWAYS FROM THE TASK1

```
In [0]: import pandas as pd
dta = [['k-nearest neighbors',1,0.29,0.57,0.7],['Logistic regression',0.01,0.9,1.0]
aa=pd.DataFrame(dta, columns=['model','BEST HYPER PARAMETER','TRAIN LOG LOSS','CV
aa
```

```
Out[4]:
```

	model	BEST HYPER PARAMETER	TRAIN LOG LOSS	CV LOG LOSS	TEST LOGLOSS
0	k-nearest neighbors	1	0.290	0.57	0.70
1	Logistic regression	0.01	0.900	1.04	1.04
2	RANDOM FOREST classifier	1000	0.045	0.10	0.15
3	XG boost with calibrated classifier cv	1000(number of estimators)	0.040	0.06	0.08
4	xgbOOST classifier WITH RANDOM SEARCH CV	2000(number of estimators),maximum depth=5,LEA...	0.030	0.09	0.09

the key take aways are due to 65536 columns instead of 256 columns the model is affected by the curse of dimensionality where we were able to achieve the 0.07 log loss in case of 256 columns we are able to achieve 0.08 log loss . the major cahnge occured in logistic regression model where due to more number of dimensions log loss is nearly 1.0 which is 0.7 in the previous case

##DOCUMENTATION

- we have the data of .byte files we have parse throught the every file and otained the bigram features of dimensions 65536.
- we have the features and aplitted the data into train,crossvalidataion and test
- we have applied models like

- knn
- logistic regression
- random forest
- xg boost with calibrated classifier
- ag boost with random search
- log loss compared to the 256 features it got degraded but still using the xgboost with calibrated classifier cv tuning the number of estimators we are able to achieve the log loss of 0.8 which is 0.7 in the previous case

#FROM HERE TASK 2 STARTS

Task 2 using the dchads github account analysis done by him and selecting the best features and models

key take aways from the dchads analysis is

- * select the best features and apply the extra trees classifier on top of it.**
- * we do this analysis selecting the 10percent best features of total features**
- * 20 percent of best features**
- * 30 percent of best features so on..**
- * we select the features based on the chi squared test performed on top of them.**
- * using the extra trees classifier we obtain the parameters of the model such as logloss.**
- * finally we select the percent of features which are best and apply models on them and predict the y labels.**

use the image features with the byte and .asm features and apply models on them like

- xgboost
- extra trees classifier
- logistic regression
- naive bayes
- knn

BASED ON THE RESULTS OF LOGLOSS AND CONFUSION MATRIX OBTAIN THE BEST MODEL.

FINALLY WE WILL GET THE OPIMUM FEATURES WHAICH ARE TO BE CONSIDERD AND THE MODEL TO BE EMPLOYED ON THE DATA.

for merging the .bytes features sharable link

<https://drive.google.com/open?id=1Gwb9wEDI4tKhSnOHVwTJyqVe1yTIDh3IEYUXBFdmB2c>
(<https://drive.google.com/open?id=1Gwb9wEDI4tKhSnOHVwTJyqVe1yTIDh3IEYUXBFdmB2c>)

#asm file sharable link https://drive.google.com/open?id=1XhfwsDj4G_R94oyW-QZJHAKfoilmKE1G
(https://drive.google.com/open?id=1XhfwsDj4G_R94oyW-QZJHAKfoilmKE1G)

as documented by dchad the .byte image features are weak learners.we are extracting the image features of the .asm files anfd teake as the features . we consider the 1000 pixel of data and append them with the .asm features and .byte features we apply model on top of them and visualise the results.

function below will parse through all the .asm files and generate the image features and converted into csv file

we will upload this file into drive and append the features and apply the models.

this is the code used to obtain the image featrues and stores in the form of csv file . later this csv file uploaded to drive and used featurise and train the model using the iamge features. as a note we have considered 0only 1000 pixel of data .

```
In [0]: from multiprocessing import Pool
import os
from csv import writer
import numpy as np
import math
import scipy.misc
import array
import time as tm
```

```
In [0]: def read_image(filename):
    #startTime = time.time()
    f=open(filename,'rb')
    ln=os.path.getsize(filename)
    width=int(ln**0.5)
    rem=ln%width
    a=array.array("B")
    a.fromfile(f,ln-rem)
    f.close()
    print(int(len(a)/width),width)
    print(type(a))
    g=np.reshape(a,(int(len(a)/width),width))
    print(g.shape)

    g=np.uint8(g)

    g=np.resize(g,(1000,))
    return list(g)
```

```
In [0]: def extract_byte_image_features(tfiles):
    asm_files = [i for i in tfiles if '.bytes' in i]
    ftot = len(asm_files)

    pid = os.getpid()
    print('Process id:', pid)
    feature_file = 'data/' + str(pid) + '-train-image-features-byte.csv'
    print('feature file:', feature_file)

    outrows = []
    with open(feature_file,'w') as f:
        fw = writer(f)
        column_names = ['filename'] + [("BYTE_{:s}".format(str(x))) for x in range(1000)]
        fw.writerow(column_names)
        for idx, fname in enumerate(asm_files):
            file_id = fname.split('.')[0]
            image_data = read_image('bbytefiles/' + fname)
            outrows.append([file_id] + image_data)
            print('yes')
            # Print progress
            if (idx+1) % 10 == 0:
                print(pid, idx + 1, 'of', ftot, 'files processed.')
                fw.writerow(outrows)
                outrows = []

        # Write remaining files
        if len(outrows) > 0:
            fw.writerow(outrows)
            outrows = []
```

```
In [0]: def extract_asm_image_features(tfiles):
asm_files = [i for i in tfiles if '.asm' in i]
ftot = len(asm_files)
count=0
pid = os.getpid()
print('Process id:', pid)
feature_file = 'data/' + str(pid) + '-test-image-featuresout-asm.csv'
print('feature file:', feature_file)

outrows = []
with open(feature_file, 'w') as f:
    fw = writer(f)
    column_names = ['filename'] + [("ASM_{:s}".format(str(x))) for x in range
fw.writerow(column_names)
    for idx, fname in enumerate(asm_files):
        file_id = fname.split('.')[0]
        image_data = read_image('basmFiles/' + fname)
        outrows.append([file_id] + image_data)
        print(count)
        count+=1
        # Print progress
        if (idx+1) % 10 == 0:
            print(pid, idx + 1, 'of', ftot, 'files processed.')
            fw.writerows(outrows)
            outrows = []

    # Write remaining files
    if len(outrows) > 0:
        fw.writerows(outrows)
        outrows = []
```

```
In [0]: # TRAIN FILES BYTE
# Now divide the train files into four groups for multiprocessing
start_time = tm.time()
ext_drive = 'basmFiles'
tfiles = os.listdir(ext_drive)
extract_asm_image_features(tfiles)
```

```
In [0]: # Code to read csv file into Colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)
```

```
100% |████████████████████████████████████████| 993kB 12.5MB/s ta 0:00:01
Building wheel for PyDrive (setup.py) ... done
```

```
In [0]: link = 'https://drive.google.com/open?id=10Nrz6CTKEG0oLmBsHy7sMHfL2v8zRjxb' # The
```



```
In [0]: fluff, id = link.split('=')
print (id) # Verify that you have everything after '='
```

```
10NrZ6CTKEG0oLmBsHy7sMHfL2v8zRjxb
```

```
In [0]: import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('result_with_size.csv')
dataframe = pd.read_csv('result_with_size.csv')
```

```
In [0]: dataframe.head(5)
```

```
Out[6]:
```

	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	1
0	0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	...	310
1	1	01IsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	...	43
2	2	01jsnpXSAlgW6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	...	224
3	3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	...	48
4	4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	...	35

5 rows × 261 columns



```
In [0]: print(dataframe.shape)

(10868, 261)
```

we got the byte features file with the class laebl we also get the .asm feature file and appen them so that we can start the process and apply chisqaure tests on the these features.

```
In [0]: # Code to read csv file into Colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)
```

```
In [0]: link = 'https://drive.google.com/open?id=1XhfwsDj4G_R94oyW-QZJHAKfoilmKE1G' # The
```

```
In [0]: fluff, id = link.split('=')
print (id) # Verify that you have everything after '='

1XhfwsDj4G_R94oyW-QZJHAKfoilmKE1G
```

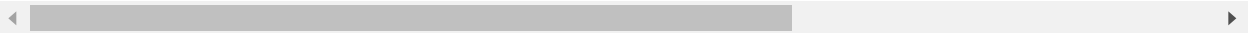
```
In [0]: import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('asmoutputfile.csv')
dat= pd.read_csv('asmoutputfile.csv')
```

```
In [0]: dat.head(5)
```

```
Out[12]:
```

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

5 rows × 52 columns



```
In [0]: print(dat.shape)
```

(10868, 52)

```
In [0]: dataframe.head(5)
```

```
Out[14]:
```

	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	310
0	0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	...	310
1	1	01IsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	...	43
2	2	01jsnpXSAIgw6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	...	224
3	3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	...	48
4	4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	...	35

5 rows × 261 columns



```
In [0]: result=pd.merge(dataframe,dat,on='ID',how='left')
result.head(5)
print(result.shape)
```

(10868, 312)

```
In [0]: result.head(5)
```

```
Out[16]:
```

	Unnamed: 0	ID	0	1	2	3	4	5	6	7	...	:dw
0	0	01azqd4lnC7m9JpocGv5	601905	3905	2816	3832	3345	3242	3650	3201	...	4
1	1	01IsoiSMh5gxyDYTI4CB	39755	8337	7249	7186	8663	6844	8420	7589	...	1
2	2	01jsnpXSAlgW6aPeDxrU	93506	9542	2568	2438	8925	9330	9007	2342	...	
3	3	01kcPWA9K2BOxQeS5Rju	21091	1213	726	817	1257	625	550	523	...	
4	4	01SuzwMJEIXsK7A8dQbl	19764	710	302	433	559	410	262	249	...	1

5 rows × 312 columns



```
In [0]: asm_y = result['Class']
asm_x = result.drop(['ID', 'Class', '.BSS:', 'rtn', '.CODE'], axis=1)
```

we have taken merged the features of asm files and the byte feature files

```
In [0]: print(asm_x.shape)
print(asm_y.shape)
```

```
(10868, 307)
(10868,)
```

```
In [0]: from sklearn.feature_selection import chi2
x_,y_=chi2(asm_x,asm_y)
```

```
In [0]: print(x_)
        print(y_)
```

```
[1.14341795e+05 3.96296388e+08 8.29424914e+07 2.11866508e+08
1.21518834e+08 1.56646103e+08 3.77045925e+07 1.50660934e+08
3.67865391e+07 3.98035682e+07 1.85680256e+07 2.05170518e+07
5.09745245e+07 3.74615486e+07 1.36008327e+07 1.64122105e+07
4.58867963e+07 1.19408013e+08 9.10233722e+07 1.19800971e+08
8.91709202e+07 3.85225315e+07 1.95004632e+07 2.61291230e+07
2.49322500e+07 4.51176763e+07 3.14994291e+07 1.76055421e+07
2.32777739e+07 2.61497954e+07 1.64436549e+07 2.40118909e+07
2.72054860e+07 8.72142804e+07 6.28374499e+07 6.88625460e+07
6.90361203e+07 2.12679337e+07 1.02287734e+07 9.25638142e+06
1.17846363e+07 9.55342397e+06 8.09213410e+06 8.56765717e+06
9.13113731e+06 1.30361363e+07 8.85251561e+06 8.74861499e+06
1.01578612e+07 9.64213901e+07 6.39678473e+07 7.63267913e+07
9.63187771e+07 6.82188948e+06 1.07745550e+07 9.76914967e+06
1.29361868e+07 1.52510927e+07 2.65284624e+07 7.55811432e+06
2.82021999e+07 1.25944700e+07 1.23870312e+07 8.60819574e+06
7.80445372e+06 2.29188963e+07 4.80024614e+06 1.07449512e+07
2.61557473e+07 1.94970061e+07 2.45290898e+07 2.80350716e+07
8.35130304e+06 2.39797988e+07 8.26110962e+06 2.21932538e+07
1.55096393e+07 6.82202812e+06 2.29193601e+07 1.17169544e+07
8.24101285e+06 2.10933456e+07 9.56477686e+06 9.40836477e+06
1.99271126e+07 8.01192604e+06 2.53034958e+07 2.67434246e+07
2.33396797e+07 1.08146024e+07 2.89749941e+07 8.74032658e+06
7.29187784e+06 8.13167262e+06 2.55645589e+07 2.81805872e+07
1.32293593e+07 6.53701558e+06 5.66472254e+06 1.29099366e+07
8.85836725e+06 9.58186054e+06 8.37632523e+06 6.64552302e+06
8.69360608e+06 2.53642176e+07 1.51844548e+07 2.43651561e+07
1.28224436e+07 1.04244086e+07 1.33841188e+07 7.74185089e+06
1.58272859e+07 7.68515776e+06 1.28486641e+07 1.12068929e+07
1.35986565e+07 1.91236330e+07 2.89523067e+07 1.09235663e+07
8.81238786e+06 8.64597076e+06 1.43645832e+07 1.14435136e+07
8.98319404e+06 7.15000999e+06 1.83905894e+07 8.05189187e+06
8.18685505e+06 1.67937230e+07 7.06479915e+06 7.81887838e+06
1.77267055e+07 1.68795689e+07 1.75058522e+07 9.72172825e+06
9.33386677e+06 1.79811326e+07 1.73483254e+07 8.83821988e+06
4.09829863e+07 7.96613660e+06 1.55432523e+07 8.70380552e+06
8.81554472e+06 1.11903584e+07 9.42852734e+06 8.80229330e+06
9.76317777e+06 8.01558554e+06 9.00134829e+06 9.11655667e+06
1.07158572e+07 8.84707967e+06 8.51767052e+06 9.74309002e+06
9.15427575e+06 9.15269961e+06 1.01901481e+07 8.91515035e+06
1.04389954e+07 7.39670301e+06 7.28951423e+06 1.04178388e+07
8.23601320e+06 8.54827896e+06 8.88133834e+06 9.93605278e+06
8.98101185e+06 8.26562167e+06 9.08747509e+06 8.95291399e+06
8.89318604e+06 8.92941778e+06 1.02959482e+07 9.34119070e+06
1.01727726e+07 8.51030036e+06 1.00736733e+07 8.90845629e+06
1.01217630e+07 8.47033925e+06 1.05480795e+07 1.82001170e+07
8.73714480e+06 8.14548029e+06 8.73526260e+06 9.17380861e+06
9.05486517e+06 8.80747396e+06 9.18908094e+06 1.02625853e+07
9.24478293e+06 1.68038410e+07 1.98113330e+07 6.62864929e+06
1.96578210e+07 1.75221891e+07 9.06427107e+06 1.74496523e+07
1.88674859e+07 7.51646764e+06 7.76110137e+06 8.46934065e+06
9.13774947e+06 2.74379005e+08 9.15097392e+06 8.43286212e+06
8.97877077e+06 7.92591348e+06 8.49485473e+06 8.13098386e+06
8.95734670e+06 8.06514164e+06 8.46204012e+06 1.30251137e+07]
```

45/91

```
In [0]: import numpy as np
print(type(y_))
a=np.argsort(y_)[::-1][:32]
print(a)
```

```
<class 'numpy.ndarray'>
[306 104  97  98  99 100 101 102 103 105  76 106 107 108 109 110 111 112
  96  95  94  93  78  79  80  81  82  83  84  85  86  87]
```

```
In [0]: #as suggested by dchad we gonna consider the models performance on 10 percent
#20 oercent and 30 percent of the total features
#for the model performance we gonna consider logloss,accuracy,confusion matrix
#we gon na see the model performance using the extra trees classifier
```

```
In [0]: traindatafor10percentfeatures=asm_x.iloc[:,a]
```

```
In [0]: traindatafor10percentfeatures.head(5)
```

```
Out[42]:
```

	eip	67	60	61	62	63	64	65	66	68	...	4d	4e	4f	50	
0	456	4375	3580	3344	3481	2945	3300	3447	2851	3730	...	2774	2732	3161	3955	45
1	227	801	711	952	733	648	1088	1591	1109	11323	...	8913	465	65918	9127	5
2	117	2883	2601	2553	2843	3055	2755	2712	2403	2748	...	2443	2464	2584	16251	27
3	29	499	545	814	399	538	1092	1041	941	1129	...	602	423	420	1108	5
4	76	365	1512	2134	484	585	368	568	347	385	...	845	290	389	562	4

5 rows × 32 columns

```
In [0]: from sklearn.model_selection import train_test_split
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(traindatafor1
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_tr
```

```

In [0]: import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
from sklearn.calibration import CalibratedClassifierCV
from sklearn.metrics import log_loss
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
from sklearn.ensemble import ExtraTreesClassifier
for i in alpha:
    print('current alpha value is ',i)
    r_cfl=ExtraTreesClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

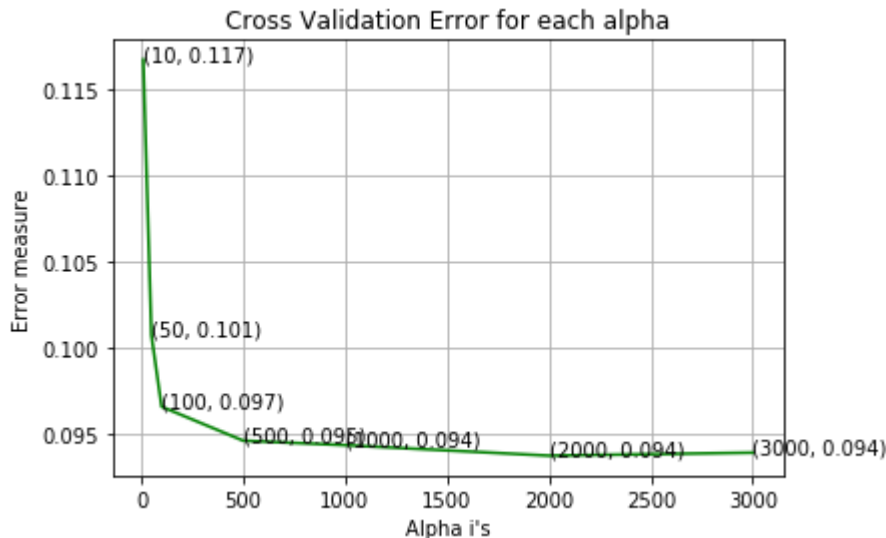
r_cfl=ExtraTreesClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1

current alpha value is 10
current alpha value is 50
current alpha value is 100
current alpha value is 500
current alpha value is 1000
current alpha value is 2000
current alpha value is 3000
log_loss for c = 10 is 0.11673606524838116
log_loss for c = 50 is 0.10075306929177441
log_loss for c = 100 is 0.09659916231276144
log_loss for c = 500 is 0.09462803478471248

```

log_loss for c = 1000 is 0.09436658540263229
 log_loss for c = 2000 is 0.093757295762736
 log_loss for c = 3000 is 0.09392943133367676



For values of best alpha = 2000 The train log loss is: 0.03354420236772946
 For values of best alpha = 2000 The cross validation log loss is: 0.093757295762736
 For values of best alpha = 2000 The test log loss is: 0.11686961272569582

considerinf the 20 percent features of total features

```
In [0]: import numpy as np
print(type(y_))
a=np.argsort(y_)[:-1][:64]
print(a)
```

```
<class 'numpy.ndarray'>
[306 104  97  98  99 100 101 102 103 105  76 106 107 108 109 110 111 112
  96  95  94  93  78  79  80  81  82  83  84  85  86  87  88  89  90  91
  92 113 114 115 134 136 137 138 139 140 141 142 143 144 145 146 147 148
 149 150 135 133 116 132 117 118 119 120]
```

```
In [0]: #as suggested by dchad we gonna consider the models performance on 10 percent
#20 oercent and 30 percent of the total features
#for the model performance we gonna consider logloss,accuracy,confusion matrix
#we gon na see the model performance using the extra trees classifier
```

```
In [0]: traindatafor10percentfeatures=asm_x.iloc[:,a]
```



```
In [0]: traindatafor10percentfeatures.head(5)
```

```
Out[30]:
```

	eip	67	60	61	62	63	64	65	66	68	...	94	95	86	84	7
0	456	4375	3580	3344	3481	2945	3300	3447	2851	3730	...	3920	2952	2796	3290	373
1	227	801	711	952	733	648	1088	1591	1109	11323	...	781	977	755	7153	6583
2	117	2883	2601	2553	2843	3055	2755	2712	2403	2748	...	2862	2510	2516	2425	263
3	29	499	545	814	399	538	1092	1041	941	1129	...	408	377	376	746	72
4	76	365	1512	2134	484	585	368	568	347	385	...	253	277	231	269	40

5 rows × 64 columns

```
In [0]: from sklearn.model_selection import train_test_split
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(traindatafor10percentfeatures, y_train_asm,
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_asm,
```

```

In [0]: import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
from sklearn.calibration import CalibratedClassifierCV
from sklearn.metrics import log_loss
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
from sklearn.ensemble import ExtraTreesClassifier
for i in alpha:
    print('current alpha value is ',i)
    r_cfl=ExtraTreesClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

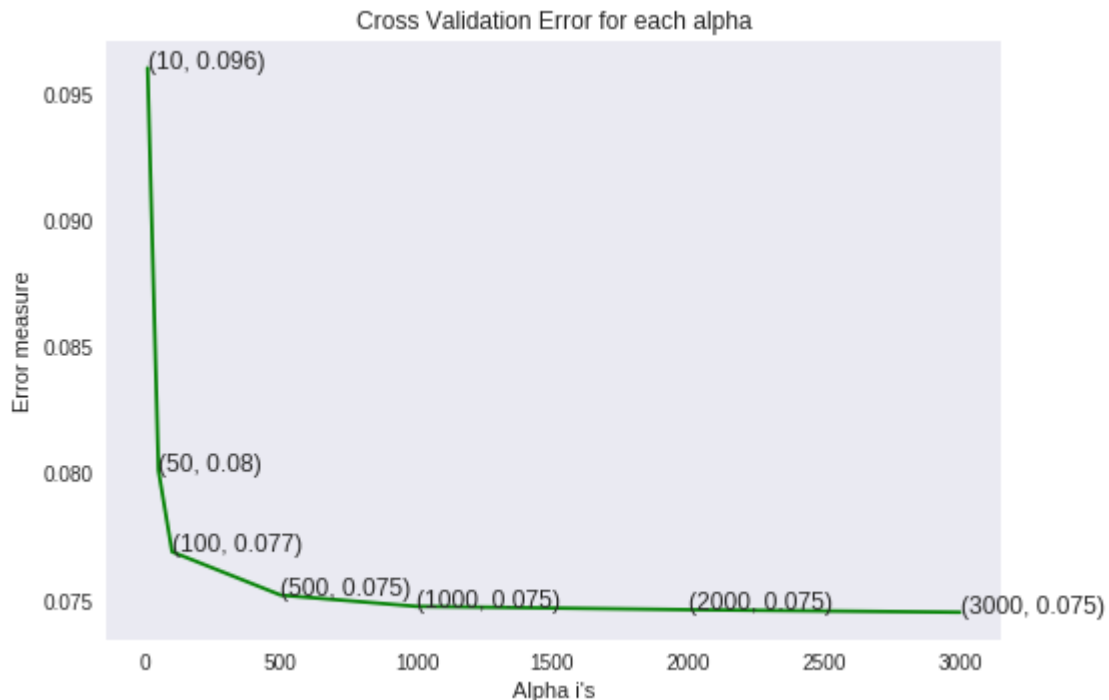
r_cfl=ExtraTreesClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1

current alpha value is 10
current alpha value is 50
current alpha value is 100
current alpha value is 500
current alpha value is 1000
current alpha value is 2000
current alpha value is 3000
log_loss for c = 10 is 0.09604268783702838
log_loss for c = 50 is 0.0801253872753246
log_loss for c = 100 is 0.07691579590699824
log_loss for c = 500 is 0.07520384984483364

```

log_loss for c = 1000 is 0.07475897100917324
 log_loss for c = 2000 is 0.0746270519739272
 log_loss for c = 3000 is 0.07452613334498641



For values of best alpha = 3000 The train log loss is: 0.02898564815591404
 For values of best alpha = 3000 The cross validation log loss is: 0.07452613334498641
 For values of best alpha = 3000 The test log loss is: 0.09296194748361782

```
In [0]: import numpy as np
print(type(y_))
a=np.argsort(y_)[::-1][:100]
print(a)
```

```
<class 'numpy.ndarray'>
[306 104  97  98  99 100 101 102 103 105  76 106 107 108 109 110 111 112
  96  95  94  93  78  79  80  81  82  83  84  85  86  87  88  89  90  91
  92 113 114 115 134 136 137 138 139 140 141 142 143 144 145 146 147 148
 149 150 135 133 116 132 117 118 119 120 121 122 123 124 125 126 127 128
 129 130 131  77  75 152  27  20  21  22  23  24  25  26  28  74  29  30
  31  32  33  34  35  19  18  17  16   1]
```

```
In [0]: #as suggested by dchad we gonna consider the models performance on 10 percent
#20 oercent and 30 percent of the total features
#for the model performance we gonna consider logloss,accuracy,confusion matrix
#we gon na see the model performance using the extra trees classifier
```

considering the 30 percent of the total features

```
In [0]: traindatafor10percentfeatures=asm_x.iloc[:,a]
```

```
In [0]: traindatafor10percentfeatures.head(5)
```

```
Out[47]:
```

	Unnamed: 0	ebx	7	64
0	0	587	3201	3300
1	1	905	7589	1088
2	2	451	2342	2755
3	3	43	523	1092
4	4	1546	249	368

```
In [0]: from sklearn.model_selection import train_test_split  
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(traindatafor10percentfeatures, y_train_asm, y_test_asm)  
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_asm, y_cv_asm, y_test_asm)
```

```

In [0]: import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
from sklearn.calibration import CalibratedClassifierCV
from sklearn.metrics import log_loss
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
from sklearn.ensemble import ExtraTreesClassifier
for i in alpha:
    print('current alpha value is ',i)
    r_cfl=ExtraTreesClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

r_cfl=ExtraTreesClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1

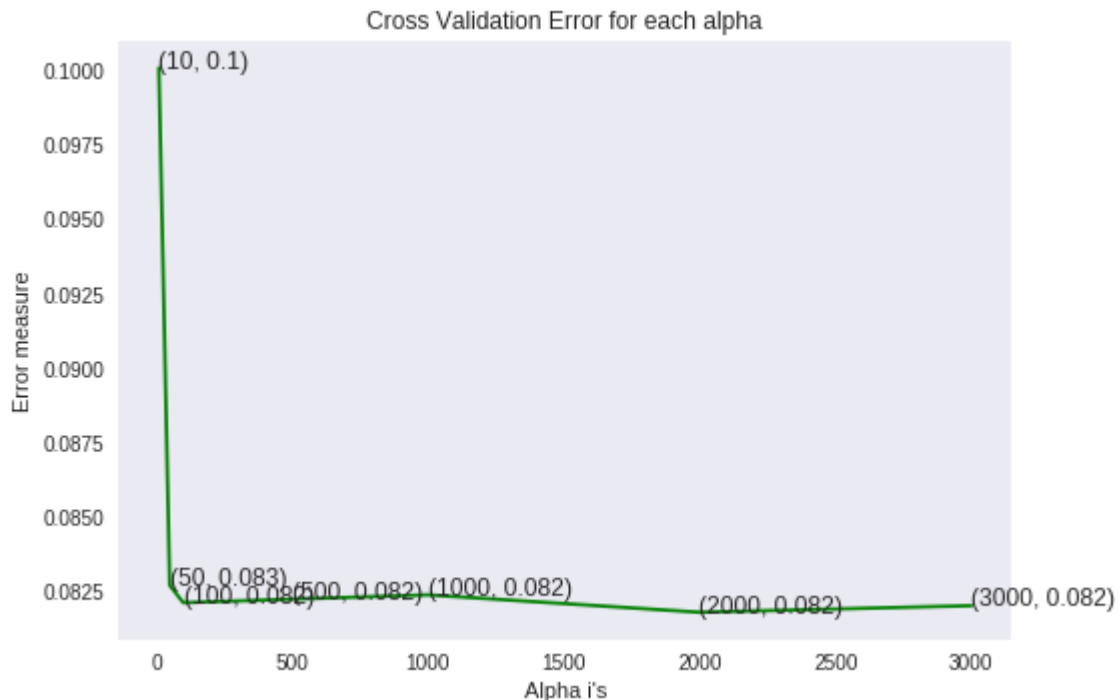
```

```

current alpha value is  10
current alpha value is  50
current alpha value is  100
current alpha value is  500
current alpha value is  1000
current alpha value is  2000
current alpha value is  3000
log_loss for c =  10 is 0.10007166534501818
log_loss for c =  50 is 0.08272449199994654
log_loss for c =  100 is 0.08214194884276566
log_loss for c =  500 is 0.08226022793142136

```

log_loss for c = 1000 is 0.08240476320614683
 log_loss for c = 2000 is 0.08182134329564608
 log_loss for c = 3000 is 0.0820416198197262



For values of best alpha = 2000 The train log loss is: 0.02871133773648726
 For values of best alpha = 2000 The cross validation log loss is: 0.08182134329564608
 For values of best alpha = 2000 The test log loss is: 0.08349945842357324

considering the all the features of the model and checking the performance metrics

```
In [0]: import numpy as np
print(type(y_))
a=np.argsort(y_)[:100]
print(a)
```

```
<class 'numpy.ndarray'>
[ 0 301  8 101]
```

```
In [0]: #as suggested by dchad we gonna consider the models performance on 10 percent
#20 oercent and 30 percent of the total features
#for the model performance we gonna consider logloss,accuracy,confusion matrix
#we gon na see the model performance using the extra trees classifier
```

```
In [0]: traindatafor10percentfeatures=asm_x.iloc[:,a]
```

```
In [0]: traindatafor10percentfeatures.head(5)
```

```
Out[40]:
```

	Unnamed: 0	cf	ce	cd	cc	cb	ca	c9	c8	c7	...	68	66	65	
0	0	3401	3347	3068	3105	2607	3103	3753	3222	4039	...	3730	2851	3447	33
1	1	613	607	480	1359	520	778	1127	662	7169	...	11323	1109	1591	10
2	2	2588	2796	2457	2641	174342	2720	2838	2685	2719	...	2748	2403	2712	27
3	3	344	363	345	433	359	370	566	609	531	...	1129	941	1041	10
4	4	188	245	237	285	228	230	212	238	375	...	385	347	568	3

5 rows × 307 columns

```
In [0]: from sklearn.model_selection import train_test_split
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(traindatafor10percentfeatures, y_train_asm,
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_train_asm,
```

```

In [0]: import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
from sklearn.calibration import CalibratedClassifierCV
from sklearn.metrics import log_loss
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
from sklearn.ensemble import ExtraTreesClassifier
for i in alpha:
    print('current alpha value is ',i)
    r_cfl=ExtraTreesClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

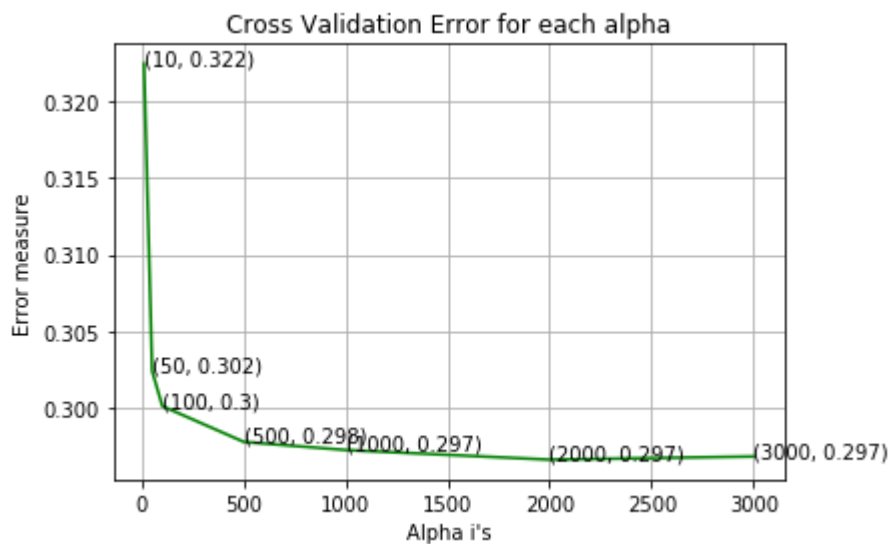
r_cfl=ExtraTreesClassifier(n_estimators=alpha[best_alpha],random_state=42,n_jobs=
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1

current alpha value is 10
current alpha value is 50
current alpha value is 100
current alpha value is 500
current alpha value is 1000
current alpha value is 2000
current alpha value is 3000
log_loss for c = 10 is 0.322454305064343
log_loss for c = 50 is 0.3023925443469619
log_loss for c = 100 is 0.3001046106544286
log_loss for c = 500 is 0.2977680388248046

```


log_loss for c = 1000 is 0.2972394059439264
log_loss for c = 2000 is 0.2966102882870639
log_loss for c = 3000 is 0.29681675133522983



For values of best alpha = 2000 The train log loss is: 0.08567315215010517
For values of best alpha = 2000 The cross validation log loss is: 0.2966102882870639
For values of best alpha = 2000 The test log loss is: 0.2873658900334469

```

In [0]: # This function plots the confusion matrices given y_i, y_i_hat.
%matplotlib inline
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    print("Number of misclassified points ",(len(test_y)-np.trace(C))/len(test_y))
    # C = 9,9 matrix, each cell (i,j) represents number of points of class i are j

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

    B = (C/C.sum(axis=0))

    plt.figure(figsize=(15,8))
    labels = [1,2,3,4,5,6,7,8,9]
    cmap=sns.light_palette("green")
    # representing A in heatmap format
    print("-"*50, "Confusion matrix", "-"*50)
    plt.figure(figsize=(10,5))
    sns.heatmap(C, annot=True, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.show()

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

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

```

using knn

```
In [0]: import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use('nbAgg')
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process# this is used for multithreading
import multiprocessing
import codecs# this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

```

In [0]: alpha = [x for x in range(1, 21,2)]
cv_log_error_array=[]
for i in alpha:
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=k_cfl.classes_)

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

best_alpha = np.argmin(cv_log_error_array)

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

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

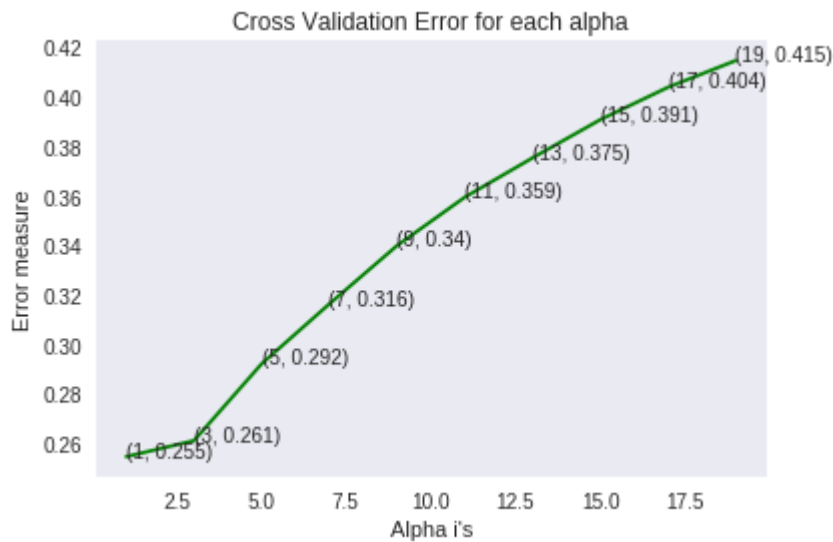
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',log_loss(y_train_asm, predict_y))
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',log_loss(y_cv_asm, predict_y))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',log_loss(y_test_asm,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

```

```

log_loss for k = 1 is 0.25464462500299606
log_loss for k = 3 is 0.2611059232576465
log_loss for k = 5 is 0.2921579026264672
log_loss for k = 7 is 0.31635093630130773
log_loss for k = 9 is 0.3399155429599249
log_loss for k = 11 is 0.3594522901540095
log_loss for k = 13 is 0.37533661571554444
log_loss for k = 15 is 0.39076336659271643
log_loss for k = 17 is 0.4038524146779985
log_loss for k = 19 is 0.414643193359781

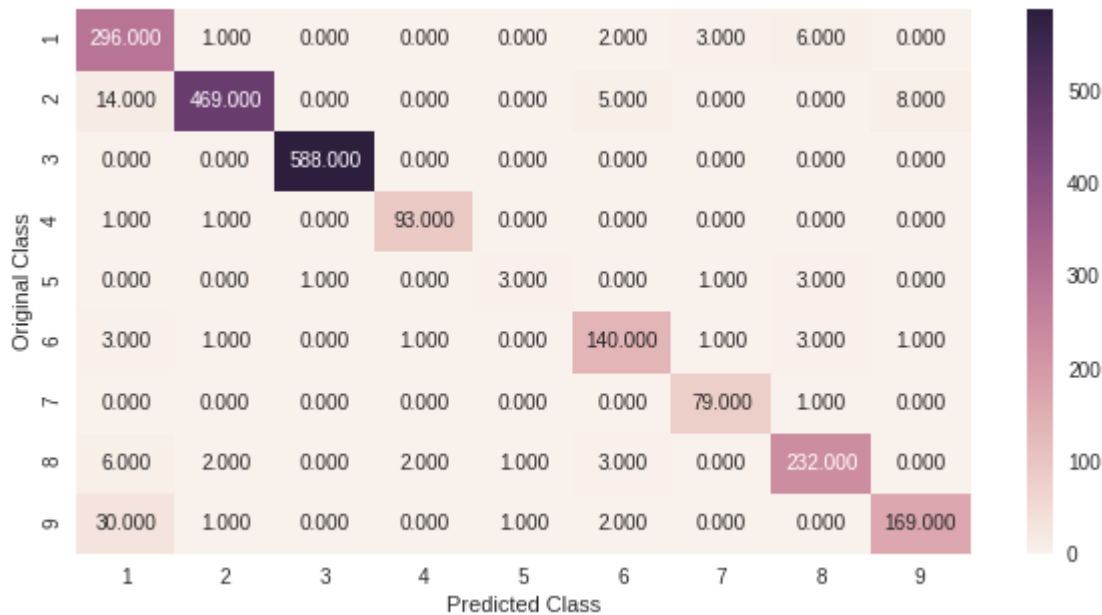
```



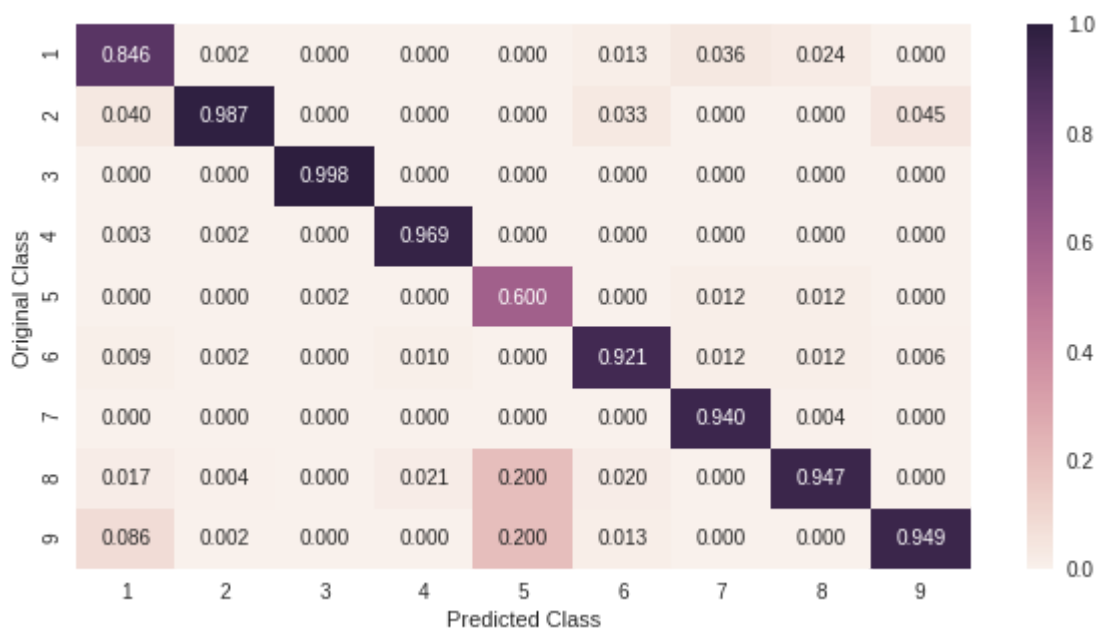
log loss for train data 0.08952853188535019
log loss for cv data 0.25464462500299606
log loss for test data 0.24004561397708524
Number of misclassified points 4.8298068077276906

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

<Figure size 1080x576 with 0 Axes>



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



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

Recall matrix -----



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

performing the logistic regression

```

In [0]: alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.clas

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

best_alpha = np.argmin(cv_log_error_array)

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

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

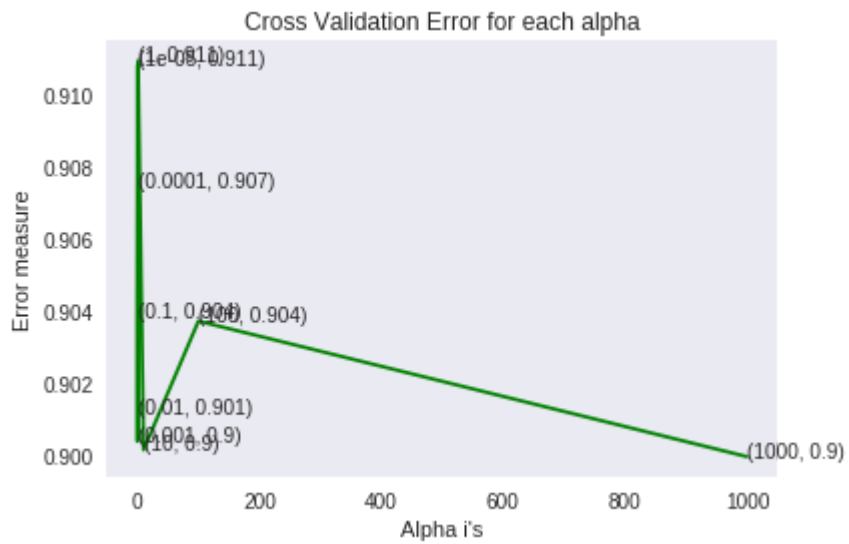
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=logisti
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=logisticR.cla
predict_y = sig_clf.predict_proba(X_test_asm)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

```

```

log_loss for c = 1e-05 is 0.9108268313408105
log_loss for c = 0.0001 is 0.9074469992749733
log_loss for c = 0.001 is 0.9003858534116725
log_loss for c = 0.01 is 0.901146026225207
log_loss for c = 0.1 is 0.9038174284100232
log_loss for c = 1 is 0.9109609189503458
log_loss for c = 10 is 0.9001506529005922
log_loss for c = 100 is 0.9037259486954758
log_loss for c = 1000 is 0.8999618815449928

```

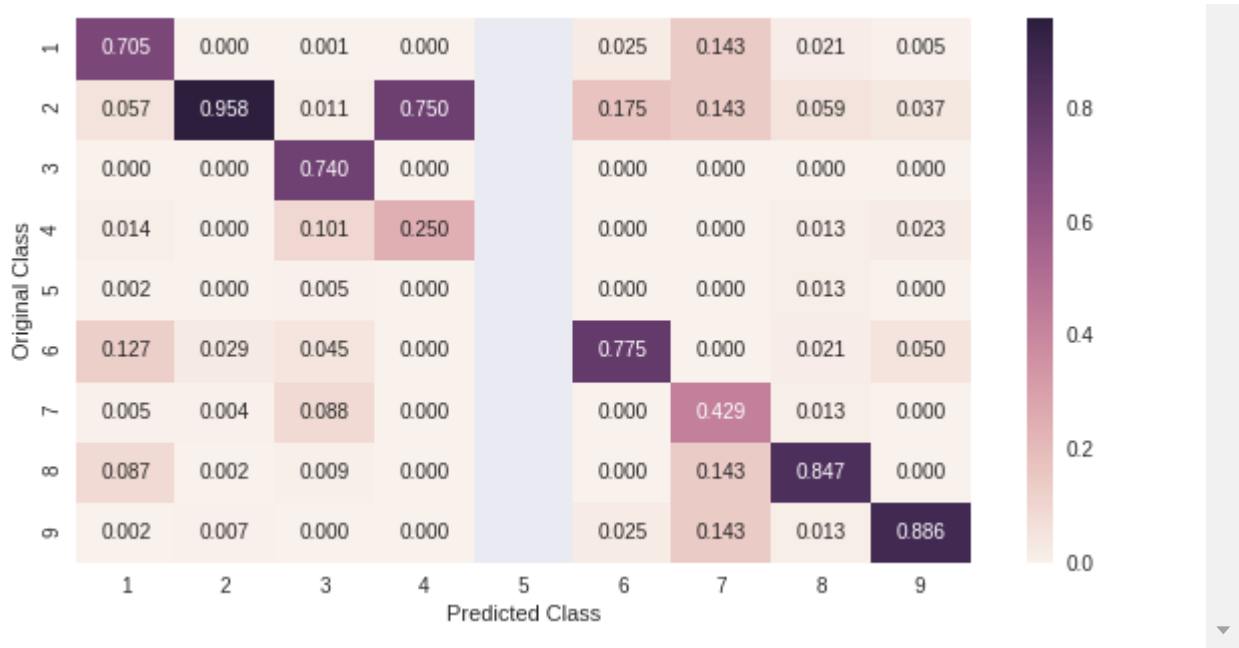


log loss for train data 0.8354959617323964
log loss for cv data 0.8999618815449928
log loss for test data 0.8868849826934874
Number of misclassified points 19.687212511499542
----- Confusion matrix -----

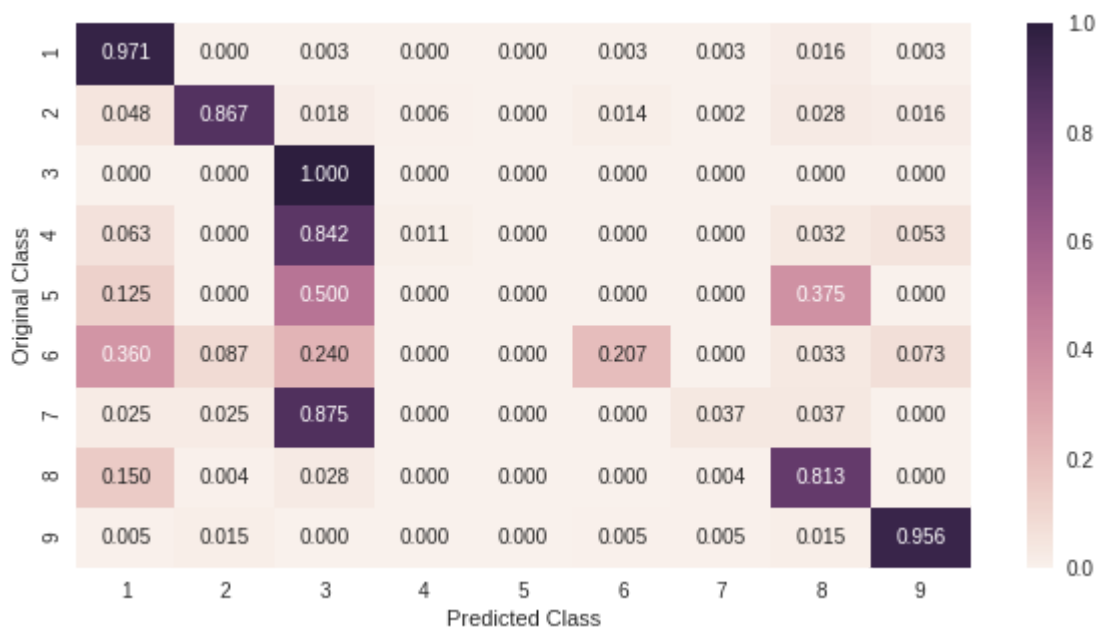
<Figure size 1080x576 with 0 Axes>



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



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



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

```
In [0]: import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use('nbAgg')
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process# this is used for multithreading
import multiprocessing
import codecs# this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

performing the xgboost classifier

```

In [0]: %matplotlib inline
alpha=[10,50,100,500,1000,2000,3000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
    x_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=x_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

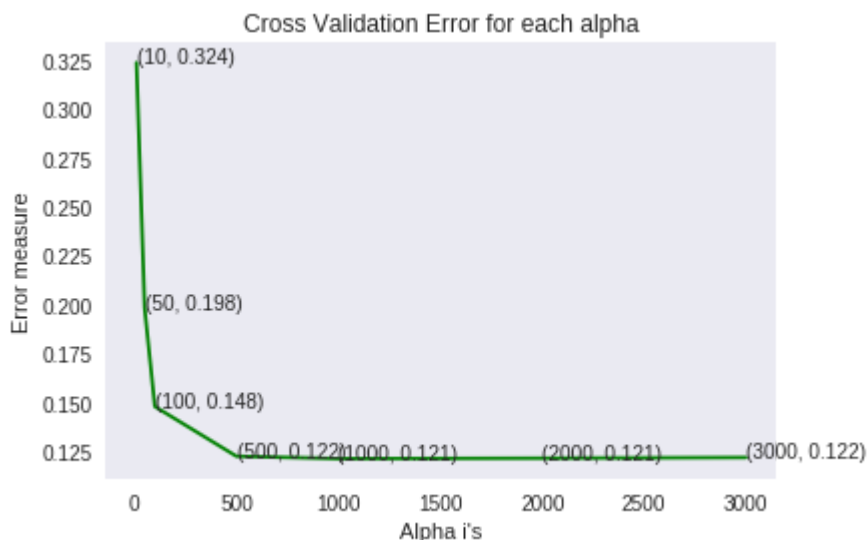
x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
x_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(x_cfl, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)

print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:")
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1
plot_confusion_matrix(y_test_asm,sig_clf.predict(X_test_asm))

log_loss for c = 10 is 0.3235558845671812
log_loss for c = 50 is 0.19808666668071795
log_loss for c = 100 is 0.14781031593982852
log_loss for c = 500 is 0.12240207360914032
log_loss for c = 1000 is 0.12115657918298504
log_loss for c = 2000 is 0.12138005117726003
log_loss for c = 3000 is 0.12173681249500563

```



For values of best alpha = 1000 The train log loss is: 0.03824069498262199
 For values of best alpha = 1000 The cross validation log loss is: 0.12115657918298504

For values of best alpha = 1000 The test log loss is: 0.13977494894836945

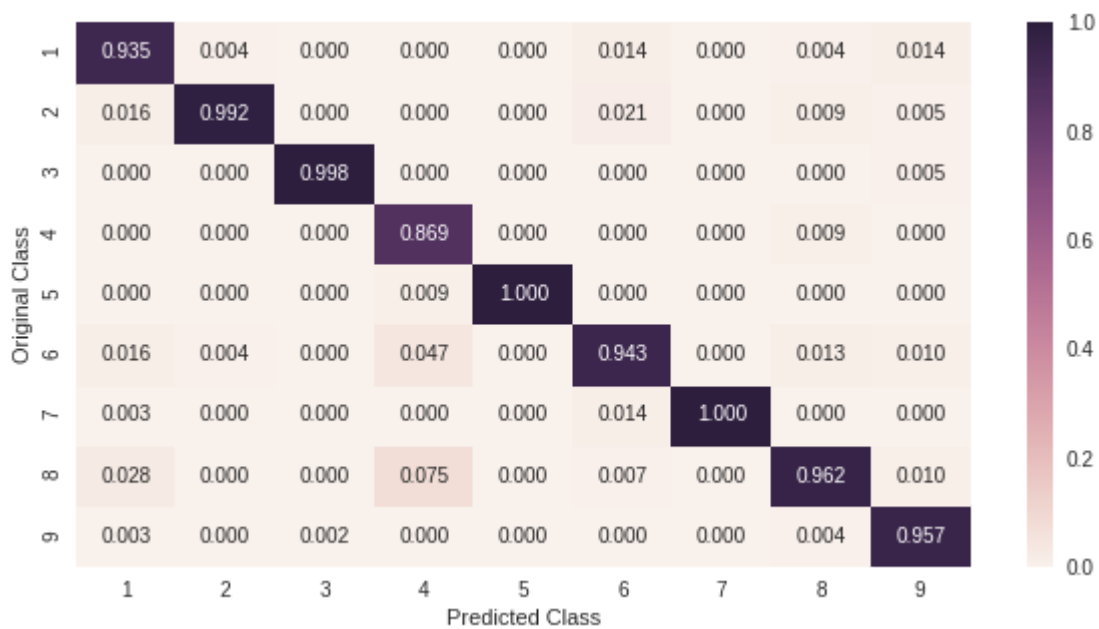
Number of misclassified points 3.035878564857406

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

<Figure size 1080x576 with 0 Axes>

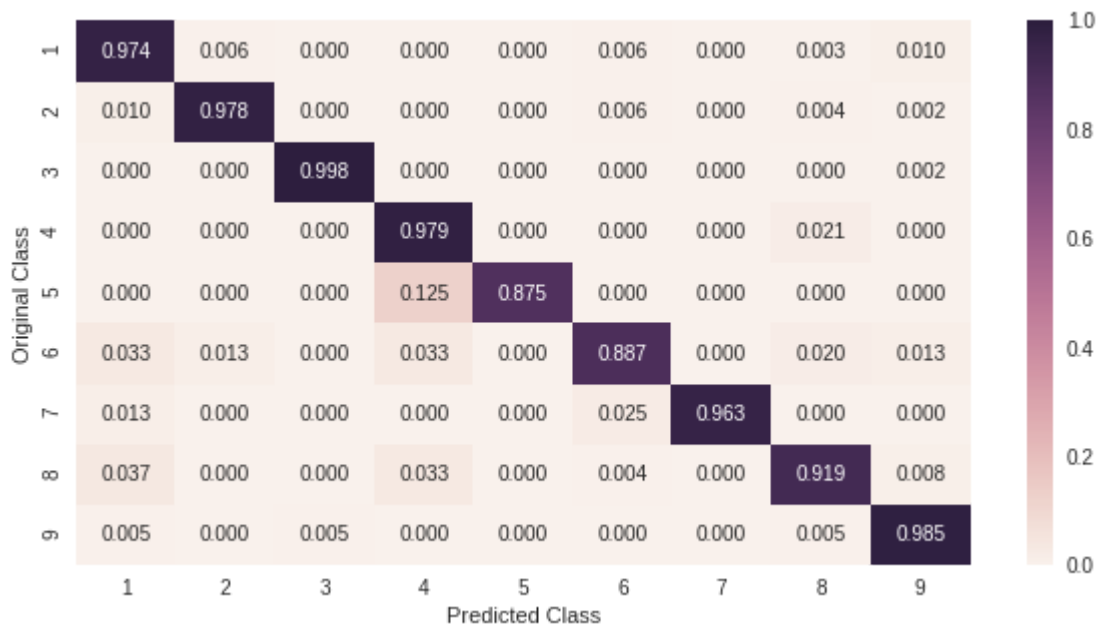


----- Precision matrix -----



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

Recall matrix



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

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

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

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

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done   1 tasks      | elapsed: 11.2min
[Parallel(n_jobs=-1)]: Done   4 tasks      | elapsed: 17.2min
[Parallel(n_jobs=-1)]: Done   9 tasks      | elapsed: 23.9min
[Parallel(n_jobs=-1)]: Done  14 tasks      | elapsed: 46.4min
[Parallel(n_jobs=-1)]: Done  21 tasks      | elapsed: 64.8min
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed: 75.3min finished
```

```
Out[46]: RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
    estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_b
ylevel=1,
    colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
    max_depth=3, min_child_weight=1, missing=None, n_estimators=100,
    n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
    reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
    silent=True, subsample=1),
    fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
    param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15,
0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10], 'co
lsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
    pre_dispatch='2*n_jobs', random_state=None, refit=True,
    return_train_score='warn', scoring=None, verbose=10)
```

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

{'subsample': 1, 'n_estimators': 100, 'max_depth': 3, 'learning_rate': 0.2, 'co
lsample_bytree': 0.3}
```

```
In [0]: x_cfl=XGBClassifier(n_estimators=100,subsample=1,learning_rate=0.2,colsample_bytree=1)
x_cfl.fit(X_train_asm,y_train_asm)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid')
c_cfl.fit(X_train_asm,y_train_asm)

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

train loss 0.04210134627502086
cv loss 0.12931697122069016
test loss 0.13685257159475508

```
In [0]: dta = [['knn',1,0.08,0.25,0.24],['logistic regression',1000,0.83,0.9,0.88],['xgb',1000,0.03,0.12,0.13]]
aa=pd.DataFrame(dta, columns=['model','best hyperparameter','train logloss','cv logloss','test logloss'])
aa
```

```
Out[75]:
```

	model	best hyperparameter	train logloss	cv logloss	test logloss
0	knn	1	0.08	0.25	0.24
1	logistic regression	1000	0.83	0.90	0.88
2	xgb classifier	1000(number of estimators)	0.03	0.12	0.13

USING THE IMAGE FEATURES WITH THE FEATURES OF .BYTE FILES AND .ASM FILES

```
In [0]: # Code to read csv file into Colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)
```

```
In [0]: link = 'https://drive.google.com/open?id=1gYNngh_VHdqW-VtQsPIYSv0_88k1b_JD' # The link to the file
```

```
In [0]: fluff, id = link.split('=')
print(id) # Verify that you have everything after '='

1gYNngh_VHdqW-VtQsPIYSv0_88k1b_JD
```

In [0]:

```
import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('8324-test-image-featuresout-asm.csv')
dataframe1 = pd.read_csv('8324-test-image-featuresout-asm.csv')
```

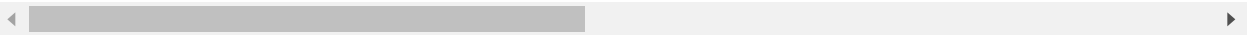
In [0]:

```
dataframe1.head(2)
```

Out[27]:

	filename	ASM_0	ASM_1	ASM_2	ASM_3	ASM_4	ASM_5	ASM_6	ASM_7	ASM_8	ASM_9
0	01azqd4InC7m9JpocGv5	72	69	65	68	69	82	58	48	48	48
1	01IsoiSMh5gxyDYTi4CB	46	116	101	120	116	58	48	48	48	48

2 rows × 1001 columns



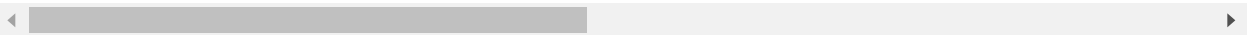
In [0]:

```
dataframe1['ID']=dataframe1['filename']
dataframe1=dataframe1.drop(columns=['filename'])
dataframe1.head(5)
```

Out[28]:

	ASM_0	ASM_1	ASM_2	ASM_3	ASM_4	ASM_5	ASM_6	ASM_7	ASM_8	ASM_9	...	ASM_999
0	72	69	65	68	69	82	58	48	48	52	...	1000
1	46	116	101	120	116	58	48	48	52	48	...	1000
2	72	69	65	68	69	82	58	48	48	52	...	1000
3	72	69	65	68	69	82	58	49	48	48	...	1000
4	72	69	65	68	69	82	58	48	48	52	...	1000

5 rows × 1001 columns



In [0]:

```
result=pd.merge(dataframe,dataframe1,on='ID',how='left')
result.head(5)
print(result.shape)
```

(10868, 1261)

In [0]:

```
asm_y = result['Class']
asm_x = result.drop(['ID', 'Class'], axis=1)
```

#we have generated the data with the image features of .asm files as we are using the byte features and .asm file features

we are doing the preprocessing step of data standardisation and apply the models.

we are applying the models

- knn
- naive bayes
- logistic regression
- random forest

- xgboost with random search cv

```
In [0]: from sklearn.preprocessing import StandardScaler
scale=StandardScaler()
asm_x=scale.fit_transform(asm_x)
```

/usr/local/lib/python3.6/dist-packages/sklearn/preprocessing/data.py:645: DataConversionWarning: Data with input dtype int64, float64 were all converted to float64 by StandardScaler.

return self.partial_fit(X, y)

/usr/local/lib/python3.6/dist-packages/sklearn/base.py:464: DataConversionWarning: Data with input dtype int64, float64 were all converted to float64 by StandardScaler.

return self.fit(X, **fit_params).transform(X)

```
In [0]: from sklearn.model_selection import train_test_split
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x,asm_y ,
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_tr
```

KNN(k nearest neighbors)

```

In [0]: alpha = [x for x in range(1, 10,2)]
cv_log_error_array=[]
for i in alpha:
    print(i)
    k_cfl=KNeighborsClassifier(n_neighbors=i)
    k_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(k_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=k_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

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

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

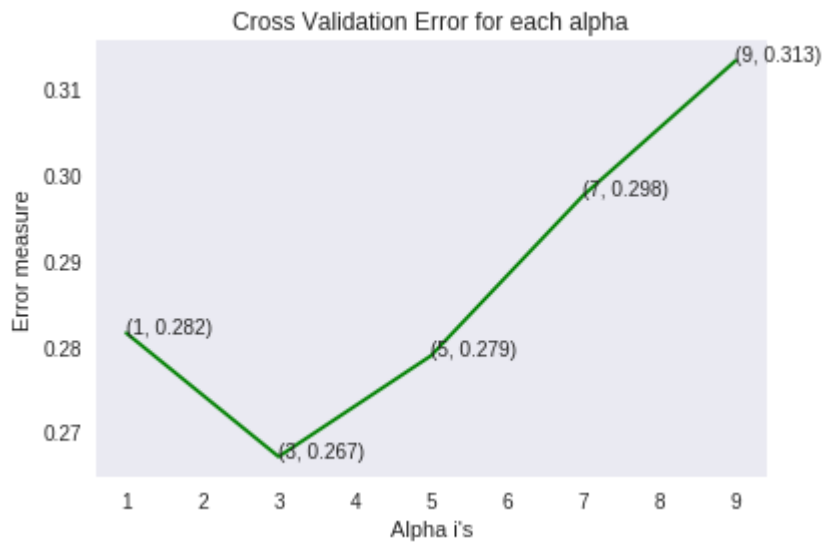
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',log_loss(y_train_asm, predict_y))
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',log_loss(y_cv_asm, predict_y))
predict_y = sig_clf.predict_proba(X_test_asm)
print ('log loss for test data',log_loss(y_test_asm,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

```

```

1
3
5
7
9
log_loss for k = 1 is 0.2815193019820606
log_loss for k = 3 is 0.2671040855888815
log_loss for k = 5 is 0.27886597408651553
log_loss for k = 7 is 0.29756526031536856
log_loss for k = 9 is 0.31326949294254613

```



log loss for train data 0.1466472782761681

log loss for cv data 0.2671040855888815

log loss for test data 0.2629322790620862

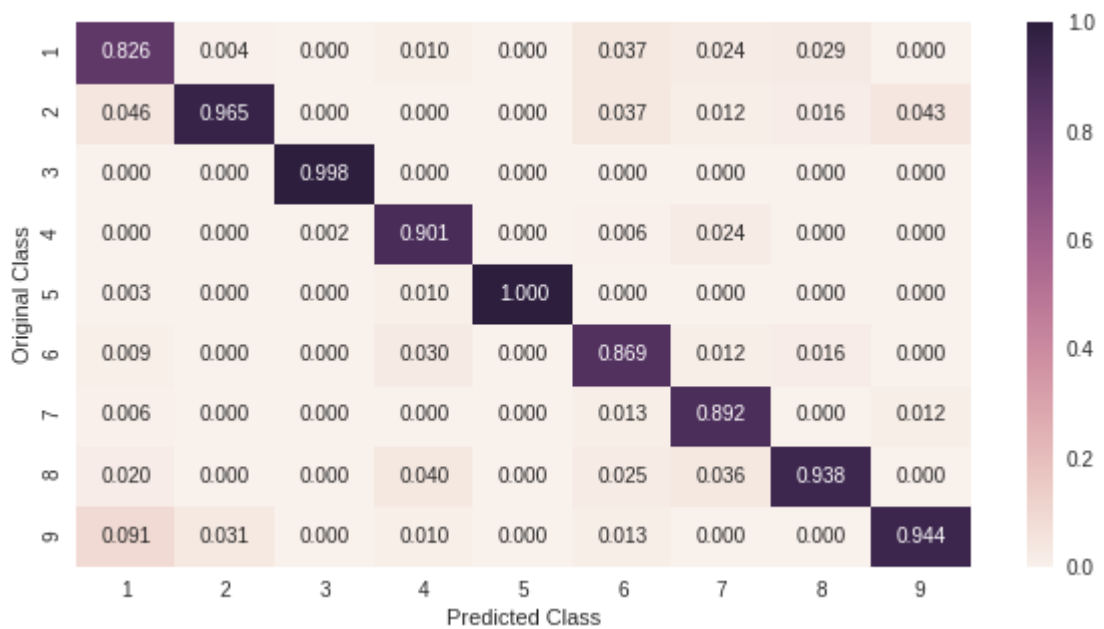
Number of misclassified points 6.5777368905243785

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

<Figure size 1080x576 with 0 Axes>

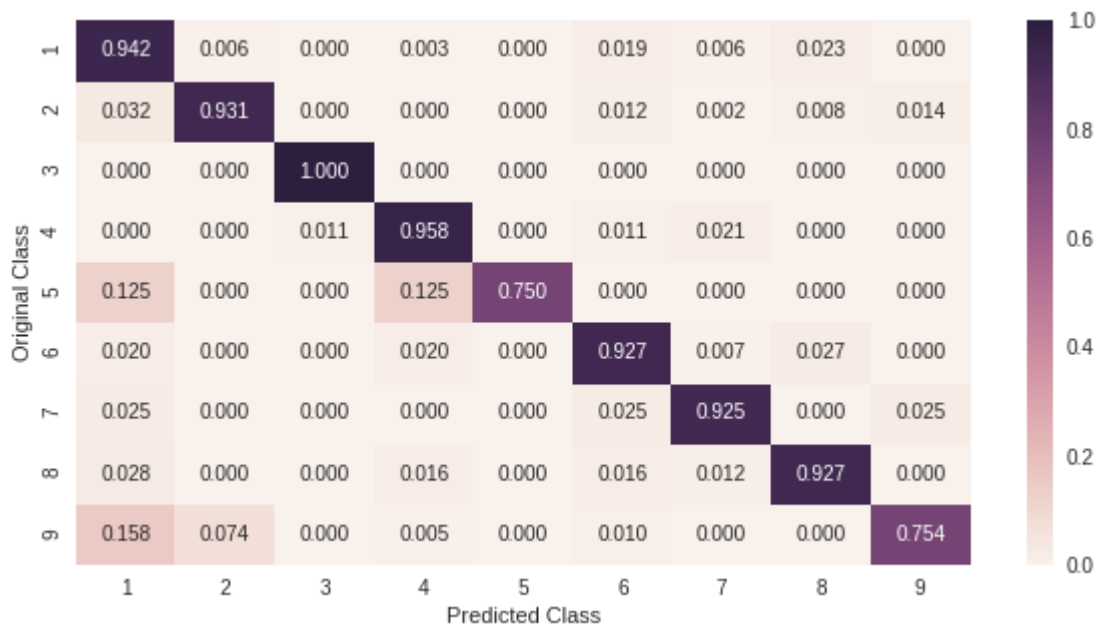


----- Precision matrix -----



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

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



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

#LOGISTIC REGRESSION

```

In [0]: alpha = [10 ** x for x in range(-5, 4)]
cv_log_error_array=[]
for i in alpha:
    logisticR=LogisticRegression(penalty='l2',C=i,class_weight='balanced')
    logisticR.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(logisticR, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=logisticR.clas

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

best_alpha = np.argmin(cv_log_error_array)

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

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

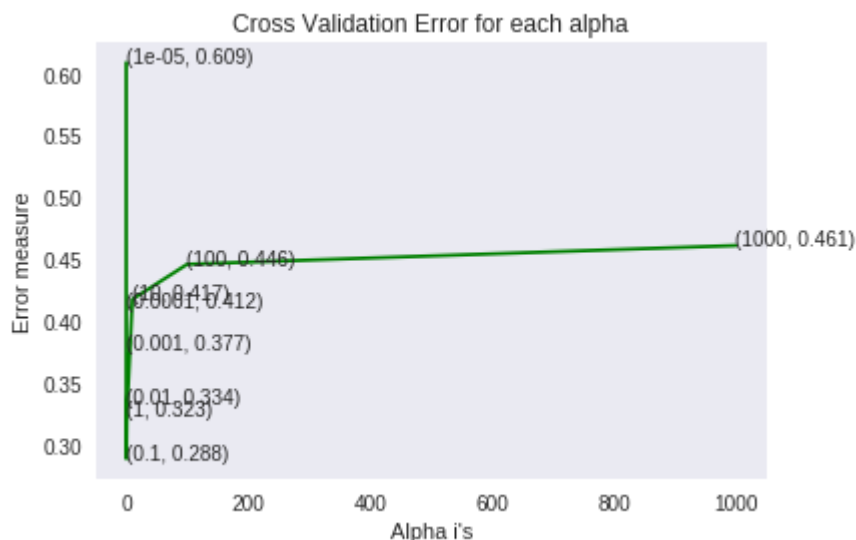
predict_y = sig_clf.predict_proba(X_train_asm)
print ('log loss for train data',(log_loss(y_train_asm, predict_y, labels=logisti
predict_y = sig_clf.predict_proba(X_cv_asm)
print ('log loss for cv data',(log_loss(y_cv_asm, predict_y, labels=logisticR.cla
predict_y = sig_clf.predict_proba(X_test_asm)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

```

```

log_loss for c = 1e-05 is 0.6090034322994382
log_loss for c = 0.0001 is 0.41189347494321016
log_loss for c = 0.001 is 0.3774405398785833
log_loss for c = 0.01 is 0.33410836190450405
log_loss for c = 0.1 is 0.28843301372023583
log_loss for c = 1 is 0.32343885524222354
log_loss for c = 10 is 0.4174086805596062
log_loss for c = 100 is 0.44578402193842614
log_loss for c = 1000 is 0.46095090847608494

```



log loss for train data 0.2640092989362992

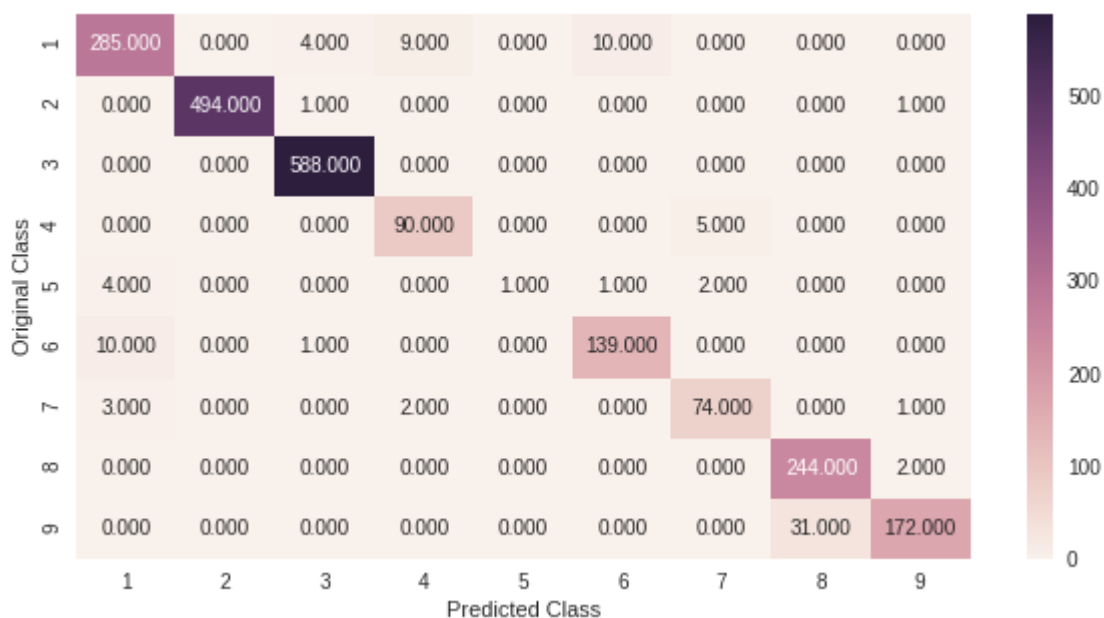
log loss for cv data 0.28843301372023583

log loss for test data 0.282203111170429

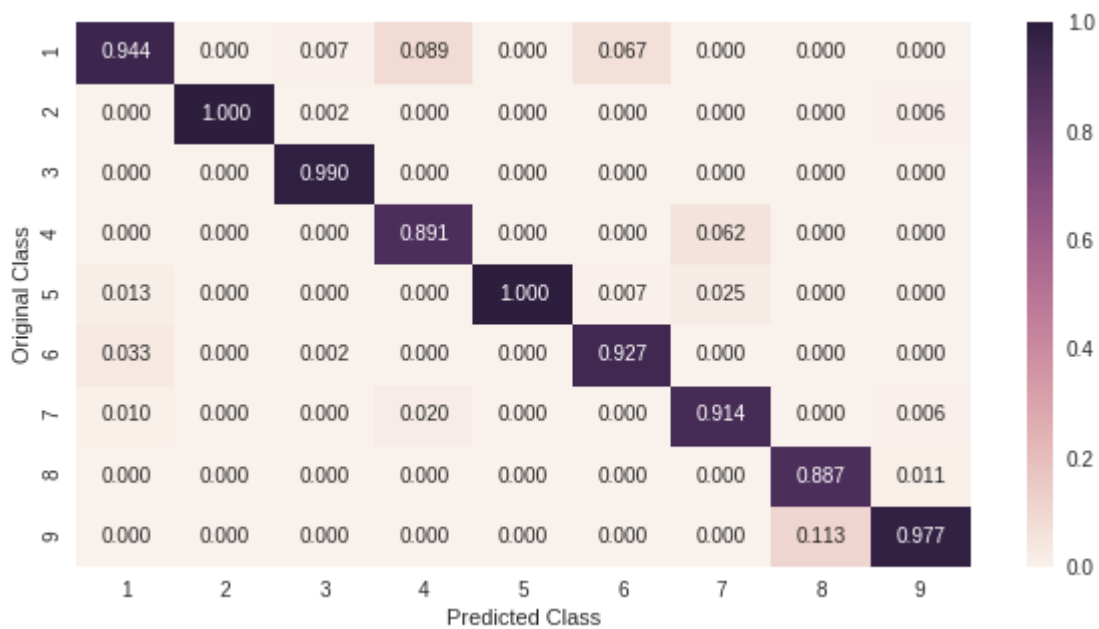
Number of misclassified points 4.001839926402944

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

<Figure size 1080x576 with 0 Axes>

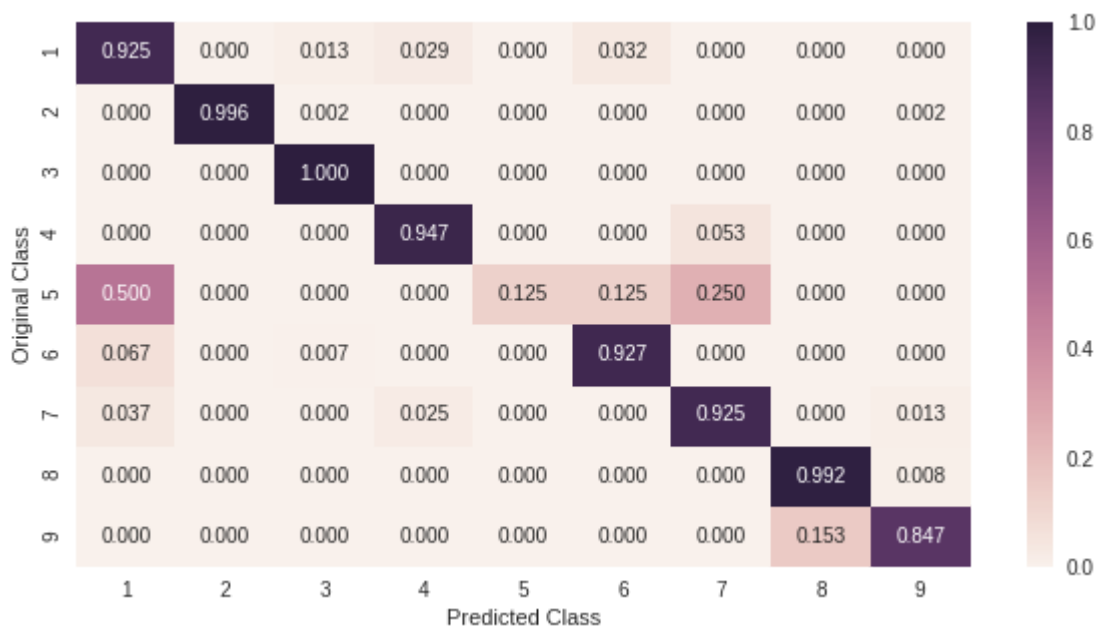


----- Precision matrix -----



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

Recall matrix



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

multinomial naive bayes does not work on negative data we used normalizer instead of using the standard scaler

```
In [0]: from sklearn.preprocessing import Normalizer
norm=Normalizer()
asm_x=norm.fit_transform(asm_x)
```

```
In [0]: from sklearn.model_selection import train_test_split
X_train_asm, X_test_asm, y_train_asm, y_test_asm = train_test_split(asm_x,asm_y ,
X_train_asm, X_cv_asm, y_train_asm, y_cv_asm = train_test_split(X_train_asm, y_tr
```



```

In [0]: from sklearn.naive_bayes import MultinomialNB
alpha = [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100, 1000]
cv_log_error_array = []
for i in alpha:
    print("for alpha =", i)
    clf = MultinomialNB(alpha=i)
    clf.fit(X_train_asm, y_train_asm)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train_asm, y_train_asm)
    sig_clf_probs = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, sig_clf_probs, labels=clf.classes_))
    # to avoid rounding error while multiplying probabilities we use log-probabilities
    print("Log Loss :", log_loss(y_cv_asm, sig_clf_probs))

fig, ax = plt.subplots()
ax.plot(np.log10(alpha), cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array, 3)):
    ax.annotate((alpha[i], str(txt)), (np.log10(alpha[i]), cv_log_error_array[i]))
plt.grid()
plt.xticks(np.log10(alpha))
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(cv_log_error_array)
clf = MultinomialNB(alpha=alpha[best_alpha])
clf.fit(X_train_asm, y_train_asm)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train_asm, y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", 1
plot_confusion_matrix(y_test_asm, sig_clf.predict(X_test_asm))

```

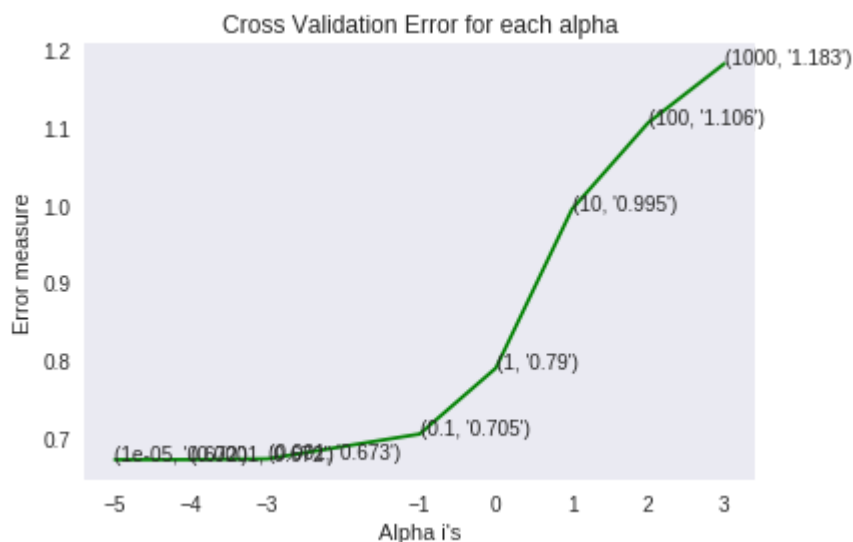
```

for alpha = 1e-05
Log Loss : 0.6715329427802057
for alpha = 0.0001
Log Loss : 0.6716322192344913
for alpha = 0.001
Log Loss : 0.6725542641805241
for alpha = 0.1
Log Loss : 0.7046067230160736
for alpha = 1
Log Loss : 0.7896935980535442
for alpha = 10
Log Loss : 0.9947952571443085
for alpha = 100
Log Loss : 1.1063123268802026

```

for alpha = 1000

Log Loss : 1.182612374427028



For values of best alpha = 1e-05 The train log loss is: 0.6639536954595033

For values of best alpha = 1e-05 The cross validation log loss is: 0.6715329427802057

For values of best alpha = 1e-05 The test log loss is: 0.6789986580234393

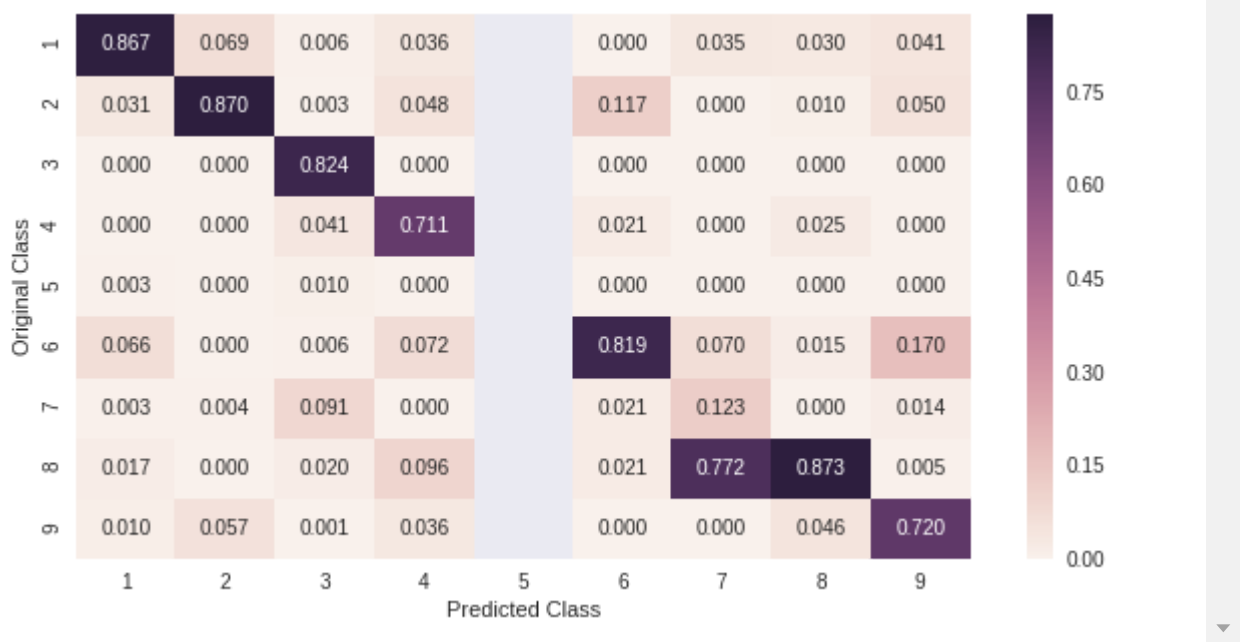
Number of misclassified points 18.813247470101196

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

<Figure size 1080x576 with 0 Axes>

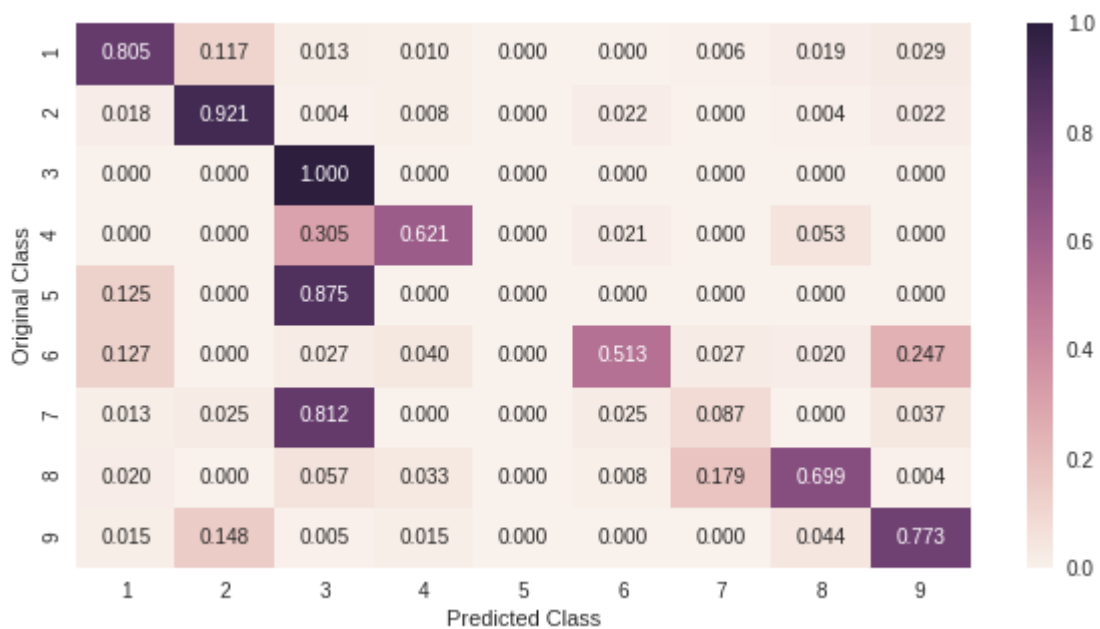


----- Precision matrix -----



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

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



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

Random Forest Classifier

```

In [0]: alpha=[10,50,100,500,1000,2000]
cv_log_error_array=[]
from sklearn.ensemble import RandomForestClassifier
for i in alpha:
    print('current value of alpha is',i)
    r_cfl=RandomForestClassifier(n_estimators=i,random_state=42,n_jobs=-1)
    r_cfl.fit(X_train_asm,y_train_asm)
    sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
    sig_clf.fit(X_train_asm,y_train_asm)
    predict_y = sig_clf.predict_proba(X_cv_asm)
    cv_log_error_array.append(log_loss(y_cv_asm, predict_y, labels=r_cfl.classes_

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

best_alpha = np.argmin(cv_log_error_array)

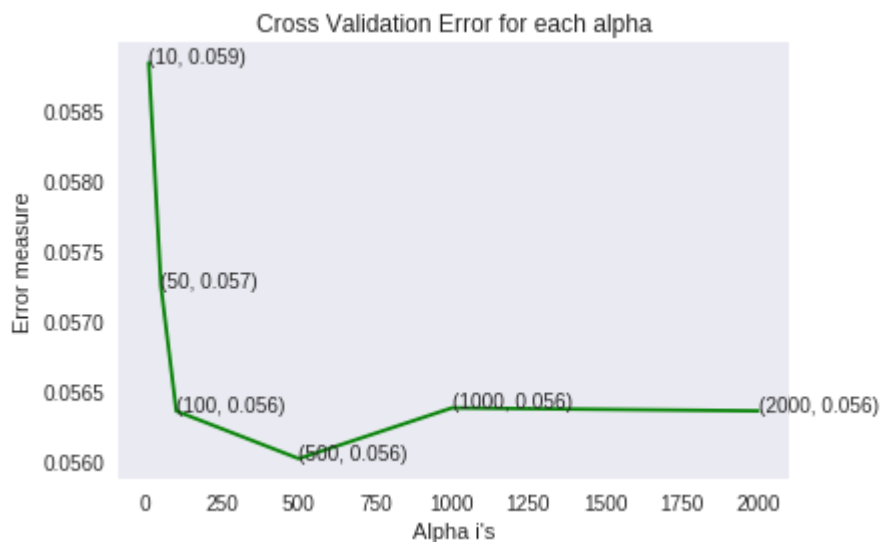
fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array,c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

r_cfl=RandomForestClassifier(n_estimators=alpha[best_alpha],random_state=42,n_job
r_cfl.fit(X_train_asm,y_train_asm)
sig_clf = CalibratedClassifierCV(r_cfl, method="sigmoid")
sig_clf.fit(X_train_asm,y_train_asm)

predict_y = sig_clf.predict_proba(X_train_asm)
print ('For values of best alpha = ', alpha[best_alpha], "The train log loss is:"
predict_y = sig_clf.predict_proba(X_cv_asm)
print('For values of best alpha = ', alpha[best_alpha], "The cross validation log
predict_y = sig_clf.predict_proba(X_test_asm)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",1

current value of alpha is 10
current value of alpha is 50
current value of alpha is 100
current value of alpha is 500
current value of alpha is 1000
current value of alpha is 2000
log_loss for c = 10 is 0.058836192465336654
log_loss for c = 50 is 0.057245687441351355
log_loss for c = 100 is 0.05635909048492949
log_loss for c = 500 is 0.056020994169041605
log_loss for c = 1000 is 0.056381575085719636
log_loss for c = 2000 is 0.05636034947761574

```



For values of best alpha = 500 The train log loss is: 0.020302926521527122

For values of best alpha = 500 The cross validation log loss is: 0.056020994169041605

For values of best alpha = 500 The test log loss is: 0.05165091509556082

XGB CLASSIFIER WITH RANDOM SEARCH CV

```
In [0]: x_cfl=XGBClassifier()
        from sklearn.model_selection import GridSearchCV

        prams={
            'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
            'n_estimators':[100,200,500,1000,2000,3000],
            'max_depth':[3,5,10,20,30],
            'colsample_bytree':[0.1,0.3,0.5,1],
            'subsample':[0.1,0.3,0.5,1]
        }

        random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,scoring='neg_log_lo
        random_cfl.fit(X_train_asm,y_train_asm)
```

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

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 14 tasks      | elapsed: 38.0min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 106.5min finished
```

```
Out[30]: RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
        estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_b
        ylevel=1,
        colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
        max_depth=3, min_child_weight=1, missing=None, n_estimators=100,
        n_jobs=1, nthread=None, objective='binary:logistic', random_state=0,
        reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
        silent=True, subsample=1),
        fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
        param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15,
        0.2], 'n_estimators': [100, 200, 500, 1000, 2000, 3000], 'max_depth': [3, 5, 1
        0, 20, 30], 'colsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.
        5, 1]}},
        pre_dispatch='2*n_jobs', random_state=None, refit=True,
        return_train_score='warn', scoring='neg_log_loss', verbose=5)
```

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

{'subsample': 0.5, 'n_estimators': 500, 'max_depth': 3, 'learning_rate': 0.2,
'colsample_bytree': 0.5}
```

```
In [0]: x_cfl=XGBClassifier(n_estimators=500,subsample=0.5,learning_rate=0.2,colsample_by
x_cfl.fit(X_train_asm,y_train_asm)
c_cfl=CalibratedClassifierCV(x_cfl,method='sigmoid',cv=3)
c_cfl.fit(X_train_asm,y_train_asm)
```

```
Out[32]: CalibratedClassifierCV(base_estimator=XGBClassifier(base_score=0.5, booster='gb
tree', colsample_bylevel=1,
        colsample_bytree=0.5, gamma=0, learning_rate=0.2, max_delta_step=0,
        max_depth=3, min_child_weight=1, missing=None, n_estimators=500,
        n_jobs=1, nthread=None, objective='multi:softprob', random_state=0,
        reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
        silent=True, subsample=0.5),
        cv=3, method='sigmoid')
```

```
In [0]: predict_y = c_cfl.predict_proba(X_train_asm)
print ('train loss',log_loss(y_train_asm, predict_y))
```

train loss 0.010377843567128

```
In [0]: predict_y = c_cfl.predict_proba(X_cv_asm)
print ('cv loss',log_loss(y_cv_asm, predict_y))
```

cv loss 0.013632749128642

```
In [0]: predict_y = c_cfl.predict_proba(X_test_asm)
print ('test loss',log_loss(y_test_asm, predict_y))
```

test loss 0.01431311263546474

```
In [0]: print(X_train_asm.shape)
print(y_train_asm.shape)
```

(6955, 1259)
(6955,)

```
In [0]: %matplotlib inline
plot_confusion_matrix(y_test_asm,c_cfl.predict(X_test_asm))
```

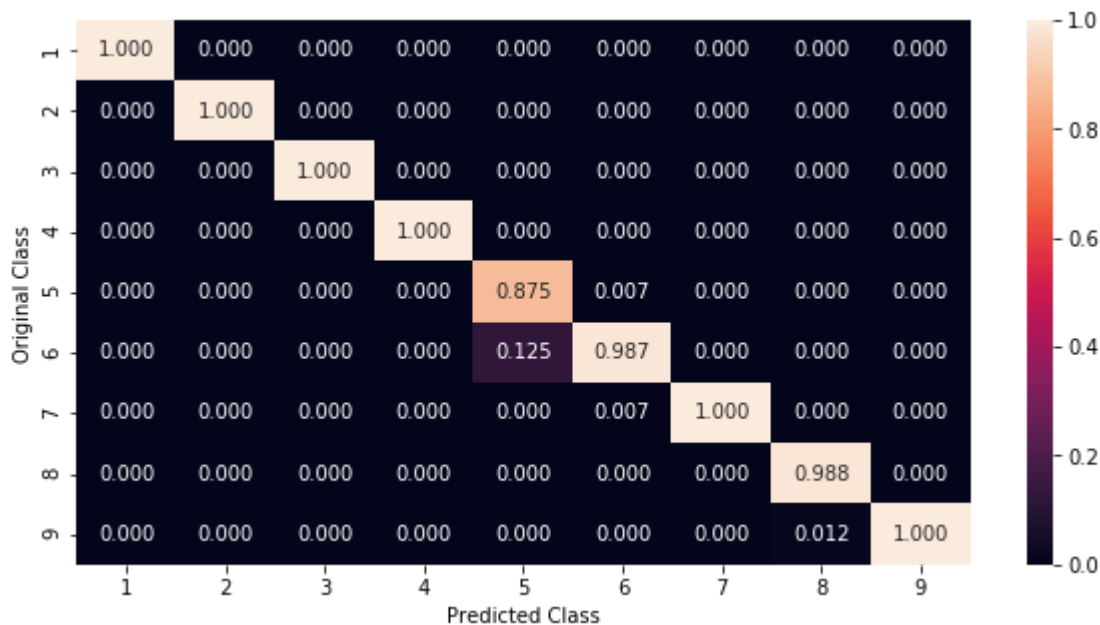
Number of misclassified points 0.27598896044158233

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

<Figure size 1080x576 with 0 Axes>



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



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

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



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

```
In [0]: # Code to read csv file into Colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)
```

```
In [0]: link = 'https://drive.google.com/open?id=12LIFrxVEcAnTSgDDsraUMAi6yzV1zGqk' # The
```

```
In [0]: fluff, id = link.split('=')
print(id) # Verify that you have everything after '='

12LIFrxVEcAnTSgDDsraUMAi6yzV1zGqk
```

```
In [0]: import pandas as pd
downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('finalsubmission - microsoftmalwareddetectiontable - Sheet1.csv')
datafra= pd.read_csv('finalsubmission - microsoftmalwareddetectiontable - Sheet1.csv')
```

```
In [0]: datafra[:5]
```

```
Out[11]:
```

	MACHINE LEARNING MODEL	BEST HYPERPARAMETER VALUE	TRAIN LOG LOSS	CROSS VALIDATION LOGLOSS	TEST LOG LOSS
0	LOGISTIC REGRESSION	0.1(ALPHA)	0.264	0.288	0.282
1	MULTINOMIAL NAIVE BAYES	0.00001(ALPHA)	0.660	0.670	0.670
2	K NEAREST NEIGHBORS	3(VALUE OF K)	0.140	0.260	0.260
3	RANDOM FOREST CLASSIFIER	500(NUMBER OF ESTIMATORS)	0.020	0.050	0.050
4	XGB WITH RANDOM SEARCH CV	NUMBER OF ESTIMATORS=1000,MAXIMUMDEPTH=10,LEAR...	0.010	0.012	0.014

we have achieved log loss of 0.01 with the xgb classifier with random search cv.

in task2 studing the anlaysis done by dchad by dchad in vision of reducing of logloss

the key take aways are

- use the image features
- using chi square method get the best features
- apply various models and reduce the log loss which are the requirements that are satisfied above in the task 2 assignment. *thank you*

using the dchad github account and analysis done by him we are able to achieve the log loss of 0.01 which is very good.

##DOCUMENTATION

- IN TASK 2 WE HAVE TAKEN THE FEATRUES OF .ASM FILES AND .BYTE FILES AND APPLY CHI2 ON TOP OF THE FEATURES AND ANALYSED THE MODELS HOW THEY ARE PERFORMING. THEN WE APPLIED THE IMAGE FEATURES ALONG WITH THE .BYTE FILES AND .ASM FILES APPLIED MODELS ON THEM LIKE
- LOGISTIC REGRESSION
- MULTINOMIAL NAIVE BAYES
- K NEAREST NEIGHBORS
- RANDOM FOREST CLASSIFIER
- XGB WITH RANDOM SEARCH CV

**OBTAINED VERY GOOD RESULTS OF TEST LOGLOSS
LIKE 0.05 AND 0.01.**

In [0]: