

Assignment 7

Joshua Graham

April 8, 2018

Abstract

All Scripts used in the assignment can be found in the A8 folder, if needed.

0.1 Problem 1

Question 1 was pretty much a manual question. Grab 20 spam emails and 20 regular emails. Divide them up into testing and training sets. This was about as hard as question 2. This is because ODU email accounts automatically delete spam emails after 30 days. I had 6 emails, 3 of which aren't usable. So if you examine the emails, you may notice references to a Karen and/or Scott Graham. These are my parents, I asked to use their email for spam messages and they gave the okay. With these and a few from one of my other gmail accounts, I managed to scrape together 20 spam emails. To keep the data set related, I also grabbed a few of their emails from the not spam sections to keep the data set from being skewed. There was no automation in this question so I have no code to show. **However**, during question 2, 19 out of 20 emails were listed as not spam, so instead of recreating the data, I swapped half the files from each data set which produced a more reasonable dataset.

0.2 Problem 2

Question 2 had code. Namely, I had to correctly implement test.py and change the code slightly in docclass.py. Starting with docclass.py, it was relatively simple. I used the glob import to find all the .txt files under each respective training directory, then put them into a variable. This is what I passed to cl.train. The changed method is shown below.

Training

```
1 def spamTrain(cl):
2     path = "Training/NotSpam/*.txt"
3     for filename in glob.glob(path):
4         with open(filename, 'r') as myfile:
5             print("Training_file:", filename)
6             data = myfile.read()
7             cl.train(data, 'not_spam')
8     path = "Training/Spam/*.txt"
9     for filename in glob.glob(path):
10        with open(filename, 'r') as myfile:
11            print("Training_file:", filename)
12            data = myfile.read()
13            cl.train(data, 'spam')
```

After changing docclass.py I had to correctly implement test.py. This wasn't a challenge, all I had to do was point at the testing files at the right time. I again used the glob import, for no reason other than I like using wildcard and it got the job done. test.py can be shown below.

test.py

```
1 import docclass
2 import glob
3 from subprocess import check_output
4
5 cl = docclass.naivebayes(docclass.getwords)
6 try:
7     check_output(['rm', 'spam.db'])
8 except:
9     print("No_file")
10 #remove previous db file
11 cl.setdb('spam.db')
12 docclass.spamTrain(cl)
13 spampath="Testing/Spam/*.txt"
14 notspampath="Testing/NotSpam/*.txt"
15
16 print("Testing_against_spam_emails.")
17 for filename in glob.glob(spampath):
18     with open(filename, 'r') as myfile:
19         data = myfile.read()
20         print("Filename:", filename, "Status:", cl.classify(data))
21 print("Testing_against_regular_emails.")
22 for filename in glob.glob(notspampath):
23     with open(filename, 'r') as myfile:
24         data = myfile.read()
25         print("Filename:", filename, "Status:", cl.classify(data))
```

All in all, I think I'm happy with how it turned out. After running the code, this was the output:

test.py

```
1 Filename: Testing/Spam/18.txt Status:  not spam
2 Filename: Testing/Spam/20.txt Status:  spam
3 Filename: Testing/Spam/3.txt Status:  spam
4 Filename: Testing/Spam/7.txt Status:  spam
5 Filename: Testing/Spam/9.txt Status:  spam
6 Filename: Testing/Spam/12.txt Status:  not spam
7 Filename: Testing/Spam/5.txt Status:  spam
8 Filename: Testing/Spam/14.txt Status:  spam
9 Filename: Testing/Spam/16.txt Status:  spam
10 Filename: Testing/Spam/1.txt Status:  spam
11 Testing against regular emails.
12 Filename: Testing/NotSpam/18.txt Status:  not spam
13 Filename: Testing/NotSpam/20.txt Status:  not spam
14 Filename: Testing/NotSpam/3.txt Status:  not spam
15 Filename: Testing/NotSpam/7.txt Status:  not spam
16 Filename: Testing/NotSpam/9.txt Status:  not spam
17 Filename: Testing/NotSpam/12.txt Status:  not spam
18 Filename: Testing/NotSpam/5.txt Status:  not spam
19 Filename: Testing/NotSpam/14.txt Status:  not spam
20 Filename: Testing/NotSpam/16.txt Status:  not spam
21 Filename: Testing/NotSpam/1.txt Status:  not spam
```

I think there was a bit of skew because of using someone else's emails. But it correctly represented all but 2 emails. On a larger data set, the error would probably be larger. However, for the given example, it appears to work decently well.