

# **CS 432: Homework Number Eight**

Due on April 13, 2019 at 4:20 PM

*Alexander Nwala*

**Tim Bruce**

## Problem 1

Create two datasets; the first called Testing, the second called Training.

The Training dataset should:

- Consist of 10 text documents for email messages you consider spam (from your spam folder)
- Consist of 10 text documents for email messages you consider not spam (from your inbox)

The Testing dataset should:

- Consist of 10 text documents for email messages you consider spam (from your spam folder)
- Consist of 10 text documents for email messages you consider not spam (from your inbox)

Upload your datasets on github Please do not include emails that contain sensitive information

### Solution

This one is pretty self explanatory. I found the emails and copy pasted them to various text files on my computer organized into four files of ten based on if they were spam. It was hard to find enough non-spam emails that were not too personal to put on GitHub. The dataset does have one notable curveball. Two emails from the same source that talk about very similar topics, one of which is spam and one of which is not. Specifically, one of them is a daily report from my favorite ski mountain, Sugarloaf. The other email is a notification that they have had a significant amount of snow. I want the snow warning, because it may change driving conditions en-route to the mountain, and skiing conditions will be significantly changed as well (not always for the better). The daily report is garbage (well written, but useless to receive every day)

## Problem 2

Using the PCI book modified docclass.py code and test.py (see Slack assignment-8 channel) Use your Training dataset to train the Naive Bayes classifier ( e.g., docclass.spamTrain() ) Use your Testing dataset to test (test.py) the Naive Bayes classifier and report the classification results.

### Solution

I got the files from the slack, and edited test.py to fit my needs. In order to pull up the dataset, I created the get\_emails function to get the dataset.

Listing 1: get\_emails function.

```
1 def get_emails(directory):
2     a = [directory + x for x in os.listdir(directory)]
3     emails = []
4     for email in a:
5         f = open(email)
6         string = ""
7         for line in f.readlines():
8             string += line
9         emails.append(string)
10        f.close()
11    return emails
```

Then, I train the Bayes classifier using the training data set, and classify the other emails. It should be noted that the training set is alternating entry, with one being spam followed by not spam. This is to prevent overloading of the last one to train for. Even though this randomization is in place, for my entire testing dataset, the classifier thinks that it is spam. There could be several causes:

- My testing and training datasets are too disjoint, and the not spam testing values are similar to my spam messages from the training.
- The to from and subject fields are messing up the content, and skewing it towards spam.
- A combination of the above two.

Either way, the classifier is not useful if everything is spam, but it does clean up your inbox really well.