

CSCI-580 Project Spam & Ham Classifier

Objective: In this project, you will implement a Naïve Bayes Classifier to classify emails into spam or ham classes.

You will be provided with four files: *training_ham.txt*, *training_spam.txt*, *testing_ham.txt* and *testing_spam.txt*. Your program must read the names of these files from the command line. In addition, your program will read an integer, *k*, from the command line. This is the command line arguments order:

```
./run training_ham.txt training_spam.txt testing_ham.txt testing_spam.txt 2
```

Here, “run” is an executable, and “2” is *k*, the integer that you will use for Laplace smoothing.

Specification Requirements:

1. Read in *training_ham.txt* and put all the words into *map<string, int>* (I will call this map *wordsHam*), where key is a word, and integer represents the count of the word in this file with *ham* emails. While reading this file, count the total number of words in this file, and count the total number of messages in this file (total lines).
2. Do the same with *training_spam.txt* file (I will call the map *wordsSpam*).
3. Calculate *a priori class probabilities*:

$$P(\text{ham}) = \frac{\text{totalMsgHam} + k}{\text{totalMsgBoth} + kn}$$

$$P(\text{spam}) = \frac{\text{totalMsgSpam} + k}{\text{totalMsgBoth} + kn}$$

where $n = 2$ (since there are two classes, namely *ham* and *spam*), and k is the last command line argument; *totalMsgBoth* is the total number of messages (lines) in both files *training_ham.txt* and *training_spam.txt*.

4. For each word in the dictionary (each word occurring in either of the files), calculate the conditional word likelihood probabilities:

$$P(w_i|\text{ham}) = \frac{\text{countHam} + k}{\text{totalWordsHam} + kn}$$

$$P(w_i|\text{spam}) = \frac{\text{countSpam} + k}{\text{totalWordsSpam} + kn}$$

Here, n is the size of the dictionary (total distinct words in both *ham* and *spam*), k is the same as before, *countHam* and *countSpam* are total of occurrences of the given word in *ham* and *spam* respectively, and *totalWordsHam* and *totalWordsSpam* are total words in *ham* and *spam* respectively.

Consider using *map<string, pair<double, double>>* to store these probabilities, where the key is a given word, and *pair<double, double>* has $P(w|\text{ham})$ as *first* and $P(w|\text{spam})$ as *second* members of pair.

5. Read in *testing_ham.txt* and *testing_spam.txt* and for each email (line), calculate the unnormalized posterior probabilities of the given message being *ham* and *spam*:

$$\begin{aligned} P(\text{ham}|\text{msg}) &= P(\text{ham}|w_1, w_2, \dots, w_n) \\ &= P(w_1, w_2, \dots, w_n|\text{ham})P(\text{ham}) \\ &= P(\text{ham})P(w_1|\text{ham})P(w_2|\text{ham}) \dots P(w_n|\text{ham}) \end{aligned}$$

$$\begin{aligned} P(\text{spam}|\text{msg}) &= P(\text{spam}|w_1, w_2, \dots, w_n) \\ &= P(w_1, w_2, \dots, w_n|\text{spam})P(\text{spam}) \\ &= P(\text{spam})P(w_1|\text{spam})P(w_2|\text{spam}) \dots P(w_n|\text{spam}) \end{aligned}$$

In case, a word is not in the dictionary (has not occurred in *training_ham.txt* or *training_spam.txt*), choose a small probability for the conditional $P(w|\text{ham})$ or $P(w|\text{spam})$:

$$P(w_i|\text{ham}) = \frac{1}{\text{totalWordsHam} + \text{totalWordsSpam}}$$

$$P(w_i|\text{spam}) = \frac{1}{\text{totalWordsHam} + \text{totalWordsSpam}}$$

To avoid underflow, instead of the product of probabilities, use the sum of their logs:

$$\begin{aligned} \log P(\text{ham}|\text{msg}) &= \log P(\text{ham}) + \log P(w_1|\text{ham}) + \log P(w_2|\text{ham}) + \dots + \log P(w_n|\text{ham}) \\ \log P(\text{spam}|\text{msg}) &= \log P(\text{spam}) + \log P(w_1|\text{spam}) + \log P(w_2|\text{spam}) + \dots + \log P(w_n|\text{spam}) \end{aligned}$$

6. Classify a message as *ham* if $\log(P(\text{ham}|\text{msg})) > \log(P(\text{spam}|\text{msg}))$, and as *spam*, otherwise.
7. Calculate *Specificity*, *Sensitivity* and *Accuracy* of the classifier:

$$\textbf{Specificity} = \frac{\text{True Negative}}{\text{True Negative} + \text{False Positive}} = \frac{\text{True Negative}}{\text{Total Negative in a Population}}$$

$$\textbf{Sensitivity} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} = \frac{\text{True Positive}}{\text{Total Positive in a Population}}$$

$$\textbf{Accuracy} = \frac{\text{True Positive} + \text{True Negative}}{\text{Total size of data}}$$

8. Output the required results using *cout*. First, for each line in *testing_ham.txt*, output on a separate line “ham” or “spam” word according to your classification, followed by the

$\log(P(\text{ham} | \text{msg}))$ and $\log(P(\text{spam} | \text{msg}))$ values. Second, repeat similar output for each message in *testing_spam.txt* file. Finally, the last line of the output will contain the values of *Specificity*, *Sensitivity* and *Accuracy* of the classifier.

Input Format:

Each of the input files above have lines inside of them, where each line represents a separate email.

Output Format:

```
ham -40.0978 -51.307
spam -62.5378 -56.8432
0.69 0.927 0.74
```

```
<word><space><logHam><space><logSpam><endl>
<specificity><space><sensitivity><space><accuracy><endl>
```

Additional Specification Requirements:

- Break your program into smaller tasks and write a separate function for each task. Your *main()* function should just declare necessary data structures (or classes) and call functions to accomplish the goal of this assignment.
- You may write a C++ class to complete your program, but please put the class's definition and description into *main.cpp* file.
- Follow the best coding practice rules such as pass-by-reference large data structures or big objects as parameters to the functions to save time/space of copying large objects; and include comments to clarify your code.

Submission:

- 1) Submit *main.cpp* to [turnin](#).
- 2) In addition to submitting files to turnin, you need to submit *main.cpp* to Blackboard. Make sure to include your name inside *main.cpp* file.

Failure to submit *main.cpp* file to Blackboard will result in 0pts for your project.