

Lab Topic 5 - Spam Filtering Using A Naïve Bayes Classifier with Laplace Smoothing CNM Peralta

I. Naïve Bayes Classifier

A Spam Filter classifies messages (commonly emails) to remove incoming spam. This is the technology that prevents you from seeing spam in your inbox folder as much as possible; the spam instead goes to the Spam folder. A common way to implement a spam filter is to represent all received emails as in a bag-of-words and compute the probability of a message being spam given the words that make it up.

One of the ways to solve the spam filtering problem is via the Naive Bayes classifier technique. This technique is based on the Bayesian theorem which states that,

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}.$$

This theorem describes the probability of an event based on prior knowledge of conditions relating to the event. The naivety of the Naive Bayes classifier comes from its strong assumption that the conditions are **independent** of each other, that is, the features/attributes describing the event.

Using the Bayesian network, the relationship between the likelihood of the presence of a word (provided a spam or ham classification) could be modeled. That is, $P(\text{Spam} | \text{message})$, where $\text{message} = w_0 w_1 w_2 \dots w_n$. The Bayes network that describes the spam filtering problem is:

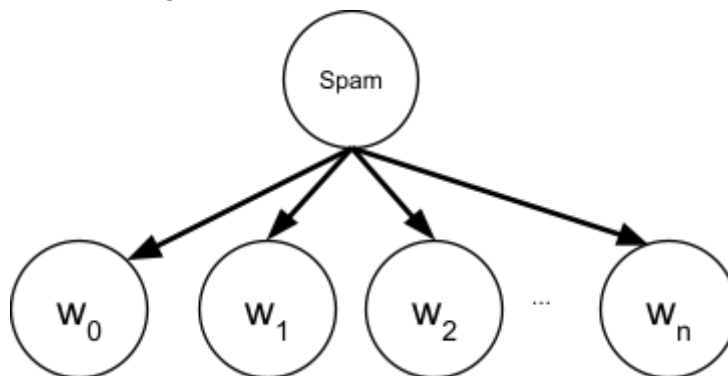


Figure 1. The Bayes* network for the Spam filtering problem.

Thus, if we apply Bayes' rule to $P(\text{Spam} | \text{message})$, we get:

$$P(\text{Spam} | \text{message}) = \frac{P(\text{message} | \text{Spam}) P(\text{Spam})}{P(\text{message})}$$

Spam filtering is a **supervised learning problem**. We are initially given a data set of both Spam and Ham messages. The data set is the set of messages, and the target labels indicate whether a message is Spam/Ham. Since *we already know whether a message is Spam/Ham*, the target labels are already given, thus making spam filtering a supervised learning problem.

Lastly, it is important to note that $\text{Ham} = \neg \text{Spam}$. We will use Ham and $\neg \text{Spam}$ interchangeably.

How do we compute $P(\text{Spam})$?

The probability of Spam can be computed as:

$$P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$$

That is, the number of spam messages over the total number of messages (both Spam and Ham). Consequently, we can compute the probability of Ham as:

$$P(Ham) = \frac{\text{count}(Ham)}{\text{count}(Spam \cup Ham)} = 1 - P(Spam)$$

How do we compute $P(message | Spam)$?

Given the Bayes network in Figure 1, we can observe that if we know the value of the Spam variable, the probabilities of the words become independent. That is,

$$P(message|Spam) = P(w_0 w_1 \dots w_n | Spam) = P(w_0 | Spam) P(w_1 | Spam) \dots P(w_n | Spam)$$

The probability of each word given Spam is given by:

$$P(w|Spam) = \frac{P(w, Spam)}{P(Spam)} = \frac{\text{count}(w \text{ in spam})}{\text{count}(\text{totality of the words in spam})}$$

Applying this, we get:
$$P(Spam | message) = \frac{P(w_0 | Spam) P(w_1 | Spam) \dots P(w_n | Spam) P(Spam)}{P(message)}$$

How do we compute $P(message)$?

The total probability of getting a message is:

$$P(message) = P(message|Spam)P(Spam) + P(message|Ham)P(Ham)$$

We already know how to compute $P(message|Spam)$, $P(Spam)$, and $P(Ham)$. So we only need to figure out how to compute $P(message|Ham)$. To do this, we only need to apply the same concept we used in computing their Spam counterpart.

We thus have:

$$P(message|Ham) = P(w_0 w_1 \dots w_n | Ham) = P(w_0 | Ham) P(w_1 | Ham) \dots P(w_n | Ham)$$

And to compute each $P(w|Ham)$, we have:

$$P(w|Ham) = \frac{P(w, Ham)}{P(H)} = \frac{\text{count}(w \text{ in Ham})}{\text{count}(\text{totality of the words in Ham})}$$

Putting it all together...

We can, therefore, compute $P(Spam|message)$ as...

$$P(S | m) = \frac{P(w_0 | S) P(w_1 | S) \dots P(w_n | S) P(S)}{P(w_0 | S) P(w_1 | S) \dots P(w_n | S) P(S) + P(w_0 | H) P(w_1 | H) \dots P(w_n | H) P(H)}$$

where S stands for Spam, m stands for the message, and H stands for Ham (to conserve space).

Lastly,

A **threshold** must be set in order to determine which among the two (2) classes does a given message fall into. For example, when a 50% probability is set in order to classify a file as a spam, once the computed $P(Spam|message)$ falls below 50%, it is automatically classified as a ham. Otherwise, the file will be classified as spam.

Where does the bag-of-words representation come in?

Given the Spam and Ham data set, **two bags-of-words must be constructed: one for Spam and one for Ham**. Each $P(w|Spam)$ and $P(w|Ham)$ is just the **frequency of the word in the data set specified by Spam/Ham**. You thus have everything you need to compute $P(Spam|message)$.

II. Augmenting the Spam Filter with Laplace Smoothing

The problem with the Naïve Bayes Classifier is that, when it encounters a word that is not in either the Spam or Ham data set, the entire message is immediately classified to the other set. For example, if we have:

Spam	Ham
Ang pogi pogi pogi ni Sir James	Ang pogi ni Sir Mir
	Ang cute ni Sir Froi
	Ang macho ni Sir JM

When we attempt to compute $P(\text{Spam} | \text{'Ang cute ni Sir James'})$, we will get:

$$P(\text{Spam} | \text{'Ang cute ni Sir James'}) = \frac{P(\text{'ang'}|\text{Spam})P(\text{'cute'}|\text{Spam})P(\text{'ni'}|\text{Spam})P(\text{'sir'}|\text{Spam})P(\text{'james'}|\text{Spam})P(\text{Spam})}{P(\text{'Ang cute ni Sir James'})}$$

And looking at the spam data set, we can see that the word 'cute' does not occur; thus, $P(\text{'cute'}|\text{Spam}) = 0$, and consequently, the entire numerator becomes 0. Since the numerator will be 0,

$$P(\text{Spam} | \text{'Ang cute ni Sir James'}) = \frac{0}{P(\text{'Ang cute ni Sir James'})} = 0$$

Thus, if a message has a word that occurs in only one data set, it is automatically classified to that data set.

Furthermore, if a word does not exist in both sets, a $\frac{0}{0}$ operation will occur, which we know to be undefined. For example:

$$P(\text{Spam} | \text{'ang gwapo'}) = \frac{P(\text{'ang'}|\text{Spam})P(\text{'gwapo'}|\text{Spam})P(\text{Spam})}{P(\text{'ang'}|\text{Spam})P(\text{'gwapo'}|\text{Spam})P(\text{Spam}) + P(\text{'ang'}|\text{Ham})P(\text{'gwapo'}|\text{Ham})P(\text{Ham})}$$

$$P(\text{Spam} | \text{'ang gwapo'}) = \frac{P(\text{'ang'}|\text{Spam}) \times 0 \times P(\text{Spam})}{P(\text{'ang'}|\text{Spam}) \times 0 \times P(\text{Spam}) + P(\text{'ang'}|\text{Ham}) \times 0 \times P(\text{Ham})}$$

$$P(\text{Spam} | \text{'ang gwapo'}) = \frac{0}{0}$$

This is a manifestation of **overfitting**; we want to be more friendly to new words to both sets.

Enter Laplace Smoothing

Laplace smoothing handles the case of new words by introducing "phantom words" to the data set. For example:

Spam	Ham
Ang pogi pogi pogi ni Sir James	Ang pogi ni Sir Mir
	Ang cute ni Sir Froi
	Ang macho ni Sir JM

Given the data set, the bags-of-words are:

Spam Bag-of-Words	
Word	Frequency
ang	1
pogi	3
ni	1
sir	1
james	1

Ham Bag-of-Words	
Word	Frequency
ang	3
pogi	1
ni	3
sir	3
mir	1
cute	1
froi	1
macho	1

		jm	1
--	--	----	---

We will modify the formulas that we introduced last meeting to accommodate new words. This is done by using a **smoothing factor**, k , that will be used in the modified formulas.

HOW DO WE COMPUTE $P(\text{Spam})$?

The probability of Spam can be computed as:

$$P(\text{Spam}) = \frac{\text{count}(\text{Spam}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

The **probability of Ham is computed in the same way**. The effect of this smoothing is that an additional k observations of both ham and spam are taken into consideration, even if they don't really exist.

HOW DO WE COMPUTE $P(\text{message} | \text{Spam})$?

Recall that the probability of each word given Spam is given by:

$$P(w|\text{Spam}) = \frac{\text{count}(w \text{ in Spam})}{\text{count}(\# \text{ of words in Spam})}$$

Applying the Laplace smoothing, we have:

$$P(w|\text{Spam}) = \frac{\text{count}(w \text{ in Spam}) + k}{\text{count}(\# \text{ of words in Spam}) + (k \times (\text{dictionarySize} + \text{count}(\text{new words})))}$$

dictionarySize = number of unique words found on both Spam and Ham
new words = words to be classified not found on both Spam and Ham

The effect of this modification is that the frequency of all words is increased by k observations each. Thus, if we compute $P('gwapo'|\text{Spam})$ from our example earlier, and $k = 2$, we will get:

$$P('gwapo'|\text{Spam}) = \frac{\text{count}('gwapo' \text{ in Spam}) + k}{\text{count}(\# \text{ of words in Spam}) + (k \times (\text{dictionarySize} + \text{count}(\text{new words})))}$$

$$P('gwapo'|\text{Spam}) = \frac{0+2}{7+(2 \times (10+1))}$$

The number of new words is 1 because 'ang gwapo' only has one new word: 'gwapo.' Thus, $P('gwapo'|\text{Spam})$ is no longer 0, and we can avoid a $\frac{0}{0}$ operation.

Formula Cheat Sheet (Naive Bayes Classifier with Laplace smoothing)

$$P(\text{Spam} | \text{message}) = \frac{P(\text{message} | \text{Spam}) P(\text{Spam})}{P(\text{message})}$$

$$P(\text{Spam}) = \frac{\text{count}(\text{Spam}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

$$P(\text{Ham}) = \frac{\text{count}(\text{Ham}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

$$P(\text{message} | \text{Spam}) = P(w_0|\text{Spam})P(w_1|\text{Spam})... P(w_n|\text{Spam})$$

$$P(w|\text{Spam}) = \frac{\text{count}(w \text{ in Spam}) + k}{\text{count}(\# \text{ of words in Spam}) + (k \times (\text{dictionarySize} + \text{count}(\text{new words})))}$$

dictionarySize: number of **unique** words found on both Spam and Ham
new words: words to be classified not found on both Spam and Ham

$$P(\text{message}) = P(\text{message}|\text{Spam})P(\text{Spam}) + P(\text{message}|\text{Ham})P(\text{Ham})$$

$$P(\text{message} | \text{Ham}) = P(w_0|\text{Ham})P(w_1|\text{Ham})... P(w_n|\text{Ham})$$

$$P(w|\text{Ham}) = \frac{\text{count}(w \text{ in Ham}) + k}{\text{count}(\# \text{ of words in Ham}) + (k \times (\text{dictionarySize} + \text{count}(\text{new words})))}$$

$$P(\text{Spam} | \text{message}) = \frac{P(w_0|\text{Spam})P(w_1|\text{Spam})...P(w_n|\text{Spam}) P(\text{Spam})}{P(w_0|\text{Spam})P(w_1|\text{Spam})...P(w_n|\text{Spam}) + P(w_0|\text{Ham})P(w_1|\text{Ham})...P(w_n|\text{Ham})}$$

Exercise

You may use your previous bag-of-words representation exercise as a starting point for this exercise. Given the spam and ham datasets provided by your lab instructors, your program should:

- Make the bag-of-words representation of the Spam data set from the given **spam folder**.
- Make the **separate** bag-of-words representation of the Ham data set from the given **ham folder**.
- Read a **folder named classify** containing a series of messages; there will be **one message per file**.
- Ask the user for the value of the **smoothing factor, k**.
- Compute the probability that each message is Spam using Naive Bayes with Laplace Smoothing. Output your answers to a text file, **classify.out**. The classify.out file must have the filename of the file to be classified, its classification, and the computed probability in one line.

```
001 SPAM 0.85122
002 HAM 0.4302
099 SPAM 0.5001
```

In your terminal output, display the dictionary size and the total number of words of each bow:

```
HAM
Dictionary Size: 2303
Total Number of Words: 14578

SPAM
Dictionary Size: 6886
Total Number of Words: 20409
```

Hints:

- Read files in Latin-1 encoding
- Check out decimal module in Python and BigDecimal in Java (to avoid division by zero error)
- $k = 5$ is used in the sample data set output (for checking)
- Bonus 1 (2pts): Make a UI for your spam filter. It should be of the following format:

Select Spam Folder		Select Ham Folder		Dictionary Size:	Total Words:
Word	Frequency	Word	Frequency		
Total Words in Spam:		Total Words in Ham:			
<div style="display: flex; justify-content: space-between; align-items: center;"> <div> Select Classify Folder k <input style="width: 50px;" type="text"/> </div> <div> <input type="button" value="Filter"/> </div> </div>					
Output					
Filename	Class	P(Spam)			

- Bonus 2 (5 pts): Look for the best value for the smoothing factor (i.e. parameter tuning) using **cross-validation**:
 - Randomly partition your data (spam and ham folder) into three (3) non-overlapping subsets: training, tuning, and test subsets.
 - The training subset should be 70% of your whole data set; the tuning subset should be 20%; and the remaining 10% should be the test subset. That is, 210 spam files and 210 ham files (total of 420) will be used for the training subset; 60 spam and 60 ham files (total of 120) for the tuning subset; and the remaining files 30 spam and 30 ham files (60 files) will be for the test subset (see video).
 - Generate the spam and ham bag of words using the files in the training subset.
 - Classify the files in the **tuning subset** using Laplace smoothing and the training subset's bow. You will start with $k=1$, and increase the value of k by 1 until the accuracy of your agent becomes greater than or equal to 90%.
 - Once you find a k with accuracy $\geq 90\%$, you will then classify the files in the **test subset** using the k you found. Compute for the accuracy of the new model given the test subset.

The results of the cross-validation above must be written in a report using LaTeX. The template (containing the guide questions) can be downloaded [here](#).

- The criteria for scoring this exercise is:

Criteria	Score
Correctly generate bag-of-words for Spam	2
Correctly generate bag-of-words for Ham	2
Correctly compute $P(\text{Spam})$ (and $P(\text{Ham})$)	2
Correctly compute $P(\text{word} \text{Spam})$	3
Correctly compute $P(\text{Spam} \text{message})$	5
Correctly output results to file	1
Total	15

References:

Stuart Russell and Peter Norvig. 2009. *Artificial Intelligence: A Modern Approach* (3rd ed.). Prentice Hall Press, Upper Saddle River, NJ, USA.