



CMSC 170: Introduction to Artificial intelligence

Week 05: Spam Filtering using a Naïve Bayes Classifier

John O-Neil V. Geronimo
Institute of Computer Science
University of the Philippines Los Baños



Content

- I. Background
- II. Implementing Spam Filter using Naïve Bayes Classifier



Content

- I. Background
- II. Implementing Spam Filter using Naïve Bayes Classifier

Background



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.

Background



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.

Background



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,

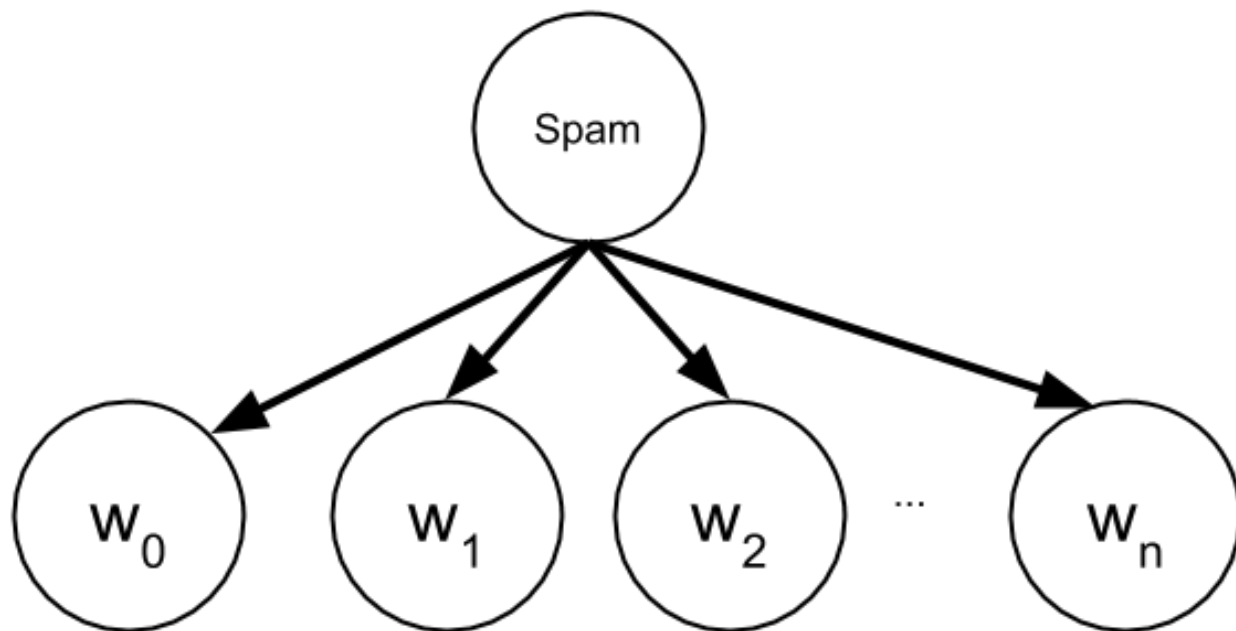
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Background



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.

Background



Naïve Bayes Theorem

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

Spam Filtering Using Naïve Bayes

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

Background

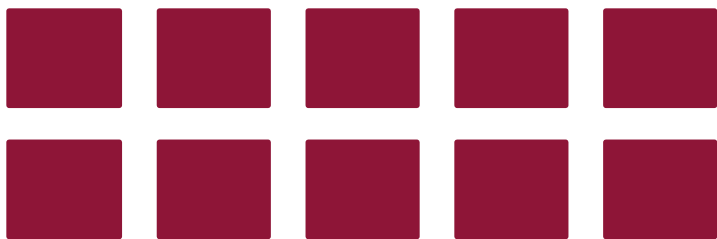


We are initially given a data set of both Spam and Ham messages. The data set is the set of messages, and the target labels indicate whether a message is Spam/Ham. Since we already know whether a message is Spam/Ham, the target labels are already given, thus making spam filtering a **supervised learning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

Background

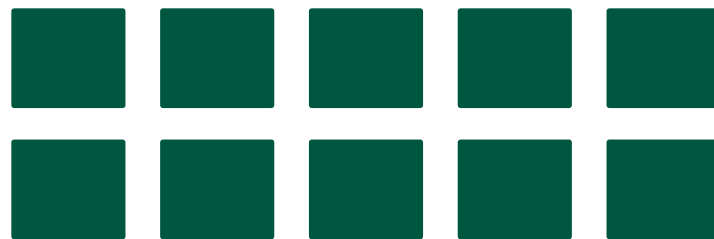


Spam Bag-of-Words



| Index | Word | Frequency |
|-------|------|-----------|
| | | |
| | | |
| | | |

Ham Bag-of-Words



| Index | Word | Frequency |
|-------|------|-----------|
| | | |
| | | |
| | | |

Formula:

$$\begin{aligned} & P(\text{Spam}|\text{message}) \\ &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \end{aligned}$$



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

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

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

Formula:

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

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



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

$$\text{message} = w_0 w_1 w_2 \dots w_n$$

$$\begin{aligned} P(\text{message}|\text{Spam}) &= P(w_0 w_1 w_2 \dots w_n | \text{Spam}) \\ &= P(w_0 | \text{Spam}) P(w_1 | \text{Spam}) \dots P(w_n | \text{Spam}) \end{aligned}$$

Formula:

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

- $P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$
- $P(\text{Ham}) = 1 - P(\text{Spam})$
- $P(\text{message}|\text{Spam})$
 $= P(w_0|\text{Spam}) \dots P(w_n|\text{Spam})$

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

$$P(w_n|\text{Spam}) = \frac{P(w, \text{Spam})}{P(\text{Spam})}$$

$$= \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$$



Formula:

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

- $P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$
- $P(\text{Ham}) = 1 - P(\text{Spam})$
- $P(\text{message}|\text{Spam})$
 $= P(w_0|\text{Spam}) \dots P(w_n|\text{Spam})$
- $P(w_n|\text{Spam})$
 $= \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$

How do we compute for $P(\text{message}|\text{Ham})$?

$$\text{message} = w_0 w_1 w_2 \dots w_n$$

$$\begin{aligned} P(\text{message}|\text{Ham}) &= P(w_0 w_1 w_2 \dots w_n|\text{Ham}) \\ &= P(w_0|\text{Ham})P(w_1|\text{Ham}) \dots P(w_n|\text{Ham}) \end{aligned}$$



Formula:

$$\begin{aligned} & P(\text{Spam}|\text{message}) \\ &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \end{aligned}$$

- $P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$
- $P(\text{Ham}) = 1 - P(\text{Spam})$
- $P(\text{message}|\text{Spam})$
 $= P(w_0|\text{Spam}) \dots P(w_n|\text{Spam})$
- $P(w_n|\text{Spam})$
 $= \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$
- $P(\text{message}|\text{Ham})$
 $= P(w_0|\text{Ham}) \dots P(w_n|\text{Ham})$

How do we compute for $P(w|\text{Ham})$?

$$\begin{aligned} P(w_n|\text{Ham}) &= \frac{P(w, \text{Ham})}{P(\text{Ham})} \\ &= \frac{\text{count}(w \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} \end{aligned}$$



Formula:

$$\begin{aligned} & P(\text{Spam}|\text{message}) \\ &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \end{aligned}$$

- $P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$
- $P(\text{Ham}) = 1 - P(\text{Spam})$
- $P(\text{message}|\text{Spam})$
 $= P(w_0|\text{Spam}) \dots P(w_n|\text{Spam})$
- $P(w_n|\text{Spam})$
 $= \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$
- $P(\text{message}|\text{Ham})$
 $= P(w_0|\text{Ham}) \dots P(w_n|\text{Ham})$
- $P(w_n|\text{Ham})$
 $= \frac{\text{count}(w \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$



How do we compute for P(message)?

$$\begin{aligned} P(\text{message}) &= \\ & P(\text{message}|\text{Spam})P(\text{Spam}) \\ & + P(\text{message}|\text{Ham})P(\text{Ham}) \end{aligned}$$

Formula:

$$\begin{aligned} & P(\text{Spam}|\text{message}) \\ &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \end{aligned}$$

- $P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$
- $P(\text{Ham}) = 1 - P(\text{Spam})$
- $P(\text{message}|\text{Spam})$
 $= P(w_0|\text{Spam}) \dots P(w_n|\text{Spam})$
- $P(w_n|\text{Spam})$
 $= \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$
- $P(\text{message}|\text{Ham})$
 $= P(w_0|\text{Ham}) \dots P(w_n|\text{Ham})$
- $P(w_n|\text{Ham})$
 $= \frac{\text{count}(w \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$
- $P(\text{message})$
 $= P(\text{message}|\text{Spam})P(\text{Spam})$
 $+ P(\text{message}|\text{Ham})P(\text{Ham})$



Putting it all together,

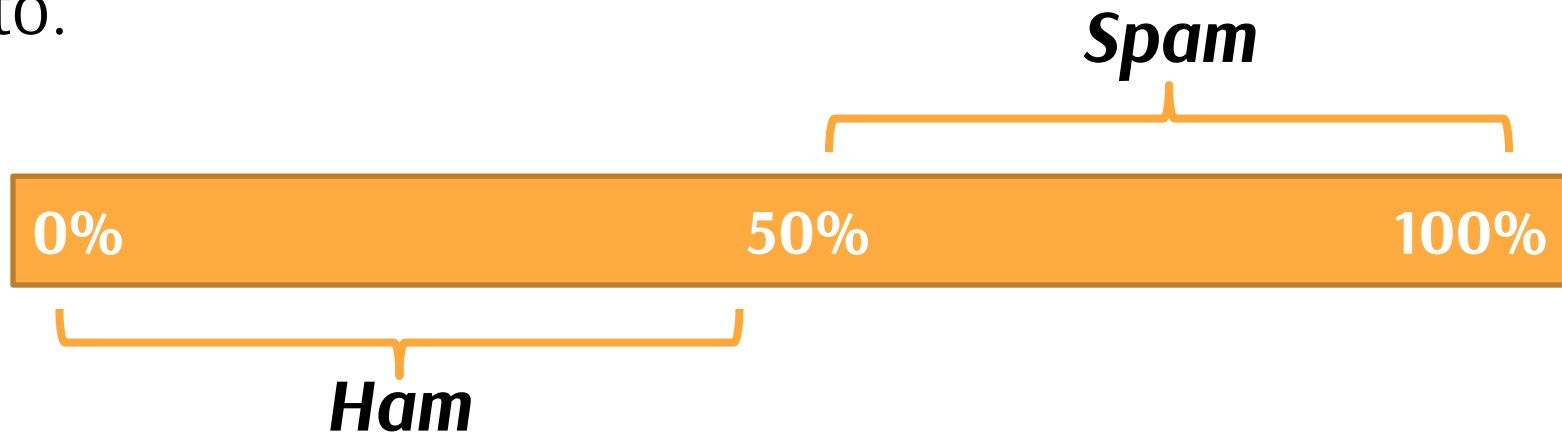
$$\begin{aligned} & P(S|m) \\ &= \frac{P(w_0|S)P(w_1|S) \dots P(w_n|S)P(S)}{P(w_0|S)P(w_1|S) \dots P(w_n|S)P(S) + P(w_0|H)P(w_1|H) \dots P(w_n|H)P(H)} \end{aligned}$$

where $m = \text{message}$, $S = \text{Spam}$, and $H = \text{Ham}$

Background



Lastly, a **threshold** must be set in order to determine which among the two (2) classes does a given message fall into.





Content

- I. Background
- II. Implementing Spam Filter using Naïve Bayes Classifier

Spam Filter



Spam Dataset



You will love this product!



I used this cream and it removed my acne!!!



Buy now!

Ham Dataset



Lucky I am in love with my bestfriend



You are still the one that I love



Could I love you anymore?

Spam Filter



Spam Bag-of-Words

| Index | Word | Frequency |
|-------|---------|-----------|
| 0 | you | 1 |
| 1 | will | 1 |
| 2 | love | 1 |
| 3 | this | 2 |
| 4 | product | 1 |
| 5 | i | 1 |
| 6 | used | 1 |
| 7 | cream | 1 |

| Index | Word | Frequency |
|-------|---------|-----------|
| 8 | and | 1 |
| 9 | it | 1 |
| 10 | removed | 1 |
| 11 | my | 1 |
| 12 | acne | 1 |
| 13 | buy | 1 |
| 14 | now | 1 |

TNOW:

16

DS:

15

Ham Bag-of-Words

| Index | Word | Frequency |
|-------|------------|-----------|
| 0 | lucky | 1 |
| 1 | i | 3 |
| 2 | am | 1 |
| 3 | in | 1 |
| 4 | love | 3 |
| 5 | with | 1 |
| 6 | my | 1 |
| 7 | bestfriend | 1 |

| Index | Word | Frequency |
|-------|---------|-----------|
| 8 | you | 2 |
| 9 | are | 1 |
| 10 | still | 1 |
| 11 | the | 1 |
| 12 | one | 1 |
| 13 | that | 1 |
| 14 | could | 1 |
| 15 | anymore | 1 |

TNOW:

21

DS:

16

Spam Filter



Is this message ham or spam?

I love you.

Spam Filter



Compute for $P(\text{Spam})$:

Since we have three emails in the spam dataset and a total of six messages for both spam and ham dataset

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

$$P(\text{Spam}) = \frac{3}{6} = \mathbf{0.5}$$

Spam Filter



Compute for $P(\text{Ham})$:

Since we have three emails in the ham dataset and a total of six messages for both spam and ham dataset

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

$$P(\text{Ham}) = \frac{3}{6} = \mathbf{0.5}$$

Compute for $P(\text{message}|\text{Spam})$:

Remember that we can compute $P(\text{message}|\text{spam})$ using this formula:

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

$P(w_n|\text{Spam})$ can be computed using this formula:

$$P(w_n|\text{Spam}) = \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$$

Spam Filter



Compute for $P(\text{message}|\text{Spam})$:

$$P("i"|\text{Spam}) = \frac{\text{count}("i" \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$$

$$P("love"|\text{Spam}) = \frac{\text{count}("love" \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$$

$$P("you"|\text{Spam}) = \frac{\text{count}("you" \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$$

Spam Bag-of-Words

| Index | Word | Frequency |
|-------|---------|-----------|
| 0 | you | 1 |
| 1 | will | 1 |
| 2 | love | 1 |
| 3 | this | 2 |
| 4 | product | 1 |
| 5 | i | 1 |
| 6 | used | 1 |
| 7 | cream | 1 |

| Index | Word | Frequency |
|-------|---------|-----------|
| 8 | and | 1 |
| 9 | it | 1 |
| 10 | removed | 1 |
| 11 | my | 1 |
| 12 | acne | 1 |
| 13 | buy | 1 |
| 14 | now | 1 |

TNOW:

16

DS:

15

Spam Filter



Compute for $P(\text{message}|\text{Spam})$:

$$P(\text{"i"}|\text{Spam}) = \frac{\text{count}(\text{"i" in Spam})}{\text{count}(\text{total no of words in Spam})} = \frac{1}{16}$$

$$P(\text{"love"}|\text{Spam}) = \frac{\text{count}(\text{"love" in Spam})}{\text{count}(\text{total no of words in Spam})} = \frac{1}{16}$$

$$P(\text{"you"}|\text{Spam}) = \frac{\text{count}(\text{"you" in Spam})}{\text{count}(\text{total no of words in Spam})} = \frac{1}{16}$$

Compute for $P(\text{message}|\text{Spam})$:

$$\begin{aligned}P(\text{message}|\text{Spam}) &= P(w_0|\text{Spam})P(w_1|\text{Spam}) \dots P(w_n|\text{Spam}) \\&= P("i"|\text{Spam})P("love"|\text{Spam})P("you"|\text{Spam}) \\&= \frac{1}{16} * \frac{1}{16} * \frac{1}{16} \\&= \frac{1}{4096}\end{aligned}$$

Compute for $P(\text{message}|\text{Ham})$:

Remember that we can compute $P(\text{message}|\text{Ham})$ using this formula:

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

$P(w_n|\text{Ham})$ can be computed using this formula:

$$P(w_n|\text{Ham}) = \frac{\text{count}(w \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$$

Spam Filter



Compute for $P(\text{message}|\text{Ham})$:

$$P("i"|\text{Ham}) = \frac{\text{count}("i" \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$$

$$P("love"|\text{Ham}) = \frac{\text{count}("love" \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$$

$$P("you"|\text{Ham}) = \frac{\text{count}("you" \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$$

Spam Filter



Ham Bag-of-Words

| Index | Word | Frequency |
|-------|------------|-----------|
| 0 | lucky | 1 |
| 1 | i | 3 |
| 2 | am | 1 |
| 3 | in | 1 |
| 4 | love | 3 |
| 5 | with | 1 |
| 6 | my | 1 |
| 7 | bestfriend | 1 |

| Index | Word | Frequency |
|-------|---------|-----------|
| 8 | you | 2 |
| 9 | are | 1 |
| 10 | still | 1 |
| 11 | the | 1 |
| 12 | one | 1 |
| 13 | that | 1 |
| 14 | could | 1 |
| 15 | anymore | 1 |

TNOW:

21

DS:

16

Spam Filter



Compute for $P(\text{message}|\text{Spam})$:

$$P(\text{"i"}|\text{Ham}) = \frac{\text{count}(\text{"i"} \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} = \frac{3}{21}$$

$$P(\text{"love"}|\text{Ham}) = \frac{\text{count}(\text{"love"} \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} = \frac{3}{21}$$

$$P(\text{"you"}|\text{Ham}) = \frac{\text{count}(\text{"you"} \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} = \frac{2}{21}$$

Compute for $P(\text{message}|\text{Ham})$:

$$\begin{aligned}P(\text{message}|\text{Ham}) &= P(w_0|\text{Ham})P(w_1|\text{Ham}) \dots P(w_n|\text{Ham}) \\&= P("i"|\text{Ham})P("love"|\text{Ham})P("you"|\text{Ham}) \\&= \frac{3}{21} * \frac{3}{21} * \frac{2}{21} \\&= \frac{2}{1029}\end{aligned}$$

Compute for $P(\text{Spam}|\text{message})$:

$$\begin{aligned}P(\text{Spam}|\text{message}) &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \\&= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message}|\text{Spam})P(\text{Spam}) + P(\text{message}|\text{Ham})P(\text{Ham})} \\&= \frac{\frac{1}{4096}(0.5)}{\frac{1}{4096}(0.5) + \frac{2}{1029}(0.5)} = 0.1115931027 = 11.16\%\end{aligned}$$

Spam Filter



Since $P(\text{Spam}|\text{message}) = 11.16\%$

The message “I love you.”
is classified as a **ham** message!



Content

- I. Background
- II. Implementing Spam Filter using Naïve Bayes Classifier



Keep safe!