

Princeton University
COS 217: Introduction to Programming Systems
Spam Filter Knowledge Module: Mathematical Foundations

First, let's define abbreviations for some logical assertions:

f_i : The message contains feature i .
 $\sim f_j$: The message does not contain feature j .
ham: The message is ham.
spam: The message is spam.

We wish to compute:

(1)

$P(\text{ham} \mid f_1, \dots, f_m, \sim f_{m+1}, \dots, \sim f_n)$

That is, the probability that the message is ham, given that the message contains some features and does not contain others.

$P(\text{spam} \mid f_1, \dots, f_m, \sim f_{m+1}, \dots, \sim f_n)$

That is, the probability that the message is spam, given that the message contains some features and does not contain others.

That notation is cumbersome. So let's abbreviate " $f_1, \dots, f_m, \sim f_{m+1}, \dots, \sim f_n$ " as " $f_1, \dots, \sim f_n$ ". So we wish to compute:

(2)

$P(\text{ham} \mid f_1, \dots, \sim f_n)$

$P(\text{spam} \mid f_1, \dots, \sim f_n)$

But we don't know how to compute those probabilities. So let's apply some mathematics...

Since the AND operator is commutative:

$$P(f_1, \dots, \sim f_n \text{ AND ham}) = P(\text{ham AND } f_1, \dots, \sim f_n)$$

By the multiplicative law of probability, $P(x \text{ AND } y) = P(x) P(y | x)$. And so:

$$P(f_1, \dots, \sim f_n) P(\text{ham} | f_1, \dots, \sim f_n) = P(\text{ham}) P(f_1, \dots, \sim f_n | \text{ham})$$

Dividing both sides of the equation by $P(f_1, \dots, \sim f_n)$, we get **Bayes' Rule**:

$$P(\text{ham} | f_1, \dots, \sim f_n) = P(\text{ham}) P(f_1, \dots, \sim f_n | \text{ham}) / P(f_1, \dots, \sim f_n)$$

Similarly:

$$\begin{aligned} P(f_1, \dots, \sim f_n \text{ AND spam}) &= P(\text{spam AND } f_1, \dots, \sim f_n) \\ P(f_1, \dots, \sim f_n) P(\text{spam} | f_1, \dots, \sim f_n) &= P(\text{spam}) P(f_1, \dots, \sim f_n | \text{spam}) \\ P(\text{spam} | f_1, \dots, \sim f_n) &= P(\text{spam}) P(f_1, \dots, \sim f_n | \text{spam}) / P(f_1, \dots, \sim f_n) \end{aligned}$$

Substituting for the above expressions, we wish to compute:

$$\frac{P(\text{ham}) P(f_1, \dots, \sim f_n | \text{ham})}{P(f_1, \dots, \sim f_n)}$$
$$\frac{P(\text{spam}) P(f_1, \dots, \sim f_n | \text{spam})}{P(f_1, \dots, \sim f_n)} \tag{3}$$

We don't know $P(\text{ham})$ or $P(\text{spam})$, and we never will. We'll need to guess them, based upon our perception of the proportion of our e-mail that is, in fact, spam.

$P(f_1, \dots, \sim f_n | \text{ham})$ and $P(f_1, \dots, \sim f_n | \text{spam})$ would be difficult to compute. Doing so would be possible only if we have examples of messages that contain (or do not contain) every combination of features. That would imply that we need a very large number of examples.

So we must make a simplifying assumption. Let's assume that f_1, \dots, f_n are independent. Then:

$$P(f_1, \dots, \sim f_n) = P(f_1) \dots P(\sim f_n)$$

and

$$\begin{aligned} P(f_1, \dots, \sim f_n \mid \text{ham}) &= P(f_1 \mid \text{ham}) \dots P(\sim f_n \mid \text{ham}) \\ P(f_1, \dots, \sim f_n \mid \text{spam}) &= P(f_1 \mid \text{spam}) \dots P(\sim f_n \mid \text{spam}) \end{aligned}$$

Substituting into expressions (3), we wish to compute:

$\frac{P(\text{ham}) P(f_1 \mid \text{ham}) \dots P(\sim f_n \mid \text{ham})}{P(f_1) \dots P(\sim f_n)}$ $\frac{P(\text{spam}) P(f_1 \mid \text{spam}) \dots P(\sim f_n \mid \text{spam})}{P(f_1) \dots P(\sim f_n)}$	(4)
--	-----

We don't know how to compute $P(f_i)$ or $P(\sim f_i)$ for any i . So we need to transform the denominators.

Certainly, any message is either ham or spam. So:

$$P(\text{ham} \mid f_1, \dots, \sim f_n) + P(\text{spam} \mid f_1, \dots, \sim f_n) = 1$$

After substituting expressions (4) into that equation:

$$\frac{P(\text{ham}) P(f_1 \mid \text{ham}) \dots P(\sim f_n \mid \text{ham})}{P(f_1) \dots P(\sim f_n)} + \frac{P(\text{spam}) P(f_1 \mid \text{spam}) \dots P(\sim f_n \mid \text{spam})}{P(f_1) \dots P(\sim f_n)} = 1$$

After multiplying both sides of the equation by $P(f_1) \dots P(\sim f_n)$:

$$P(f_1) \dots P(\sim f_n) = P(\text{ham}) P(f_1 \mid \text{ham}) \dots P(\sim f_n \mid \text{ham}) + P(\text{spam}) P(f_1 \mid \text{spam}) \dots P(\sim f_n \mid \text{spam})$$

Substituting into expressions (4), we wish to compute:

(5)

$$P(\text{ham}) P(f_1 | \text{ham}) \dots P(\sim f_n | \text{ham})$$

$$P(\text{ham}) P(f_1 | \text{ham}) \dots P(\sim f_n | \text{ham}) + P(\text{spam}) P(f_1 | \text{spam}) \dots P(\sim f_n | \text{spam})$$

$$P(\text{spam}) P(f_1 | \text{spam}) \dots P(\sim f_n | \text{spam})$$

$$P(\text{ham}) P(f_1 | \text{ham}) \dots P(\sim f_n | \text{ham}) + P(\text{spam}) P(f_1 | \text{spam}) \dots P(\sim f_n | \text{spam})$$

As noted previously, the sum of those two expressions is 1. So it would be sufficient to compute only one of them; we easily could derive the other. So, let's compute only the second expression:

(6)

$$P(\text{spam}) P(f_1 | \text{spam}) \dots P(\sim f_n | \text{spam})$$

$$P(\text{ham}) P(f_1 | \text{ham}) \dots P(\sim f_n | \text{ham}) + P(\text{spam}) P(f_1 | \text{spam}) \dots P(\sim f_n | \text{spam})$$

Note that we can estimate all components of that expression. Specifically, we can estimate:

- $P(\text{ham})$ based upon our perception of the proportion of our e-mail that is, in fact, ham.
- $P(\text{spam})$ based upon our perception of the proportion of our e-mail that is, in fact, spam. Note that $P(\text{ham}) + P(\text{spam}) = 1$.
- $P(f_i | \text{ham})$ by examining many ham messages, and determining the proportion of them that contain feature f_i .
- $P(f_i | \text{spam})$ by examining many spam messages, and determining the proportion of them that contain feature f_i .
- $P(\sim f_i | \text{ham})$ by examining many ham messages, and determining the proportion of them that do not contain feature f_i . Or we could compute it as $1 - P(f_i | \text{ham})$.
- $P(\sim f_i | \text{spam})$ by examining many spam messages, and determining the proportion of them that do not contain feature f_i . Or we could compute it as $1 - P(f_i | \text{spam})$.

So, in *theory*, we can use expression (6) to produce the results that we wish.

However, in *practice* the products may become very small (i.e. close to 0), and thus cause may cause floating-point underflow. So, relying upon the equality:

$$x * y = \exp(\log(x) + \log(y))$$

let's compute sums of logarithms instead of products. That is, we wish to compute:

(7)
$\frac{\exp(\log(P(\text{spam})) + \log(P(f_1 \text{spam})) + \dots + \log(P(\sim f_n \text{spam})))}{\exp(\log(P(\text{ham})) + \log(P(f_1 \text{ham})) + \dots + \log(P(\sim f_n \text{ham}))) + \exp(\log(P(\text{spam})) + \log(P(f_1 \text{spam})) + \dots + \log(P(\sim f_n \text{spam})))}$

The logarithms will be negative. So the sums of the logarithms may be large negative numbers (i.e. far from 0). So applying the exp operation to those sums may cause precisely the same floating-point underflow that motivated us to use logarithms in the first place.

Relying upon this equality:

$$\frac{\exp(\log(x))}{\exp(\log(x)) + \exp(\log(y))} = \frac{\exp(\log(x) + k)}{\exp(\log(x) + k) + \exp(\log(y) + k)}$$

let's add some number k to each sum-of-logs before applying the exp operator. So we wish to compute:

(8)
$\frac{\exp(\log(P(\text{spam})) + \log(P(f_1 \text{spam})) + \dots + \log(P(\sim f_n \text{spam})) + k)}{\exp(\log(P(\text{ham})) + \log(P(\sim f_1 \text{ham})) + \dots + \log(P(\sim f_n \text{ham})) + k) + \exp(\log(P(\text{spam})) + \log(P(f_1 \text{spam})) + \dots + \log(P(\sim f_n \text{spam})) + k)}$

For k, a good choice would be:

$$-(\log(P(\text{spam})) + \log(P(f_1 | \text{spam})) + \dots + \log(P(\sim f_n | \text{spam})))$$

that is,

$$-\log(P(\text{spam})) - \log(P(f_1 | \text{spam})) - \dots - \log(P(\sim f_n | \text{spam}))$$

It's a good choice because it makes the numerator equal $\exp(0)$, that is, 1. It also makes the second term of the denominator equal $\exp(0)$, that is, 1. So, we wish to compute:

(9)

1

$$\exp(\log(P(\text{ham})) + \log(P(f_1 | \text{ham})) + \dots + \log(P(\sim f_n | \text{ham})) - \log(P(\text{spam})) - \log(P(f_1 | \text{spam})) - \dots - \log(P(\sim f_n | \text{spam}))) + 1$$

Note that only one sum-of-logs remains, and it will not evaluate to a large negative number (i.e. far from 0). So applying the exp operation to that sum will not cause floating-point underflow. Thus we have an expression that is correct in theory and in practice.

Copyright © 2005 by Robert M. Dondero, Jr.