

# Chapter 1

## Text Classification

### 1.1 The Bag of Words

We normally view text as a sequence of words, or we can impose additional structure on it like syntax. Or, we can dissolve some structure, namely, the order of the words. This gives a *bag of words*, which simply counts how many times each word occurs in a sentence (or document).

the cat sat on the mat → {cat, mat, on, sat, the, the}

Why? First, this representation is computationally extremely easy to work with. Second, this representation melts a text down into bits of meaning and serves as a crude way of capturing what the text is “about.” Crude, but often effective.

**Tokenization** We assume that all the sentences/documents have been tokenized so that the word boundaries are unambiguous. A commonly used English tokenizer is part of Stanford CoreNLP,<sup>1</sup> and LDC has a simple rule-based tokenizer.<sup>2</sup> Both are implemented/wrapped in Python by NLTK.<sup>3</sup>

**Stemming** It is also common in bag-of-word approaches to do morphological *stemming*, that is, removing affixes like *-ed*, *-ing*, etc. A classic stemmer for English is the Porter stemmer,<sup>4</sup> and another is the Lancaster (Paice/Husk) stemmer.<sup>5</sup> Both are implemented/wrapped in Python by NLTK.<sup>6</sup>

**Stop Words** Finally, it's also common to remove high-frequency but low-content words, like (Manning, Raghavan, and Schütze, 2008, p. 26):

a an and are as at be by for from has he in  
is it its of on that the to was were will with

---

<sup>1</sup><http://nlp.stanford.edu/software/corenlp.shtml>

<sup>2</sup><http://www.cis.upenn.edu/~treebank/tokenizer.sed>

<sup>3</sup><http://www.nltk.org/api/nltk.tokenize.html>

<sup>4</sup><http://tartarus.org/martin/PorterStemmer/>

<sup>5</sup><http://www.comp.lancs.ac.uk/computing/research/stemming/index.htm>

<sup>6</sup><http://www.nltk.org/api/nltk.stem.html>

## 1.2 Classification Problem

Suppose we have a collection of documents in different classes, and we want to learn to classify new documents. For example:

- We have a collection of e-mails that are labeled either as “spam” or “ham” and we want to learn to classify new e-mails.
- We have some texts whose author is known (e.g., Alexander Hamilton, James Madison, and John Jay) and want to classify anonymous texts (the Federalist Papers).
- We have a database of product reviews together with star ratings, and we want to be able to predict star ratings based on reviews alone.
- We have samples of (transcribed) speech from children diagnosed as autistic or not autistic, and we want to automatically diagnose other children using their speech.
- We have tweets that are known to be in various languages and want to automatically identify the language of new tweets.

**Question 1.** What are some other possible applications?

We’ll continue working with the example of spam filtering, but just remember that this is only one of many possible applications. This is a fairly easy problem, but occasionally difficult even for humans. For example:

From: fas@nsf.gov  
Subject: Payment  
To: chiang@isi.edu

DFM has approved your requested payment and has asked the U.S. Treasury to issue a payment to you within the next 4 working days. This payment is being sent directly to the Bank or Financial Institution identified by you for this purpose. If you are an NSF employee, this payment is being sent directly to the bank/financial institution where your bi-weekly pay is being deposited.

This payment for 560.00 is 1099 reportable if it totals \$600 or more for the year (i.e. You will receive an IRS 1099-Misc form from NSF)P131318.

Head, Accounts Payable Section.  
-----

More formally, we are given documents  $d_1, \dots, d_n$  together with their correct classes  $k_1, \dots, k_n$ . We want to learn a model  $P(k | d)$ , where  $k$  is a class and  $d$  is a document, and given a new document  $d$ , we want to be able to find, with high accuracy,

$$k^* = \arg \max_k P(k | d). \quad (1.1)$$

## 1.3 Naïve Bayes

**Question 2.** What problem would you run into if you used the model

$$P(k | d) = \frac{c(k, d)}{c(d)}?$$

What might you do to repair the model?

This route turns out to be difficult (though we will return to it below, in Section 2.3). Instead of thinking about how to classify a document, let's think about how the document came to be. First, someone decided to write an e-mail; he was either a spammer or a "hammer." Then, this person authored a document. More formally:

$$\arg \max_k P(k | d) = \arg \max_k P(d)P(k | d) \quad (1.2)$$

$$= \arg \max_k P(k, d) \quad (1.3)$$

$$= \arg \max_k P(k)P(d | k). \quad (1.4)$$

The advantage of thinking about it this way is that it is much easier to write down a model for  $P(d | k)$  than for  $P(k | d)$ .

### 1.3.1 Model

We'll consider just the simplest, which is called a *naïve Bayes* classifier and looks like this:

$$P(k, d) = p(k) \prod_{w \in d} p(w | k), \quad (1.5)$$

where the  $p(k)$  and  $p(w | k)$  are the parameters of the model. (Let's adopt the convention that  $P(\text{something})$  and  $p(\text{something})$  both denote the probability of something, but  $P(\text{something})$  might be defined in terms of other probabilities, whereas  $p(\text{something})$  is a parameter of the model that needs to be estimated.)

Naïve Bayes is called naïve because it naïvely assumes that all the words in a document are independent of each other. All it captures is that spammers are more likely to use certain words and hammers are more likely to use certain words.

It is possible to count other things besides words. For example, if we're trying to classify a document as English or French, then the presence of an  $\acute{e}$  is a pretty good sign that it's in french. So character probabilities like  $p(\acute{e} | k)$  might be useful in addition to word probabilities. Or, if we're trying to predict the positivity or negativity of a product review, then the occurrence of the word *bad* might be good sign of a negative review, unless it is immediately preceded by the word *not*. So probabilities of *bigrams* like  $p(\text{not bad} | k)$  are important. Note that if we do this, however, the distribution  $P(k, d)$  no longer sums to one, because there's no such thing as a document that contains the words {dog} but the characters {c, a, t}. So, if we only sum over feasible documents  $d$ , then  $P(k, d)$  sums to less than one, that is, it is *deficient*, which is in general not considered a good thing. This is even more naïve, but might work okay in practice.

**Question 3.** What are some other features that might be helpful?

### 1.3.2 Training

Training the model, or estimating the parameters  $p(k)$  and  $p(w | k)$ , is easy. It's just:

$$\begin{aligned} p(k) &= \frac{c(k)}{\sum_k c(k)} \\ p(w | k) &= \frac{c(k, w)}{\sum_{w'} c(k, w')}. \end{aligned} \tag{1.6}$$

This is known as *relative frequency estimation* (or “count and divide”). It's so intuitive that it may seem odd to give it a name, but it's worth dwelling a bit on why we estimate probabilities this way.

Consider just  $p(k)$ . Suppose we have just two theories, or *models*, about how prevalent spam is. Model 1, or  $m_1$ , says that 10% of e-mails are spam, whereas Model 2, or  $m_2$ , says that 90% of e-mails are spam.

$$\begin{aligned} p(\text{spam} | m_1) &= \frac{1}{10} & p(\text{spam} | m_2) &= \frac{9}{10} \\ p(\text{ham} | m_1) &= \frac{9}{10} & p(\text{ham} | m_2) &= \frac{1}{10} \end{aligned}$$

Before looking at the data, suppose that we think these two models are equally valid.

$$p(m_1) = \frac{1}{2} \qquad p(m_2) = \frac{1}{2}$$

Then, suppose we look at some data and see that, out of 10 e-mails, 7 are spam and 3 are ham. Now we want to know which model is better,  $m_1$  or  $m_2$ ?

$$\begin{aligned} P(m_1 | \text{data}) &= \frac{p(m_1)P(\text{data} | m_1)}{P(\text{data})} & P(m_2 | \text{data}) &= \frac{p(m_2)P(\text{data} | m_2)}{P(\text{data})} \\ &= \frac{\frac{1}{2} \left(\frac{1}{10}\right)^7 \left(\frac{9}{10}\right)^3}{P(\text{data})} & &= \frac{\frac{1}{2} \left(\frac{9}{10}\right)^7 \left(\frac{1}{10}\right)^3}{P(\text{data})} \\ &= \frac{9^3}{2 \cdot 10^{10} \cdot P(\text{data})} & &= \frac{9^7}{2 \cdot 10^{10} \cdot P(\text{data})} \end{aligned}$$

Since  $P(m_2 | \text{data})$  is bigger, we can conclude that  $m_2$  is the better model.

Now suppose that we compare not two models, but an infinite number of models,

$$\begin{aligned} P(\text{spam} | \theta) &= \theta \\ P(\text{ham} | \theta) &= 1 - \theta \end{aligned}$$

for all  $\theta \in [0, 1]$ . The same reasoning still holds: we want to find the model that maximizes

$$\underbrace{P(\theta | \text{data})}_{\text{posterior}} = \frac{\overbrace{P(\theta) P(\text{data} | \theta)}^{\text{prior likelihood}}}{\underbrace{P(\text{data})}_{\text{evidence}}}. \tag{1.7}$$

Note that the denominator is independent of  $\theta$ , and if we assume that the prior  $P(\theta)$  is uniform, then the only factor that matters is the likelihood. We therefore call this *maximum likelihood estimation*.

Going back to the full naïve Bayes model: we want to maximize the likelihood, which is

$$L = \prod_i P(k_i, d_i) \quad (1.8)$$

$$= \prod_i p(k_i) \prod_{w \in d_i} p(w | k_i). \quad (1.9)$$

And this maximization problem has a closed-form solution, namely (1.6). See Section 1.3.7 for more details on how to derive that.

### 1.3.3 Classification

Once we've trained the model, finding the most probable class of a new document is also easy:

$$k^* = \arg \max_k P(k | d) \quad (1.10)$$

$$= \arg \max_k P(k, d) \quad (1.11)$$

$$= \arg \max_k p(k) \prod_{w \in d} p(w | k). \quad (1.12)$$

### 1.3.4 Smoothing

What happens if we encounter a word we've never seen before? Then we would have  $p(w | k) = 0$ , and therefore  $P(k, d) = 0$ , for all  $k$ . The presence of a single unknown word zeros out the whole document probability and makes it impossible to choose a class.

The standard solution is *smoothing*. We'll talk about the simplest smoothing method now, and more advanced smoothing methods later (when we talk about language models). In *add-one smoothing*, invented by Laplace in the 18th century, we add one to the count of every event, including unseen events. Thus we would estimate  $p(w | k)$  as:

$$p(w | k) = \frac{c(k, w) + 1}{\sum_{w'} c(k, w') + |V|} \quad (1.13)$$

where  $|V|$  is the size of the vocabulary, including unseen words. But what if we don't know how many unseen words there are? A typical (though not entirely correct) thing to do is to set  $|V|$  to the number of seen word types, plus one for a generic unseen word.

We can also add any value  $\delta > 0$  to the counts instead of one:

$$p(w | k) = \frac{c(k, w) + \delta}{\sum_{w'} c(k, w') + |V|\delta}. \quad (1.14)$$

The increment that we add (whether 1 or  $\delta$ ) is known as a *pseudocount*. Typically,  $\delta$  would be set by trial and error. If you have a good reason to, you can also add a different pseudocount to different word types.

These methods are not hacks; they can all be thought of as different choices of the prior probability that we saw in Equation (1.7).

**Question 4.** If our training data is:

ham	spam	spam
please pass on to your groups	we deliver to your door within 24 hours	please update your account details with citibank

and we train with add-one smoothing, what are  $P(\text{spam} | d)$  and  $P(\text{ham} | d)$  for the document below?

please forward to  
your groups

### 1.3.5 Log-probabilities

Another practical problem we run into when multiplying together many probabilities is that the numbers quickly become very small. An IEEE double only goes down to  $10^{-308}$ , and the probability of a long document could easily be less than that, causing underflow. To fix this problem, the typical solution is to use *log-probabilities*, that is, instead of storing the probabilities themselves, storing their logarithms.

Then, instead of multiplying probabilities, we add log-probabilities (because  $\log pq = \log p + \log q$ ). In the old days, this was seen as a bonus, because addition used to be faster than multiplication. Also, comparing log-probabilities is trivial, since  $\log p > \log q$  if and only if  $p > q$ .

If we want to compute  $\log(p + q)$  given  $x = \log p$  and  $y = \log q$ , we need to be careful. We can't do this:

$$\log(p + q) = \log(\exp \log p + \exp \log q)$$

because either of the exp's might cause an underflow. Instead, assume that  $p > q$ ; if not, swap them. Then, observe that:

$$\begin{aligned} \log(p + q) &= \log(\exp \log p + \exp \log q) \\ &= \log\left(\exp \log p \left(1 + \frac{\exp \log q}{\exp \log p}\right)\right) \\ &= \log p + \log\left(1 + \exp(\log q - \log p)\right) \end{aligned}$$

Now, the exp could still cause an underflow, but the underflow is harmless. (Why?) For an extra little boost in accuracy, you can use the `log1p` function, found in nearly all standard libraries, which computes  $\log(1 + x)$  but is accurate for small  $x$ .

### 1.3.6 Experiment

Let's tackle a more difficult classification task than spam filtering: Given a blog post, can you predict whether the author is male or female? A dataset is provided by **mukherjee+liu:2010** consisting of 3227 posts.<sup>7</sup>

We split the data into two parts, *training* and *test data*. We train the model on the training data, then test the model on the test data by classifying the documents and measuring the percent accuracy. We could have also set aside a third part, called *development data*, as a proxy for the test data while we are fiddling with the model. This way, we don't artificially inflate our score on the test data by lucky modifications. Here, we used 80% for training data and 20% for testing. Here's the first male-authored blog post:

<sup>7</sup><http://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html>

long time no see . like always i was rewriting it from scratch a couple of times . but nevertheless it's still java and now it uses metropolis sampling to help that poor path tracing converge . btw . i did mlt on yesterday evening after 2 beers ( it had to be ballmer peak<sup>8</sup> ) . although the implementation is still very fresh it easily outperforms standard path tracing , what is to be seen especially when difficult caustics are involved . i've implemented spectral rendering too , it was very easy actually , cause all computations on wavelengths are linear just like rgb . but then i realised that even if it does feel more physically correct to do so , whats the point ? 3d applications are operating in rgb color space , and because i cant represent a rgb color as spectrum interchangeably i have to approximate it , so as long as i'm not running a physical simulation or something i don't see the benefits ( please correct me if i'm wrong ) , thus i abandoned that .

And the first female-authored blog post:

liam is nothing like natalie . natalie never went through draws or cabinets . she was always so good . liam is a bit more adventurous . and he always picks the cabinets that are the hardest to put back together . hubby just started nutrisystem and he has his own little section of the kitchen . well , liam doesn't like the order of it all and trashes my set up at least 10 times a day . yes , i need to buy one of those locks for the doors . we have one for the chemicals . but none on the food . my mistake for sure ! in the meantime , i better document it because before i know it this boy will be in college and i will be missing cleaning up after him . yeah . . that's not me . but i still decided that i wanted to start doing yoga . i've been thinking about it for quite some time . so , this weekend i bought a pair of work out shorts and a tank and off i went . the only class that was close to me was birkham yoga , the hot one . the 100 degree sweaty room one . i . almost . died . the yoga itself was wonderful . i will continue that for sure . but the heat of that room was unbearable . the best part is , the teacher specifically came up to me and said ,

When we train the naïve Bayes classifier on the training data (with add-one smoothing) and test it on the test data, we get an accuracy of 65.1%. Note that we could have gotten 50% accuracy simply by guessing randomly. So our performance is better than random, but not all that much better.

What happens if we count both words and bigrams (pairs of consecutive words)? When counting bigrams, we add a fake word `<s>` at the beginning and a fake word `</s>` at the end. This gives an accuracy of 66.5%.

### 1.3.7 Optional: Deriving relative frequency estimation

Here's how to derive the relative frequency estimate, taking  $p(k)$  as an example. We want to maximize the likelihood  $L$ , or, equivalently, the log-likelihood

$$\log L = \sum_k c(k) \log p(k) \quad (1.15)$$

subject to the constraint

$$\sum_k p(k) = 1. \quad (1.16)$$

---

<sup>8</sup><http://xkcd.com/323/>

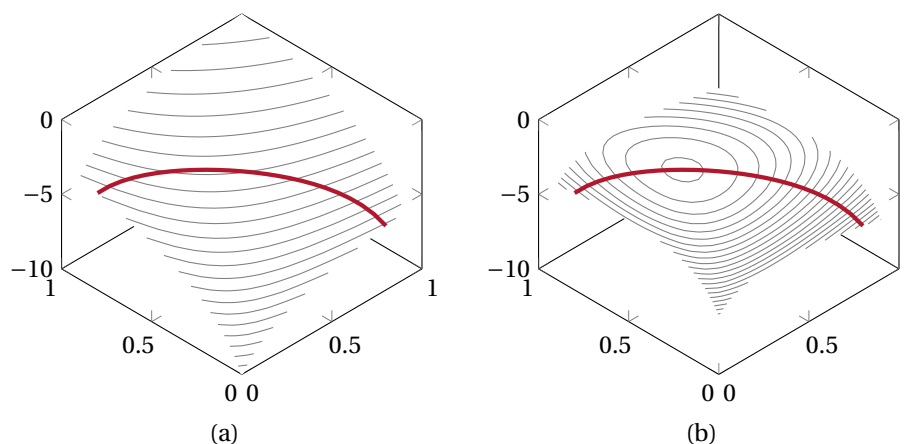


Figure 1.1: Relative frequency estimation. (a) We need to maximize the likelihood subject to the sum-to-one constraint (red line). (b) There is a value of  $\lambda$  that makes the Lagrangian have a maximum that does satisfy the sum-to-one constraint (red line).

We can't just set the partial derivatives to zero because, although the surface defined by  $\log L$  will be flat at its maximum along the line defined by (1.16), it will in general be sloped in the direction perpendicular to the line (see Figure 1.1). So we create a new function  $\mathcal{L}$ , called the *Lagrangian*, that can level the surface out:

$$\mathcal{L} = \sum_k c(k) \log p(k) - \lambda \left( \sum_k p(k) - 1 \right). \quad (1.17)$$

This is our original objective function plus a new term. This term is zero along the line defined by (1.16), and sloped in the direction perpendicular to the line, with a slope determined by the *Lagrange multiplier*  $\lambda$ . At the maximum we are seeking, there is some value of  $\lambda$  that makes all the partial derivatives with respect to the parameters  $p(k)$  zero.

The partial derivatives of  $\mathcal{L}$  are:

$$\frac{\partial \mathcal{L}}{\partial p(k)} = \frac{c(k)}{p(k)} - \lambda, \quad (1.18)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = - \left( \sum_k p(k) - 1 \right). \quad (1.19)$$

Note that the partial derivative with respect to  $\lambda$  is zero just in case the constraint (1.16) is satisfied. So the maximum we are seeking is the point where all the partial derivatives are zero. If we set them to zero and solve for the  $p(k)$ , we get:

$$p(k) = \frac{c(k)}{\sum_k c(k)}. \quad (1.20)$$

Further work (!) would be needed to verify that this is in fact the global maximum.



## References

Manning, Christopher D., Prabhakar Raghavan, and Hinrich Schütze (2008). *Introduction to Information Retrieval*. Cambridge University Press. URL: <http://www-nlp.stanford.edu/IR-book/>.