

almost always do better than GDA. For this reason, in practice logistic regression is used more often than GDA. (Some related considerations about discriminative vs. generative models also apply for the Naive Bayes algorithm that we discuss next, but the Naive Bayes algorithm is still considered a very good, and is certainly also a very popular, classification algorithm.)

## 2 Naive Bayes

In GDA, the feature vectors  $x$  were continuous, real-valued vectors. Lets now talk about a different learning algorithm in which the  $x_i$ 's are discrete-valued.

For our motivating example, consider building an email spam filter using machine learning. Here, we wish to classify messages according to whether they are unsolicited commercial (spam) email, or non-spam email. After learning to do this, we can then have our mail reader automatically filter out the spam messages and perhaps place them in a separate mail folder. Classifying emails is one example of a broader set of problems called **text classification**.

Lets say we have a training set (a set of emails labeled as spam or non-spam). We'll begin our construction of our spam filter by specifying the features  $x_i$  used to represent an email.

We will represent an email via a feature vector whose length is equal to the number of words in the dictionary. Specifically, if an email contains the  $i$ -th word of the dictionary, then we will set  $x_i = 1$ ; otherwise, we let  $x_i = 0$ . For instance, the vector

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

is used to represent an email that contains the words “a” and “buy,” but not “aardvark,” “aardwolf” or “zygmurgy.”<sup>2</sup> The set of words encoded into the

---

<sup>2</sup>Actually, rather than looking through an english dictionary for the list of all english words, in practice it is more common to look through our training set and encode in our feature vector only the words that occur at least once there. Apart from reducing the number of words modeled and hence reducing our computational and space requirements,

feature vector is called the **vocabulary**, so the dimension of  $x$  is equal to the size of the vocabulary.

Having chosen our feature vector, we now want to build a discriminative model. So, we have to model  $p(x|y)$ . But if we have, say, a vocabulary of 50000 words, then  $x \in \{0, 1\}^{50000}$  ( $x$  is a 50000-dimensional vector of 0's and 1's), and if we were to model  $x$  explicitly with a multinomial distribution over the  $2^{50000}$  possible outcomes, then we'd end up with a  $(2^{50000} - 1)$ -dimensional parameter vector. This is clearly too many parameters.

To model  $p(x|y)$ , we will therefore make a very strong assumption. We will assume that the  $x_i$ 's are conditionally independent given  $y$ . This assumption is called the **Naive Bayes (NB) assumption**, and the resulting algorithm is called the **Naive Bayes classifier**. For instance, if  $y = 1$  means spam email; “buy” is word 2087 and “price” is word 39831; then we are assuming that if I tell you  $y = 1$  (that a particular piece of email is spam), then knowledge of  $x_{2087}$  (knowledge of whether “buy” appears in the message) will have no effect on your beliefs about the value of  $x_{39831}$  (whether “price” appears). More formally, this can be written  $p(x_{2087}|y) = p(x_{2087}|y, x_{39831})$ . (Note that this is *not* the same as saying that  $x_{2087}$  and  $x_{39831}$  are independent, which would have been written “ $p(x_{2087}) = p(x_{2087}|x_{39831})$ ”; rather, we are only assuming that  $x_{2087}$  and  $x_{39831}$  are conditionally independent *given*  $y$ .)

We now have:

$$\begin{aligned} p(x_1, \dots, x_{50000}|y) &= p(x_1|y)p(x_2|y, x_1)p(x_3|y, x_1, x_2) \cdots p(x_{50000}|y, x_1, \dots, x_{49999}) \\ &= p(x_1|y)p(x_2|y)p(x_3|y) \cdots p(x_{50000}|y) \\ &= \prod_{i=1}^n p(x_i|y) \end{aligned}$$

The first equality simply follows from the usual properties of probabilities, and the second equality used the NB assumption. We note that even though the Naive Bayes assumption is an extremely strong assumptions, the resulting algorithm works well on many problems.

Our model is parameterized by  $\phi_{i|y=1} = p(x_i = 1|y = 1)$ ,  $\phi_{i|y=0} = p(x_i = 1|y = 0)$ , and  $\phi_y = p(y = 1)$ . As usual, given a training set  $\{(x^{(i)}, y^{(i)}); i =$

---

this also has the advantage of allowing us to model/include as a feature many words that may appear in your email (such as “cs229”) but that you won't find in a dictionary. Sometimes (as in the homework), we also exclude the very high frequency words (which will be words like “the,” “of,” “and,”; these high frequency, “content free” words are called **stop words**) since they occur in so many documents and do little to indicate whether an email is spam or non-spam.

$1, \dots, m\}$ , we can write down the joint likelihood of the data:

$$\mathcal{L}(\phi_y, \phi_{i|y=0}, \phi_{i|y=1}) = \prod_{i=1}^m p(x^{(i)}, y^{(i)}).$$

Maximizing this with respect to  $\phi_y, \phi_{i|y=0}$  and  $\phi_{i|y=1}$  gives the maximum likelihood estimates:

$$\begin{aligned}\phi_{j|y=1} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \\ \phi_{j|y=0} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \\ \phi_y &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m}\end{aligned}$$

In the equations above, the “ $\wedge$ ” symbol means “and.” The parameters have a very natural interpretation. For instance,  $\phi_{j|y=1}$  is just the fraction of the spam ( $y = 1$ ) emails in which word  $j$  does appear.

Having fit all these parameters, to make a prediction on a new example with features  $x$ , we then simply calculate

$$\begin{aligned}p(y = 1|x) &= \frac{p(x|y = 1)p(y = 1)}{p(x)} \\ &= \frac{(\prod_{i=1}^n p(x_i|y = 1))p(y = 1)}{(\prod_{i=1}^n p(x_i|y = 1))p(y = 1) + (\prod_{i=1}^n p(x_i|y = 0))p(y = 0)},\end{aligned}$$

and pick whichever class has the higher posterior probability.

Lastly, we note that while we have developed the Naive Bayes algorithm mainly for the case of problems where the features  $x_i$  are binary-valued, the generalization to where  $x_i$  can take values in  $\{1, 2, \dots, k_i\}$  is straightforward. Here, we would simply model  $p(x_i|y)$  as multinomial rather than as Bernoulli. Indeed, even if some original input attribute (say, the living area of a house, as in our earlier example) were continuous valued, it is quite common to **discretize** it—that is, turn it into a small set of discrete values—and apply Naive Bayes. For instance, if we use some feature  $x_i$  to represent living area, we might discretize the continuous values as follows:

Living area (sq. feet)	< 400	400-800	800-1200	1200-1600	>1600
$x_i$	1	2	3	4	5

Thus, for a house with living area 890 square feet, we would set the value of the corresponding feature  $x_i$  to 3. We can then apply the Naive Bayes

algorithm, and model  $p(x_i|y)$  with a multinomial distribution, as described previously. When the original, continuous-valued attributes are not well-modeled by a multivariate normal distribution, discretizing the features and using Naive Bayes (instead of GDA) will often result in a better classifier.

## 2.1 Laplace smoothing

The Naive Bayes algorithm as we have described it will work fairly well for many problems, but there is a simple change that makes it work much better, especially for text classification. Lets briefly discuss a problem with the algorithm in its current form, and then talk about how we can fix it.

Consider spam/email classification, and lets suppose that, after completing CS229 and having done excellent work on the project, you decide around June 2003 to submit the work you did to the NIPS conference for publication. (NIPS is one of the top machine learning conferences, and the deadline for submitting a paper is typically in late June or early July.) Because you end up discussing the conference in your emails, you also start getting messages with the word “nips” in it. But this is your first NIPS paper, and until this time, you had not previously seen any emails containing the word “nips”; in particular “nips” did not ever appear in your training set of spam/non-spam emails. Assuming that “nips” was the 35000th word in the dictionary, your Naive Bayes spam filter therefore had picked its maximum likelihood estimates of the parameters  $\phi_{35000|y}$  to be

$$\begin{aligned}\phi_{35000|y=1} &= \frac{\sum_{i=1}^m 1\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} = 0 \\ \phi_{35000|y=0} &= \frac{\sum_{i=1}^m 1\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} = 0\end{aligned}$$

I.e., because it has never seen “nips” before in either spam or non-spam training examples, it thinks the probability of seeing it in either type of email is zero. Hence, when trying to decide if one of these messages containing “nips” is spam, it calculates the class posterior probabilities, and obtains

$$\begin{aligned}p(y = 1|x) &= \frac{\prod_{i=1}^n p(x_i|y = 1)p(y = 1)}{\prod_{i=1}^n p(x_i|y = 1)p(y = 1) + \prod_{i=1}^n p(x_i|y = 0)p(y = 0)} \\ &= \frac{0}{0}.\end{aligned}$$

This is because each of the terms “ $\prod_{i=1}^n p(x_i|y)$ ” includes a term  $p(x_{35000}|y) = 0$  that is multiplied into it. Hence, our algorithm obtains 0/0, and doesn’t know how to make a prediction.

Stating the problem more broadly, it is statistically a bad idea to estimate the probability of some event to be zero just because you haven't seen it before in your finite training set. Take the problem of estimating the mean of a multinomial random variable  $z$  taking values in  $\{1, \dots, k\}$ . We can parameterize our multinomial with  $\phi_i = p(z = i)$ . Given a set of  $m$  independent observations  $\{z^{(1)}, \dots, z^{(m)}\}$ , the maximum likelihood estimates are given by

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\}}{m}.$$

As we saw previously, if we were to use these maximum likelihood estimates, then some of the  $\phi_j$ 's might end up as zero, which was a problem. To avoid this, we can use **Laplace smoothing**, which replaces the above estimate with

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\} + 1}{m + k}.$$

Here, we've added 1 to the numerator, and  $k$  to the denominator. Note that  $\sum_{j=1}^k \phi_j = 1$  still holds (check this yourself!), which is a desirable property since the  $\phi_j$ 's are estimates for probabilities that we know must sum to 1. Also,  $\phi_j \neq 0$  for all values of  $j$ , solving our problem of probabilities being estimated as zero. Under certain (arguably quite strong) conditions, it can be shown that the Laplace smoothing actually gives the optimal estimator of the  $\phi_j$ 's.

Returning to our Naive Bayes classifier, with Laplace smoothing, we therefore obtain the following estimates of the parameters:

$$\begin{aligned} \phi_{j|y=1} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 1\} + 2} \\ \phi_{j|y=0} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 0\} + 2} \end{aligned}$$

(In practice, it usually doesn't matter much whether we apply Laplace smoothing to  $\phi_y$  or not, since we will typically have a fair fraction each of spam and non-spam messages, so  $\phi_y$  will be a reasonable estimate of  $p(y = 1)$  and will be quite far from 0 anyway.)

## 2.2 Event models for text classification

To close off our discussion of generative learning algorithms, let's talk about one more model that is specifically for text classification. While Naive Bayes

as we've presented it will work well for many classification problems, for text classification, there is a related model that does even better.

In the specific context of text classification, Naive Bayes as presented uses the what's called the **multi-variate Bernoulli event model**. In this model, we assumed that the way an email is generated is that first it is randomly determined (according to the class priors  $p(y)$ ) whether a spammer or non-spammer will send you your next message. Then, the person sending the email runs through the dictionary, deciding whether to include each word  $i$  in that email independently and according to the probabilities  $p(x_i = 1|y) = \phi_{i|y}$ . Thus, the probability of a message was given by  $p(y) \prod_{i=1}^n p(x_i|y)$ .

Here's a different model, called the **multinomial event model**. To describe this model, we will use a different notation and set of features for representing emails. We let  $x_i$  denote the identity of the  $i$ -th word in the email. Thus,  $x_i$  is now an integer taking values in  $\{1, \dots, |V|\}$ , where  $|V|$  is the size of our vocabulary (dictionary). An email of  $n$  words is now represented by a vector  $(x_1, x_2, \dots, x_n)$  of length  $n$ ; note that  $n$  can vary for different documents. For instance, if an email starts with "A NIPS ...," then  $x_1 = 1$  ("a" is the first word in the dictionary), and  $x_2 = 35000$  (if "nips" is the 35000th word in the dictionary).

In the multinomial event model, we assume that the way an email is generated is via a random process in which spam/non-spam is first determined (according to  $p(y)$ ) as before. Then, the sender of the email writes the email by first generating  $x_1$  from some multinomial distribution over words ( $p(x_1|y)$ ). Next, the second word  $x_2$  is chosen independently of  $x_1$  but from the same multinomial distribution, and similarly for  $x_3$ ,  $x_4$ , and so on, until all  $n$  words of the email have been generated. Thus, the overall probability of a message is given by  $p(y) \prod_{i=1}^n p(x_i|y)$ . Note that this formula looks like the one we had earlier for the probability of a message under the multi-variate Bernoulli event model, but that the terms in the formula now mean very different things. In particular  $x_i|y$  is now a multinomial, rather than a Bernoulli distribution.

The parameters for our new model are  $\phi_y = p(y)$  as before,  $\phi_{i|y=1} = p(x_j = i|y = 1)$  (for any  $j$ ) and  $\phi_{i|y=0} = p(x_j = i|y = 0)$ . Note that we have assumed that  $p(x_j|y)$  is the same for all values of  $j$  (i.e., that the distribution according to which a word is generated does not depend on its position  $j$  within the email).

If we are given a training set  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$  where  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_{n_i}^{(i)})$  (here,  $n_i$  is the number of words in the  $i$ -training example),

the likelihood of the data is given by

$$\begin{aligned}\mathcal{L}(\phi, \phi_{i|y=0}, \phi_{i|y=1}) &= \prod_{i=1}^m p(x^{(i)}, y^{(i)}) \\ &= \prod_{i=1}^m \left( \prod_{j=1}^{n_i} p(x_j^{(i)} | y; \phi_{i|y=0}, \phi_{i|y=1}) \right) p(y^{(i)}; \phi_y).\end{aligned}$$

Maximizing this yields the maximum likelihood estimates of the parameters:

$$\begin{aligned}\phi_{k|y=1} &= \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\} n_i} \\ \phi_{k|y=0} &= \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\} n_i} \\ \phi_y &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m}.\end{aligned}$$

If we were to apply Laplace smoothing (which needed in practice for good performance) when estimating  $\phi_{k|y=0}$  and  $\phi_{k|y=1}$ , we add 1 to the numerators and  $|V|$  to the denominators, and obtain:

$$\begin{aligned}\phi_{k|y=1} &= \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 1\} n_i + |V|} \\ \phi_{k|y=0} &= \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 0\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 0\} n_i + |V|}.\end{aligned}$$

While not necessarily the very best classification algorithm, the Naive Bayes classifier often works surprisingly well. It is often also a very good “first thing to try,” given its simplicity and ease of implementation.