

# THE NAIVE BAYES CLASSIFIER

The naive Bayes classifier is another generative method for classification.

## The Naive Bayes Assumption

Let  $X = [X^{(1)} \dots X^{(d)}]^T \in \mathbb{R}^d$  denote the random feature vector in a classification problem, and  $Y$  the corresponding label.

The naive Bayes classifier assumes that, given  $Y$ ,

④  $X^{(1)}, \dots, X^{(d)}$  are \_\_\_\_\_.

Although this assumption is rarely met in practice, it can still lead to reasonable, and sometimes very good, classification performance.

The reason is that in classification, we really don't need to model the class-conditional densities well, just the decision boundary of the Bayes classifier.

## Estimation

The major advantage of NB is that we only  
② need to estimate \_\_\_\_\_ densities.

Let  $g_k(x)$  be the probability law (density or mass function) of  $X|Y=k$ ,  $k=1, \dots, K$

By the NB assumption,

$$\textcircled{c} \quad g_k(x) =$$

Let  $(x_i, y_i)_{i=1}^n$  be training data, and let

$$\hat{\pi}_k = \frac{|\{i: y_i = k\}|}{n}$$

$\hat{g}_k^{(j)}$  = estimate of  $g_k^{(j)}$  based  
on  $\{x_i^{(j)}: y_i = k\}$

Then the NB classifier is

③

So it remains to determine  $\hat{g}_k^{(j)}$ .

Another advantage of NB is that it easily handles the case where some  $X^{(j)}$  are continuous and others are discrete.

Continuous  $X^{(j)} | Y=k$

- Gaussian MLE
- kernel density estimate
- quantize to discrete variable

Discrete  $X^{(j)} | Y=k$

Suppose that given  $Y=k$ ,  $X^{(j)}$  takes on the values  $z_1, \dots, z_L$ . Denote

$$n_{\mathbf{k}} = |\{i : y_i = k\}|$$

$$n_{\mathbf{k}l}^{(j)} = |\{i : y_i = k \wedge x_i^{(j)} = z_l\}|$$

Then the natural (and maximum likelihood) estimate of  $\Pr \{X^{(j)} = z_l | Y=k\}$  is

Note:  $z_l$  depends on  $j, k$ , but this dependence is omitted for simplicity.

That is,

$$\hat{g}_k^{(j)}(z_e) = \frac{n_{ke}^{(j)}}{n_k}$$

It is possible that when we apply NB to a test pattern,  $X, X^{(i)}$  may take on a value  $z'$  not observed in the training data. In this case, the MLE  $\hat{g}_k$  is undefined, which is undesirable.

Example] Text classification

$X^{(j)}$  = # of times word  $j$  occurs in a document,  
 $j=1, \dots, d$ , where  $d$  = total # of words in training data.  
In general,  $d <$  # words in the vocabulary.

To avoid this problem it is  $c$  to estimate

Ⓕ  $\hat{g}_k^{(j)}(z_e) =$

which corresponds to a Bayesian estimate (of multinomial parameters with a Dirichlet prior)

Key

A. independent

B. scalar / univariate

$$C. g_k(x) = \prod_{j=1}^d g_k^{(j)}(x^{(j)})$$

$$D. \hat{f}(x) = \arg \max_{k=1, \dots, K} \hat{\pi}_k \hat{g}_k(x)$$

$$E. n_{ke}^{(j)} / n_k$$

$$F. \frac{n_{ke}^{(j)} + 1}{n_k + L} \quad (\text{note } L \text{ depends on } j, k)$$