# THE NAIVE BAYES CLASSIFIER

The naive Bayes classifier is another generative method for classification.

> The Naive Bayes Assumption

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

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

(A) $X^{(1)}, \ldots, X^{(d)}$ 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

(B) need to estimate _____ densities.

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

By the NB assumption,

(C)
$$g_k(x) =$$

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

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

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

Then the NB classifier is

(D)

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 $z_1, \cdots, z_L$. Denote

> Note: $z_\ell$ depends on $j, k$, but this dependence is omitted for simplicity.

$$n_k = |\{i : y_i = k\}|$$

$$n_{k\ell}^{(j)} = |\{i : y_i = k \wedge x_i^{(j)} = z_\ell\}|$$

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

Ⓔ

That is,

$$\hat{g}_k^{(j)}(z_\ell) = \frac{n_{k\ell}^{(j)}}{n_k}.$$

It is possible that when we apply NB to a test pattern $X$, $X^{(j)}$ may take on a value $z'$ not observed in the training data for some class. Then $\hat{g}_k^{(j)}(z') = 0$ and hence $\hat{g}_k = 0$, so that class will never be predicted. This is undesirable.

## Example | Document classification

Suppose you wish to classify documents, e.g. $1 = $ politics, $2 = $ sports, $3 = $ finance, etc. One simple representation of a document is

$$X = [X^{(1)} \cdots X^{(d)}]^T \quad \text{where } d \text{ is the}$$

number of words in the vocabulary, and

$$X^{(j)} = \begin{cases} 1 & \text{if } j^{th} \text{ word occurs in document} \\ 0 & \text{otherwise.} \end{cases}$$

It could happen that all training documents contain the word "ball." But we don't want to exclude the possibility that some other sports article does not contain this word. //

To avoid this problem, it is common to estimate

$$\hat{q}_k^{(j)}(z_\ell) =$$

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

A. independent

B. scalar / univariate

C. $\quad g_k(x) = \prod\limits_{j=1}^{d} g_k^{(j)}(x^{(j)})$

D.

$$\hat{f}(x) = \underset{k=1,\dots,K}{\arg\max} \quad \hat{\pi}_k \, \hat{g}_k(x)$$

E. $\quad n_{k\ell}^{(j)} / n_k$

F. $\quad \dfrac{n_{k\ell}^{(j)} + 1}{n_k + L}$ $\qquad$ (note $L$ depends on $j, k$)