

As this is valid for any Borel-measurable set  $C$ , the distribution of  $(X, Y)$  is determined by  $(\mu, \eta)$ . The function  $\eta$  is sometimes called the *a posteriori probability*.

Any function  $g : \mathcal{R}^d \rightarrow \{0, 1\}$  defines a *classifier* or a *decision function*. The error probability of  $g$  is  $L(g) = \mathbf{P}\{g(X) \neq Y\}$ . Of particular interest is the Bayes decision function

$$g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

This decision function minimizes the error probability.

**Theorem 2.1.** For any decision function  $g : \mathcal{R}^d \rightarrow \{0, 1\}$ ,

$$\mathbf{P}\{g^*(X) \neq Y\} \leq \mathbf{P}\{g(X) \neq Y\},$$

that is,  $g^*$  is the optimal decision.

**PROOF.** Given  $X = x$ , the conditional error probability of any decision  $g$  may be expressed as

$$\begin{aligned} \mathbf{P}\{g(X) \neq Y | X = x\} &= 1 - \mathbf{P}\{Y = g(X) | X = x\} \\ &= 1 - (\mathbf{P}\{Y = 1, g(X) = 1 | X = x\} + \mathbf{P}\{Y = 0, g(X) = 0 | X = x\}) \\ &= 1 - (I_{\{g(x)=1\}} \mathbf{P}\{Y = 1 | X = x\} + I_{\{g(x)=0\}} \mathbf{P}\{Y = 0 | X = x\}) \\ &= 1 - (I_{\{g(x)=1\}} \eta(x) + I_{\{g(x)=0\}} (1 - \eta(x))), \end{aligned}$$

where  $I_A$  denotes the indicator of the set  $A$ . Thus, for every  $x \in \mathcal{R}^d$ ,

$$\begin{aligned} \mathbf{P}\{g(X) \neq Y | X = x\} - \mathbf{P}\{g^*(X) \neq Y | X = x\} &= \eta(x) (I_{\{g^*(x)=1\}} - I_{\{g(x)=1\}}) + (1 - \eta(x)) (I_{\{g^*(x)=0\}} - I_{\{g(x)=0\}}) \\ &= (2\eta(x) - 1) (I_{\{g^*(x)=1\}} - I_{\{g(x)=1\}}) \\ &\geq 0 \end{aligned}$$

by the definition of  $g^*$ . The statement now follows by integrating both sides with respect to  $\mu(dx)$ .  $\square$

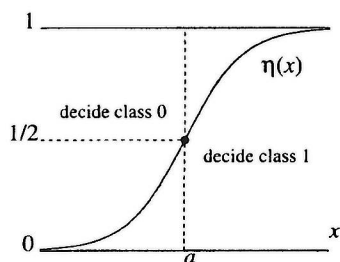


FIGURE 2.1. The Bayes decision in the example on the left is 1 if  $x > a$ , and 0 otherwise.

**REMARK.**  $g^*$  is called the Bayes decision function. The Bayes probability of error reveals that

$$L(g) = 1 - \mathbf{E} \{I_{\{\eta(X) > 1/2\}}\}$$

and in particular,

$$L^* = 1 - \mathbf{E} \{I_{\{\eta(X) > 1/2\}}\}$$

We observe that the a posteriori probability

$$\eta(x) = \mathbf{P}\{Y = 1 | X = x\}$$

minimizes the squared error loss. The Bayes decision function  $g^* : \mathcal{R}^d \rightarrow \mathcal{R}$ :

$$\mathbf{E} \{(\eta(X) - f(X))^2 | X = x\}$$

To see why the above inequality holds, note that

$$\begin{aligned} \mathbf{E} \{(\eta(X) - f(X))^2 | X = x\} &= \mathbf{E} \{(\eta(x) - f(x))^2 | X = x\} \\ &= (\eta(x) - f(x))^2 \\ &= (\eta(x) - f(x))^2 \\ &= (\eta(x) - f(x))^2 \end{aligned}$$

The conditional median, i.e.,  $\eta(x)$ , is even more closely related to the Bayes decision function.

## 2.2 A Simple Example

Let us consider the prediction problem when given a number of independent observations  $Y_1, \dots, Y_n$  of a random variable  $Y$ . If  $Y = 0$  stand for failure. The probability of failure per week. This, in itself, is not a problem because for that we would need to know the mind, health, and social habits of the person. If the probability of failure is probably monotonically increasing with age  $c > 0$ , say, our problem would be to predict the probability of failure given the age.

$$g^*(x) = \begin{cases} 1 & \text{if } x > a \\ 0 & \text{otherwise} \end{cases}$$