

# Classifiers Based on Bayes Decision Theory

## 2.1 INTRODUCTION

This is the first chapter, out of three, dealing with the design of the classifier in a pattern recognition system. The approach to be followed builds upon probabilistic arguments stemming from the statistical nature of the generated features. As has already been pointed out in the introductory chapter, this is due to the statistical variation of the patterns as well as to the noise in the measuring sensors. Adopting this reasoning as our kickoff point, we will design classifiers that classify an unknown pattern in the most probable of the classes. Thus, our task now becomes that of defining what "most probable" means.

Given a classification task of  $M$  classes,  $\omega_1, \omega_2, \dots, \omega_M$ , and an unknown pattern, which is represented by a feature vector  $\mathbf{x}$ , we form the  $M$  conditional probabilities  $P(\omega_i|\mathbf{x}), i = 1, 2, \dots, M$ . Sometimes, these are also referred to as *a posteriori probabilities*. In words, each of them represents the probability that the unknown pattern belongs to the respective class  $\omega_i$ , given that the corresponding feature vector takes the value  $\mathbf{x}$ . Who could then argue that these conditional probabilities are not sensible choices to quantify the term *most probable*? Indeed, the classifiers to be considered in this chapter compute either the maximum of these  $M$  values or, equivalently, the maximum of an appropriately defined function of them. The unknown pattern is then assigned to the class corresponding to this maximum.

The first task we are faced with is the computation of the conditional probabilities. The Bayes rule will once more prove its usefulness! A major effort in this chapter will be devoted to techniques for estimating probability density functions (pdf), based on the available experimental evidence, that is, the feature vectors corresponding to the patterns of the training set.

## 2.2 BAYES DECISION THEORY

We will initially focus on the two-class case. Let  $\omega_1, \omega_2$  be the two classes in which our patterns belong. In the sequel, we assume that the *a priori probabilities*

$P(\omega_1), P(\omega_2)$  are known. This is a very reasonable assumption, because even if they are not known, they can easily be estimated from the available training feature vectors. Indeed, if  $N$  is the total number of available training patterns, and  $N_1, N_2$  of them belong to  $\omega_1$  and  $\omega_2$ , respectively, then  $P(\omega_1) \approx N_1/N$  and  $P(\omega_2) \approx N_2/N$ .

The other statistical quantities assumed to be known are the class-conditional probability density functions  $p(\mathbf{x}|\omega_i), i = 1, 2$ , describing the distribution of the feature vectors in each of the classes. If these are not known, they can also be estimated from the available training data, as we will discuss later on in this chapter. The pdf  $p(\mathbf{x}|\omega_i)$  is sometimes referred to as the *likelihood function of  $\omega_i$  with respect to  $\mathbf{x}$* . Here we should stress the fact that an implicit assumption has been made. That is, the feature vectors can take any value in the  $l$ -dimensional feature space. In the case that feature vectors can take only discrete values, density functions  $p(\mathbf{x}|\omega_i)$  become probabilities and will be denoted by  $P(\mathbf{x}|\omega_i)$ .

We now have all the ingredients to compute our conditional probabilities, as stated in the introduction. To this end, let us recall from our probability course basics the Bayes rule (Appendix A)

$$P(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})} \quad (2.1)$$

where  $p(\mathbf{x})$  is the pdf of  $\mathbf{x}$  and for which we have (Appendix A)

$$p(\mathbf{x}) = \sum_{i=1}^2 p(\mathbf{x}|\omega_i)P(\omega_i) \quad (2.2)$$

The Bayes classification rule can now be stated as

$$\begin{aligned} \text{If } P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x}), \quad \mathbf{x} \text{ is classified to } \omega_1 \\ \text{If } P(\omega_1|\mathbf{x}) < P(\omega_2|\mathbf{x}), \quad \mathbf{x} \text{ is classified to } \omega_2 \end{aligned} \quad (2.3)$$

The case of equality is detrimental and the pattern can be assigned to either of the two classes. Using (2.1), the decision can equivalently be based on the inequalities

$$p(\mathbf{x}|\omega_1)P(\omega_1) \geq p(\mathbf{x}|\omega_2)P(\omega_2) \quad (2.4)$$

$p(\mathbf{x})$  is not taken into account, because it is the same for all classes and it does not affect the decision. Furthermore, if the *a priori* probabilities are equal, that is,  $P(\omega_1) = P(\omega_2) = 1/2$ , Eq. (2.4) becomes

$$p(\mathbf{x}|\omega_1) \geq p(\mathbf{x}|\omega_2) \quad (2.5)$$

Thus, the search for the maximum now rests on the values of the conditional pdfs evaluated at  $\mathbf{x}$ . Figure 2.1 presents an example of two equiprobable classes and shows the variations of  $p(\mathbf{x}|\omega_i), i = 1, 2$ , as functions of  $x$  for the simple case of a single feature ( $l = 1$ ). The dotted line at  $x_0$  is a threshold partitioning the feature space into two regions,  $R_1$  and  $R_2$ . According to the Bayes decision rule, for all values of  $x$  in  $R_1$  the classifier decides  $\omega_1$  and for all values in  $R_2$  it decides  $\omega_2$ . However, it is obvious from the figure that decision errors are unavoidable. Indeed, there is

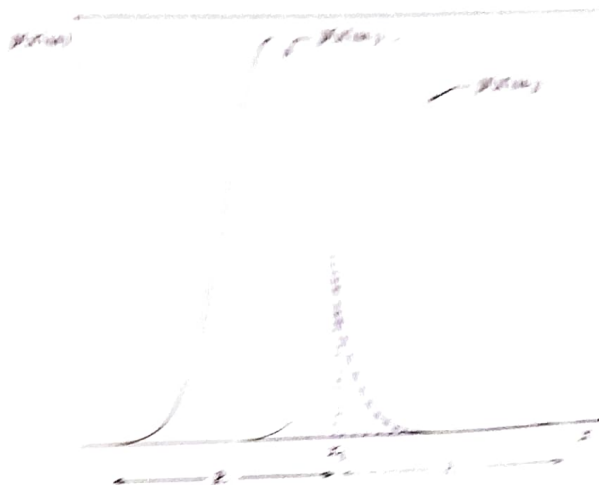


FIGURE 2.1

Example of the two regions  $R_1$  and  $R_2$  formed by the Bayesian classifier for the case of two equiprobable classes.

a finite probability for an  $x$  is to be in the  $R_2$  region and at the same time is being in class  $w_1$ . Then our decision is in error. The same is true for points originating from class  $w_2$ . It does not take much thought to see that the total probability  $P_e$  of committing a decision error for the case of two equiprobable classes is given by

$$P_e = \frac{1}{2} \int_{-\infty}^{x_0} p(x|w_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(x|w_1) dx \quad (2.6)$$

which is equal to the total shaded area under the curves in Figure 2.1. We have now touched on a very important issue. Our starting point is arrived at the Bayes classification rule was rather empirical via our interpretation of the term *most probable*. We will now see that this classification test, though simple in its formulation has a sounder mathematical interpretation.

### Minimizing the Classification Error Probability

We will show that the *Bayesian classifier is optimal with respect to minimizing the classification error probability*. Indeed, the reader can easily verify as an exercise that moving the threshold away from  $x_0$  in Figure 2.1, always increases the corresponding shaded area under the curves. Let us now proceed with a more formal proof.

**Proof.** Let  $R_1$  be the region of the feature space in which we decide in favor of  $w_1$  and  $R_2$  be the corresponding region for  $w_2$ . Then an error is made if  $x \in R_1$ , although it belongs to  $w_2$  or if  $x \in R_2$  although it belongs to  $w_1$ . That is

$$P_e = P(x \in R_2, w_1) + P(x \in R_1, w_2) \quad (2.7)$$

Total Error



$$P(A \cap B) = P(A, B)$$

$$P(A|B)P(B) =$$

$$P(B|A)P(A)$$

where  $P(\cdot, \cdot)$  is the joint probability of two events. Recalling, once more, our probability basics (Appendix A), this becomes

$$\begin{aligned} P_e &= P(\mathbf{x} \in R_2 | \omega_1)P(\omega_1) + P(\mathbf{x} \in R_1 | \omega_2)P(\omega_2) \\ &= P(\omega_1) \int_{R_2} p(\mathbf{x} | \omega_1) d\mathbf{x} + P(\omega_2) \int_{R_1} p(\mathbf{x} | \omega_2) d\mathbf{x} \end{aligned} \quad (2.8)$$

or using the Bayes rule

$$P_e = \int_{R_2} P(\omega_1 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} + \int_{R_1} P(\omega_2 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (2.9)$$

It is now easy to see that the error is minimized if the partitioning regions  $R_1$  and  $R_2$  of the feature space are chosen so that

$$\begin{aligned} R_1: P(\omega_1 | \mathbf{x}) &> P(\omega_2 | \mathbf{x}) \\ R_2: P(\omega_2 | \mathbf{x}) &> P(\omega_1 | \mathbf{x}) \end{aligned} \quad (2.10)$$

Indeed, since the union of the regions  $R_1, R_2$  covers all the space, from the definition of a probability density function we have that

$$\int_{R_1} P(\omega_1 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} + \int_{R_2} P(\omega_1 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} = P(\omega_1) \quad (2.11)$$

Combining Eqs. (2.9) and (2.11), we get

$$P_e = P(\omega_1) - \int_{R_1} (P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x})) p(\mathbf{x}) d\mathbf{x} \quad (2.12)$$

This suggests that the probability of error is minimized if  $R_1$  is the region of space in which  $P(\omega_1 | \mathbf{x}) > P(\omega_2 | \mathbf{x})$ . Then,  $R_2$  becomes the region where the reverse is true.  $\square$

So far, we have dealt with the simple case of two classes. Generalizations to the multiclass case are straightforward. In a classification task with  $M$  classes,  $\omega_1, \omega_2, \dots, \omega_M$ , an unknown pattern, represented by the feature vector  $\mathbf{x}$ , is assigned to class  $\omega_i$  if

$$P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x}) \quad \forall j \neq i \quad (2.13)$$

It turns out that such a choice also minimizes the classification error probability (Problem 2.1).

### Minimizing the Average Risk

The classification error probability is not always the best criterion to be adopted for minimization. This is because it assigns the same importance to all errors. However, there are cases in which some wrong decisions may have more serious implications than others. For example, it is much more serious for a doctor to make a wrong decision and a malignant tumor to be diagnosed as a benign one, than the other way round. If a benign tumor is diagnosed as a malignant one, the wrong decision will be cleared out during subsequent clinical examinations. However, the results

malignant tumor more harmful than benign.

from the wrong decision concerning a malignant tumor may be fatal. Thus, in such cases it is more appropriate to assign a penalty term to weigh each error. For our example, let us denote by  $\omega_1$  the class of malignant tumors and as  $\omega_2$  the class of the benign ones. Let, also,  $R_1, R_2$  be the regions in the feature space where we decide in favor of  $\omega_1$  and  $\omega_2$ , respectively. The error probability  $P_e$  is given by Eq. (2.8). Instead of selecting  $R_1$  and  $R_2$  so that  $P_e$  is minimized, we will now try to minimize a modified version of it, that is,

$$r = \lambda_{12} P(\omega_1) \int_{R_2} p(x|\omega_1) dx + \lambda_{21} P(\omega_2) \int_{R_1} p(x|\omega_2) dx \quad (2.14)$$

where each of the two terms that contributes to the overall error probability is weighted according to its significance. For our case, the reasonable choice would be to have  $\lambda_{12} > \lambda_{21}$ . Thus errors due to the assignment of patterns originating from class  $\omega_1$  to class  $\omega_2$  will have a larger effect on the cost function than the errors associated with the second term in the summation.

Let us now consider an  $M$ -class problem and let  $R_j, j = 1, 2, \dots, M$ , be the regions of the feature space assigned to classes  $\omega_i$ , respectively. Assume now that a feature vector  $x$  that belongs to class  $\omega_k$  lies in  $R_i, i \neq k$ . Then this vector is misclassified in  $\omega_i$  and an error is committed. A penalty term  $\lambda_{ki}$ , known as loss, is associated with this wrong decision. The matrix  $L$ , which has at its  $(k, i)$  location the corresponding penalty term, is known as the loss matrix.<sup>1</sup> Observe that in contrast to the philosophy behind Eq. (2.14), we have now allowed weights across the diagonal of the loss matrix ( $\lambda_{kk}$ ), which correspond to correct decisions. In practice, these are usually set equal to zero, although we have considered them here for the sake of generality. The risk or loss associated with  $\omega_k$  is defined as

$$r_k = \sum_{i=1}^M \lambda_{ki} \int_{R_i} p(x|\omega_k) dx \quad (2.15)$$

Observe that the integral is the overall probability of a feature vector from class  $\omega_k$  being classified in  $\omega_i$ . This probability is weighted by  $\lambda_{ki}$ . Our goal now is to choose the partitioning regions  $R_j$  so that the average risk

$$r = \sum_{k=1}^M r_k P(\omega_k)$$

$$= \sum_{i=1}^M \int_{R_i} \left( \sum_{k=1}^M \lambda_{ki} p(x|\omega_k) P(\omega_k) \right) dx \quad (2.16)$$

is minimized. This is achieved if each of the integrals is minimized, which is equivalent to selecting partitioning regions so that

$$x \in R_i \quad \text{if} \quad l_i \equiv \sum_{k=1}^M \lambda_{ki} p(x|\omega_k) P(\omega_k) < l_j \equiv \sum_{k=1}^M \lambda_{kj} p(x|\omega_k) P(\omega_k) \quad \forall j \neq i \quad (2.17)$$

<sup>1</sup> The terminology comes from the general decision theory.

$\omega_1 \rightarrow$  malignant  
 $\omega_2 \rightarrow$  benign

$\lambda_{12} \rightarrow \omega_1$  of  
Class 1 but  
assigning it to  
Class 2

$x \rightarrow \omega_k$  but  
assign to  $R_i$   
 $i \neq k$   
 $i \neq k$

It is obvious that if  $\lambda_{ki} = 1 - \delta_{ki}$ , where  $\delta_{ki}$  is *Kronecker's delta* (0 if  $k \neq i$  and 1 if  $k = i$ ), then minimizing the average risk becomes equivalent to minimizing the classification error probability.

*The two-class case.* For this specific case we obtain

$$\begin{aligned} I_1 &= \lambda_{11} p(\mathbf{x}|\omega_1)P(\omega_1) + \lambda_{21} p(\mathbf{x}|\omega_2)P(\omega_2) \\ I_2 &= \lambda_{12} p(\mathbf{x}|\omega_1)P(\omega_1) + \lambda_{22} p(\mathbf{x}|\omega_2)P(\omega_2) \end{aligned} \quad (2.18)$$

We assign  $\mathbf{x}$  to  $\omega_1$  if  $I_1 < I_2$ , that is,

$$(\lambda_{21} - \lambda_{22})p(\mathbf{x}|\omega_2)P(\omega_2) < (\lambda_{12} - \lambda_{11})p(\mathbf{x}|\omega_1)P(\omega_1) \quad (2.19)$$

It is natural to assume that  $\lambda_{ij} > \lambda_{ji}$  (correct decisions are penalized much less than wrong ones). Adopting this assumption, the decision rule (2.17) for the two-class case now becomes

$$\mathbf{x} \in \omega_1(\omega_2) \text{ if } I_{12} \equiv \frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} > (<) \frac{P(\omega_2)\lambda_{21} - \lambda_{22}}{P(\omega_1)\lambda_{12} - \lambda_{11}} \quad (2.20)$$

The ratio  $I_{12}$  is known as the *likelihood ratio* and the preceding test as the *likelihood ratio test*. Let us now investigate Eq. (2.20) a little further and consider the case of Figure 2.1. Assume that the loss matrix is of the form

$$L = \begin{bmatrix} 0 & \lambda_{12} \\ \lambda_{21} & 0 \end{bmatrix}$$

If misclassification of patterns that come from  $\omega_2$  is considered to have serious consequences, then we must choose  $\lambda_{21} > \lambda_{12}$ . Thus, patterns are assigned to class  $\omega_2$  if

$$p(\mathbf{x}|\omega_2) > p(\mathbf{x}|\omega_1) \frac{\lambda_{12}}{\lambda_{21}}$$

where  $P(\omega_1) = P(\omega_2) = 1/2$  has been assumed. That is,  $p(\mathbf{x}|\omega_1)$  is multiplied by a factor less than 1 and the effect of this is to move the threshold in Figure 2.1 to the left of  $x_0$ . In other words, region  $R_2$  is increased while  $R_1$  is decreased. The opposite would be true if  $\lambda_{21} < \lambda_{12}$ .

An alternative cost that sometimes is used for two class problems is the Neyman-Pearson criterion. The error for one of the classes is now constrained to be fixed and equal to a chosen value (Problem 2.6). Such a decision rule has been used, for example, in radar detection problems. The task there is to detect a target in the presence of noise. One type of error is the so-called *false alarm*—that is, to mistake the noise for a signal (target) present. Of course, the other type of error is to miss the signal and to decide in favor of the noise (*missed detection*). In many cases the error probability of false alarm is set equal to a predetermined threshold.



**Example 2.1**

In a two-class problem with a single feature  $x$  the pdfs are Gaussians with variance  $\sigma^2 = 1/2$  for both classes and mean values 0 and 1, respectively, that is,

$$p(x|\omega_1) = \frac{1}{\sqrt{\pi}} \exp(-x^2)$$

$$p(x|\omega_2) = \frac{1}{\sqrt{\pi}} \exp(-(x-1)^2)$$

If  $P(\omega_1) = P(\omega_2) = 1/2$ , compute the threshold value  $x_0$  (a) for minimum error probability and (b) for minimum risk if the loss matrix is

$$L = \begin{bmatrix} 0 & 0.5 \\ 1.0 & 0 \end{bmatrix}$$

Taking into account the shape of the Gaussian function graph (Appendix A), the threshold for the minimum probability case will be

$$x_0: \exp(-x^2) = \exp(-(x-1)^2)$$

Taking the logarithm of both sides, we end up with  $x_0 = 1/2$ . In the minimum risk case we get

$$x_0: \exp(-x^2) = 2 \exp(-(x-1)^2)$$

or  $x_0 = (1 - \ln 2)/2 < 1/2$ ; that is, the threshold moves to the left of  $1/2$ . If the two classes are not equiprobable, then it is easily verified that if  $P(\omega_1) > (<) P(\omega_2)$  the threshold moves to the right (left). That is, we expand the region in which we decide in favor of the most probable class, since it is better to make fewer errors for the most probable class.

## 2.3 DISCRIMINANT FUNCTIONS AND DECISION SURFACES

It is by now clear that minimizing either the risk or the error probability or the Neyman-Pearson criterion is equivalent to partitioning the feature space into  $M$  regions, for a task with  $M$  classes. If regions  $R_i, R_j$  happen to be contiguous, then they are separated by a decision surface in the multidimensional feature space. For the minimum error probability case, this is described by the equation

$$P(\omega_i|x) - P(\omega_j|x) = 0 \quad (2.21)$$

From the one side of the surface this difference is positive, and from the other it is negative. Sometimes, instead of working directly with probabilities (or risk functions), it may be more convenient, from a mathematical point of view, to work with an equivalent function of them, for example,  $g_i(x) \equiv f(P(\omega_i|x))$ , where  $f(\cdot)$  is a monotonically increasing function.  $g_i(x)$  is known as a discriminant function. The decision test (2.13) is now stated as

$$\text{classify } x \text{ in } \omega_i \text{ if } g_i(x) > g_j(x) \quad \forall j \neq i \quad (2.22)$$

The decision surfaces, separating contiguous regions, are described by

$$g_i(x) \equiv g_i(x) - g_j(x) = 0, \quad i, j = 1, 2, \dots, M, \quad i \neq j \quad (2.23)$$

For minimum risk classifier:  $R(\alpha_i/x)$

$$g_i(x) = -R(\alpha_i/x)$$

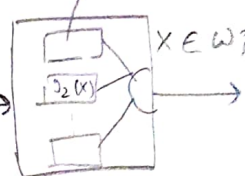
Min error rate classification

$$g_i(x) = P(\omega_i|x)$$

$$P(\omega_i|x)$$

(functional units)

$$g_i(x)$$



Classifier (black box)

$$g_1(x), g_2(x), \dots, g_n(x)$$

Discriminant function.

$$g_i(x) \equiv f(P(\omega_i|x))$$

$f(\cdot) \rightarrow$  monotonically inc. fun

$$g_i(x) > g_j(x) \quad \forall j \neq i$$

$$x \in \omega_i$$