

## 4.4 ESTIMATION OF PROBABILITY DENSITY FUNCTIONS

The most important task in implementing a statistical approach for solving pattern classification problems is estimating the density functions  $p(\mathbf{x}|C_i)$ ,  $1 \leq i \leq m$ . We will first show how to use the *maximum entropy principle* to obtain the *form* of probability density functions.

### 4.4.1 Form of the Density Function

The principle of maximum entropy states that in the case where a probability density function of a random variable is not known, the function which maximizes the entropy of this variable subject to known specified constraints is an appropriate choice. Any other choice would show a bias to some information obtained from the given data. The maximum entropy solution is easily derived when the constraints are given in the form of averages associated with the probability density function. Given a probability density function  $p(\mathbf{x})$ , the associated entropy is

$$E = - \int_{\mathbf{x}} p(\mathbf{x}) \ln[p(\mathbf{x})] d\mathbf{x} \quad (4.4.1)$$

and we assume the constraints

$$\int_{\mathbf{x}} f_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \alpha_i, \quad 0 \leq i \leq M \quad (4.4.2)$$

where  $f_0(\mathbf{x}) = 1$  and  $\alpha_0 = 1$ . We wish to obtain  $p(\mathbf{x})$  which satisfies Eq. (4.4.2) while minimizing the entropy  $E$  of Eq. (4.4.1). This is done using Lagrange multipliers. Define

$$E_1 = E + \sum_{i=0}^M \lambda_i \left[ \int_{\mathbf{x}} f_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} - \alpha_i \right] \quad (4.4.3)$$

where the constants  $\lambda_i$ ,  $0 \leq i \leq M$  are yet to be determined. By virtue of Eq. (4.4.1) we get

$$E_1 = - \int_x p(x) [\ln[p(x)] - \sum_{i=0}^M \lambda_i f_i(x)] dx - \sum_{i=0}^M \lambda_i \alpha_i \quad (4.4.4)$$

The partial derivative of  $E_1$  with respect to  $p(x)$  is

$$\frac{\partial E_1}{\partial [p(x)]} = - \int_x [\ln[p(x)] - \sum_{i=0}^M \lambda_i f_i(x) + 1] dx \quad (4.4.5)$$

and to obtain the maximum entropy solution the integrand must vanish, i.e.

$$p(x) = \exp \left[ \sum_{i=0}^M \lambda_i f_i(x) - 1 \right] \quad (4.4.6)$$

We still have freedom of choosing  $\lambda_i$ ,  $a \leq i \leq m$  and these coefficients are chosen so that Eq. (4.4.2) holds. Once the form of the probability density function is known, we may turn and perform the next step: estimating the parameters of this density.

■ **Example 4.4.1** Consider a random variable  $x$  which is characterized by

$$a < x < b, \quad \int_0^{\infty} p(x) dx = 1$$

By virtue of Eq. (4.4.6) we obtain

$$p(x) = \exp(\lambda_0 - 1), \quad \int_a^b \exp(\lambda_0 - 1) dx = 1$$

and therefore

---

$$p(x) = \begin{cases} \frac{1}{b-a} & , \quad a < x < b \\ 0 & , \quad \text{otherwise} \end{cases}$$



- **Example 4.4.2** Assume that the *a priori* information about  $x$  is

$$x \geq 0 \quad , \quad \int_0^{\infty} p(x) dx = 1 \quad , \quad \int_0^{\infty} xp(x) dx = \mu$$

Using Eq. (4.4.6) we obtain

$$p(x) = \exp(\lambda_0 - 1 + \lambda_1 x)$$

and in order to satisfy the two constraints the density function must be

$$p(x) = \begin{cases} (1/\mu) \exp(-x/\mu) & , \quad x \geq 0 \\ 0 & , \quad \text{otherwise} \end{cases}$$



#### 4.4.2 Estimating the Mean Vector and Covariance Matrix

Consider a pattern population with probability density function  $p(\mathbf{x})$ . The mean vector of this population is given by

$$\boldsymbol{\mu} = E(\mathbf{x}) = \int_{\mathbf{x}} \mathbf{x} p(\mathbf{x}) d\mathbf{x} \quad (4.4.7)$$

If the patterns are in  $R^n$ , then  $\boldsymbol{\mu}$  is a vector with  $n$  components  $(\mu_1, \dots, \mu_n)$ . Let  $\{\mathbf{x}_i\}_{i=1}^N$  denote the given patterns. An approximate to  $\boldsymbol{\mu}$  is simply

$$\boldsymbol{\mu} \approx \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \quad (4.4.8)$$

The covariance matrix  $\mathbf{C} = (c_{jk})$ ,  $1 \leq j, k \leq n$  satisfies

$$c_{jk} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_j - \mu_j)(x_k - \mu_k) p(x_j, x_k) dx_j dx_k \quad (4.4.9)$$

We can also rewrite  $\mathbf{C}$  as

$$\begin{aligned} \mathbf{C} &= E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \\ &= E[\mathbf{x}\mathbf{x}^T - 2\mathbf{x}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T] \\ &= E[\mathbf{x}\mathbf{x}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T \end{aligned} \quad (4.4.10)$$

and use the new expression to approximate  $\mathbf{C}$  as

$$\mathbf{C} \approx \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T - \boldsymbol{\mu}\boldsymbol{\mu}^T \quad (4.4.11)$$

Both estimates for  $\boldsymbol{\mu}$  and for  $\mathbf{C}$  can be conveniently used in a recursive manner. Let  $N$  be the current number of sample patterns and assume on additional incoming pattern. Denote by  $\boldsymbol{\mu}(N)$ ,  $\mathbf{C}(N)$  the current mean vector and covariance matrix. Then

$$\begin{aligned} \boldsymbol{\mu}(N+1) &= \frac{1}{N+1} \sum_{i=1}^{N+1} \mathbf{x}_i \\ &= \frac{1}{N+1} \left( \sum_{i=1}^N \mathbf{x}_i + \mathbf{x}_{N+1} \right) \\ &= \frac{1}{N+1} (N\boldsymbol{\mu}(N) + \mathbf{x}_{N+1}) \end{aligned} \quad (4.4.12)$$

where  $\boldsymbol{\mu}(1) = \mathbf{x}_1$ . This recursive expression updates the mean vector.

In the case of the covariance matrix we obtain

$$\begin{aligned}
\mathbf{C}(N+1) &= \frac{1}{N+1} \sum_{i=1}^{N+1} \mathbf{x}_i \mathbf{x}_i^T - \boldsymbol{\mu}(N+1) \boldsymbol{\mu}^T(N+1) \\
&= \frac{1}{N+1} \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + \mathbf{x}_{N+1} \mathbf{x}_{N+1}^T \right) - \boldsymbol{\mu}(N+1) \boldsymbol{\mu}^T(N+1) \\
&= \frac{1}{N+1} (N\mathbf{C}(N) + N\boldsymbol{\mu}(N) \boldsymbol{\mu}^T(N) + \mathbf{x}_{N+1} \mathbf{x}_{N+1}^T) \\
&\quad - \frac{1}{(N+1)^2} (N\boldsymbol{\mu}(N) + \mathbf{x}_{N+1})(N\boldsymbol{\mu}^T(N) + \mathbf{x}_{N+1}^T) \quad (4.4.13)
\end{aligned}$$

To start the calculation of  $\mathbf{C}(N)$  we use the relation

$$\mathbf{C}(1) = \mathbf{x}_1 \mathbf{x}_1^T - \boldsymbol{\mu}(1) \boldsymbol{\mu}^T(1)$$

to obtain  $\mathbf{C}(1)=0$ .

■ **Example 4.4.3** Consider the sample patterns

$$\mathbf{x}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{x}_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \mathbf{x}_4 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

To start the recursive procedure we set

$$\boldsymbol{\mu}(1) = \mathbf{x}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{C}_1 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

and then, using Eqs. (4.4.12-13), obtain

$$\begin{aligned}
\boldsymbol{\mu}(2) &= \begin{pmatrix} 1/2 \\ 0 \end{pmatrix}, \quad \boldsymbol{\mu}(3) = \begin{pmatrix} 2/3 \\ 1/3 \end{pmatrix}, \quad \boldsymbol{\mu}(4) = \begin{pmatrix} 3/4 \\ 3/4 \end{pmatrix}, \\
\mathbf{C}(2) &= \begin{pmatrix} 1/4 & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{C}(3) = \begin{pmatrix} 2/9 & 1/9 \\ 1/9 & 2/9 \end{pmatrix}, \quad \mathbf{C}(4) = \begin{pmatrix} 3/16 & 3/16 \\ 3/16 & 11/16 \end{pmatrix}
\end{aligned}$$



### 4.4.3 Estimation by Functional Approximation

If the form of the density function is not known we may estimate it directly using functional approximation.

Let  $p(\mathbf{x})$  denote the probability density function  $p(\mathbf{x}|C)$  and consider an approximate  $\tilde{p}(\mathbf{x})$  given by

$$\tilde{p}(\mathbf{x}) = \sum_{i=1}^m a_i \phi_i(\mathbf{x}) \quad (4.4.14)$$

where  $\{\phi_i(\mathbf{x})\}_{i=1}^m$  are specified basis functions. We wish to minimize

$$E = \int_{\mathbf{x}} w(\mathbf{x}) [p(\mathbf{x}) - \tilde{p}(\mathbf{x})]^2 d\mathbf{x} \quad (4.4.15)$$

or

$$E = \int_{\mathbf{x}} w(\mathbf{x}) \left[ p(\mathbf{x}) - \sum_{i=1}^m a_i \phi_i(\mathbf{x}) \right]^2 d\mathbf{x} \quad (4.4.16)$$

where  $w(\mathbf{x})$  is a specified weight function. Solving the system

$$\frac{\partial E}{\partial a_i} = 0, \quad i = 1, \dots, m \quad (4.4.17)$$

provides a set of linear equations

$$\sum_{j=1}^m a_j \int_{\mathbf{x}} w(\mathbf{x}) \phi_j(\mathbf{x}) \phi_i(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{x}} w(\mathbf{x}) \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}, \quad 1 \leq i \leq m \quad (4.4.18)$$

The right-hand sides of these equations are simply the expected values of  $w(\mathbf{x})\phi_i(\mathbf{x})$ ,  $1 \leq i \leq m$ . If  $\{\mathbf{x}_k\}_{k=1}^N$  are given samples which belong to  $C$ , these expected values can be estimated as

$$\int_{\mathbf{x}} w(\mathbf{x}) \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N} \sum_{k=1}^N w(\mathbf{x}_k) \phi_i(\mathbf{x}_k) \quad (4.4.19)$$


---

We thus replace the system given by Eq. (4.4.18) by

$$\sum_{j=1}^m a_j \int_{\mathbf{x}} w(\mathbf{x}) \phi_j(\mathbf{x}) \phi_i(\mathbf{x}) d\mathbf{x} = \frac{1}{N} \sum_{k=1}^N w(\mathbf{x}_k) \phi_i(\mathbf{x}_k), \quad 1 \leq i \leq m \quad (4.4.20)$$

In the particular case where  $\phi_i(\mathbf{x})$ ,  $1 \leq i \leq k$  are orthogonal with respect to  $w(\mathbf{x})$ , we have

$$\int_{\mathbf{x}} w(\mathbf{x}) \phi_j(\mathbf{x}) \phi_i(\mathbf{x}) d\mathbf{x} = \begin{cases} A_i, & j = i \\ 0, & j \neq i \end{cases} \quad (4.4.21)$$

and consequently

$$a_i = \frac{1}{NA_i} \sum_{k=1}^N w(\mathbf{x}_k) \phi_i(\mathbf{x}_k) \quad (4.4.22)$$

The expression given by Eq. (4.4.22) provides an easy way to obtain  $a_i(N+1)$  from  $a_i(N)$ . Indeed

$$\begin{aligned} a_i(N+1) &= \frac{1}{(N+1)A_i} \sum_{k=1}^{N+1} w(\mathbf{x}_k) \phi_i(\mathbf{x}_k) \\ &= \frac{1}{(N+1)A_i} [NA_i a_i(N) + w(\mathbf{x}_{N+1}) \phi_i(\mathbf{x}_{N+1})] \end{aligned} \quad (4.4.23)$$

In decision making problems, since the terms  $w(\mathbf{x}_k)$  in Eq. (4.4.22) are independent of  $i$  and are therefore common to all the coefficients, they can usually be eliminated from the process, without violating the discriminatory characteristics of the coefficients. We may usually therefore, for such problems, apply a simpler relation

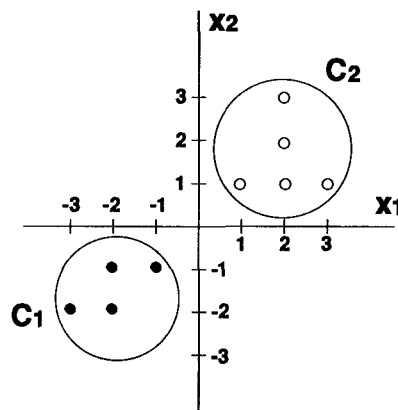
$$a_i = \frac{1}{N} \sum_{k=1}^N \phi_i(\mathbf{x}_k) \quad (4.4.24)$$


---

i.e. relate to the orthogonal basis functions, as if they were orthonormal and simplifying the computations.

Since  $p(x)$  is not known, one may not be able to decide how large  $m$  should be. Usually we start with a prefixed  $m$  and experiment with the training set to determine whether  $\tilde{p}(x)$  is an acceptable approximate to  $p(x)$ . If the classification performance of  $\tilde{p}(x)$  is poor, we increase  $m$  until we reach a 'saturation' state, i.e. until adding new terms has no effect on the classification quality of  $\tilde{p}(x)$ . It can be shown that in general  $\tilde{p}(x)$  approaches  $p(x)$  as  $m \rightarrow \infty$  and  $N \rightarrow \infty$ .

■ **Example 4.4.4** Consider the two-class classification problem given in Fig. 4.4.1



■ **Figure 4.4.1** Bayes classification using functional approximation.

Assuming that the entire domain of the patterns is the whole plane, one is tempted to use the Hermite polynomials which are orthogonal over the interval  $(-\infty, \infty)$ . The first two polynomials are  $H_0(x) = 1$ ,  $H_1(x) = 2x$ . If four basis functions are considered we may choose



$$\phi_1(\mathbf{x}) = H_0(x_1)H_0(x_2) = 1$$

$$\phi_2(\mathbf{x}) = H_1(x_1)H_0(x_2) = 2x_1$$

$$\phi_3(\mathbf{x}) = H_0(x_1)H_1(x_2) = 2x_2$$

$$\phi_4(\mathbf{x}) = H_1(x_1)H_1(x_2) = 4x_1x_2$$

We may treat the functions as if they were orthonormal and obtain

$$a_l^{(i)} = \frac{1}{N_l} \sum_{k=1}^{N_l} \phi_l(\mathbf{x}_k^{(i)}), \quad 1 \leq l \leq 2$$

where  $a_l^{(i)}$  are the coefficients associated with class  $l$ ,  $N_l$  – the number of patterns in class  $l$ , and  $\mathbf{x}_k^{(i)}$  are the patterns in class  $l$ . Thus

$$a_1^{(1)} = \frac{1}{4}(1+1+1+1) = 1, \quad a_2^{(1)} = \frac{1}{4}(-2-4-4-6) = -4$$

$$a_3^{(1)} = \frac{1}{4}(-2-2-4-4) = -3, \quad a_4^{(1)} = \frac{1}{4}(4+8+16+24) = 13$$

$$a_1^{(2)} = \frac{1}{5}(1+1+1+1+1) = 1, \quad a_2^{(2)} = \frac{1}{5}(2+4+4+4+6) = 4$$

$$a_3^{(2)} = \frac{1}{5}(2+2+2+4+6) = 3.2, \quad a_4^{(2)} = \frac{1}{5}(4+8+12+16+24) = 12.8$$

and consequently

$$\tilde{p}(\mathbf{x} | C_1) = 1 - 8x_1 - 6x_2 + 52x_1x_2$$

$$\tilde{p}(\mathbf{x} | C_2) = 1 + 8x_1 + 6.4x_2 + 51.2x_1x_2$$

The decision functions are

$$d_1(\mathbf{x}) = \tilde{p}(\mathbf{x} | C_1)p(C_1)$$

$$d_2(\mathbf{x}) = \tilde{p}(\mathbf{x} | C_2)p(C_2)$$

and by assuming  $p(C_1) = p(C_2) = 1/2$  we get

$$d_1(\mathbf{x}) = \frac{1}{2} - 4x_1 - 3x_2 + 26x_1x_2$$

$$d_2(\mathbf{x}) = \frac{1}{2} + 4x_1 + 3.2x_2 + 25.6x_1x_2$$

The decision boundary is therefore

$$d_{12}(\mathbf{x}) = d_1(\mathbf{x}) - d_2(\mathbf{x}) = -8x_1 - 6.2x_2 + 0.4x_1x_2 = 0$$



## PROBLEMS

1. Use the maximum entropy principle to obtain the probability density function if the information

$$-\infty < x < \infty, \int_{-\infty}^{\infty} p(x)dx = 1, \int_{-\infty}^{\infty} xp(x)dx = \mu, \int_{-\infty}^{\infty} x^2 p(x)dx = \sigma^2$$

is *a priori* known.

2. Given the sample patterns

$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 2 \\ 0 \\ 0 \end{pmatrix}, \quad \mathbf{x}_3 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, \quad \mathbf{x}_4 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

Use Eqs. (4.4.12-13) to estimate  $\mu(i)$ ,  $C(i)$ ,  $1 \leq i \leq 4$ .

3. Apply the method of functional approximation to get estimates of  $p(x|C_1)$  and  $p(x|C_2)$  where

$$C_1 = \{(1,0)^T, (1,1)^T, (2,1)^T, (3,0)^T, (4,1)^T\}$$

$$C_2 = \{(-1,0)^T, (-2,0)^T, (-2,-1)^T, (-3,1)^T, (-3,2)^T\}$$

Use the first three 2-D Hermite polynomials and Eq. (4.4.24) to obtain the coefficients.

4. Repeat problem 3 but use the first four Hermite polynomials.
5. Repeat problems 3 and 4 using Hermite orthonormal functions. Replace  $H_k(x)$  by

$$H_k^*(x) = \frac{\exp(-x^2/2)}{\sqrt{2^k k! \sqrt{\pi}}} H_k(x)$$

but still use Eq. (4.4.24).

---