

adaptive filtering owes much to the classic paper of Widrow and Hoff (1960) for pioneering the so-called *least-mean-square (LMS) algorithm*, also known as the *delta rule*. The LMS algorithm is simple to implement yet highly effective in application. Indeed, it is the workhorse of *linear* adaptive filtering, linear in the sense that the neuron operates in its linear mode. Adaptive filters have been successfully applied in such diverse fields as antennas, communication systems, control systems, radar, sonar, seismology, and biomedical engineering (Widrow and Stearns, 1985; Haykin, 1996).

The LMS algorithm and the perceptron are naturally related. It is therefore proper for us to study them together in one chapter.

### Organization of the Chapter

The chapter is organized in two parts. The first part, consisting of Sections 3.2 through 3.7 deals with linear adaptive filters and the LMS algorithm. The second part, consisting of Sections 3.8 through 3.10 deals with Rosenblatt's perceptron. From a presentation point of view, we find it more convenient to discuss linear adaptive filters first and then Rosenblatt's perceptron, reversing the historical order in which they appeared.

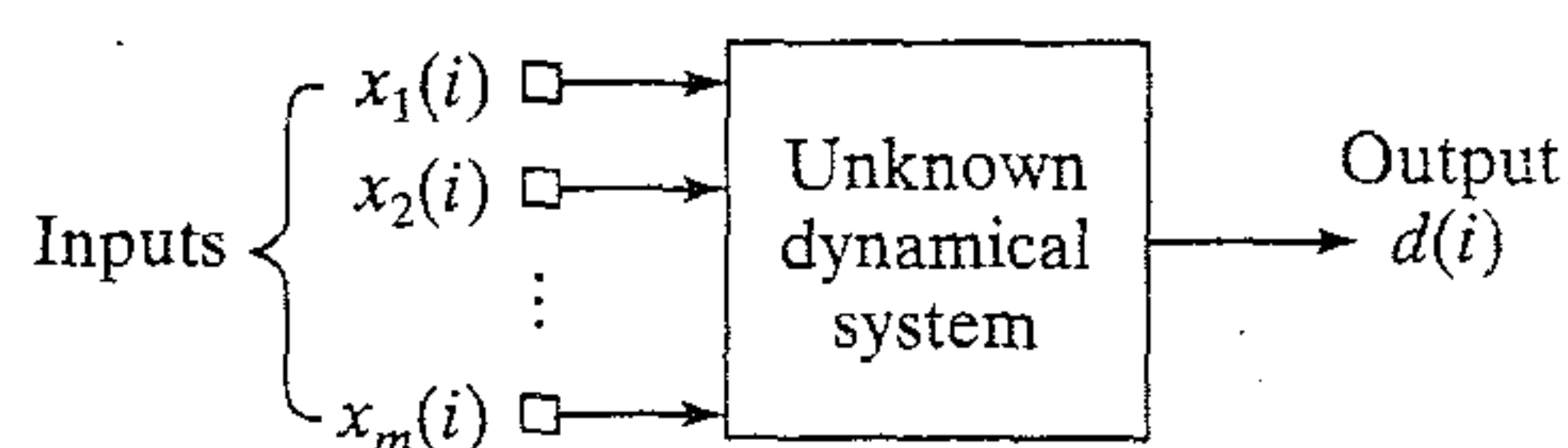
In Section 3.2 we address the adaptive filtering problem, followed by Section 3.3, a review of three unconstrained optimization techniques: the method of steepest descent, Newton's method, and Gauss-Newton method, which are particularly relevant to the study of adaptive filters. In Section 3.4 we discuss a linear least-squares filter, which asymptotically approaches the Wiener filter as the data length increases. The Wiener filter provides an ideal framework for the performance of linear adaptive filters operating in a stationary environment. In Section 3.5 we describe the LMS algorithm, including a discussion of its virtues and limitations. In Section 3.6 we explore the idea of learning curves commonly used to assess the performance of adaptive filters. This is followed by a discussion of annealing schedules for the LMS algorithm in Section 3.7.

Then moving on to Rosenblatt's perceptron, Section 3.8 presents some basic considerations involved in its operation. In Section 3.9 we describe the algorithm for adjusting the synaptic weight vector of the perceptron for pattern classification of linearly separable classes, and demonstrate convergence of the algorithm. In Section 3.10 we consider the relationship between the perceptron and the Bayes classifier for a Gaussian environment.

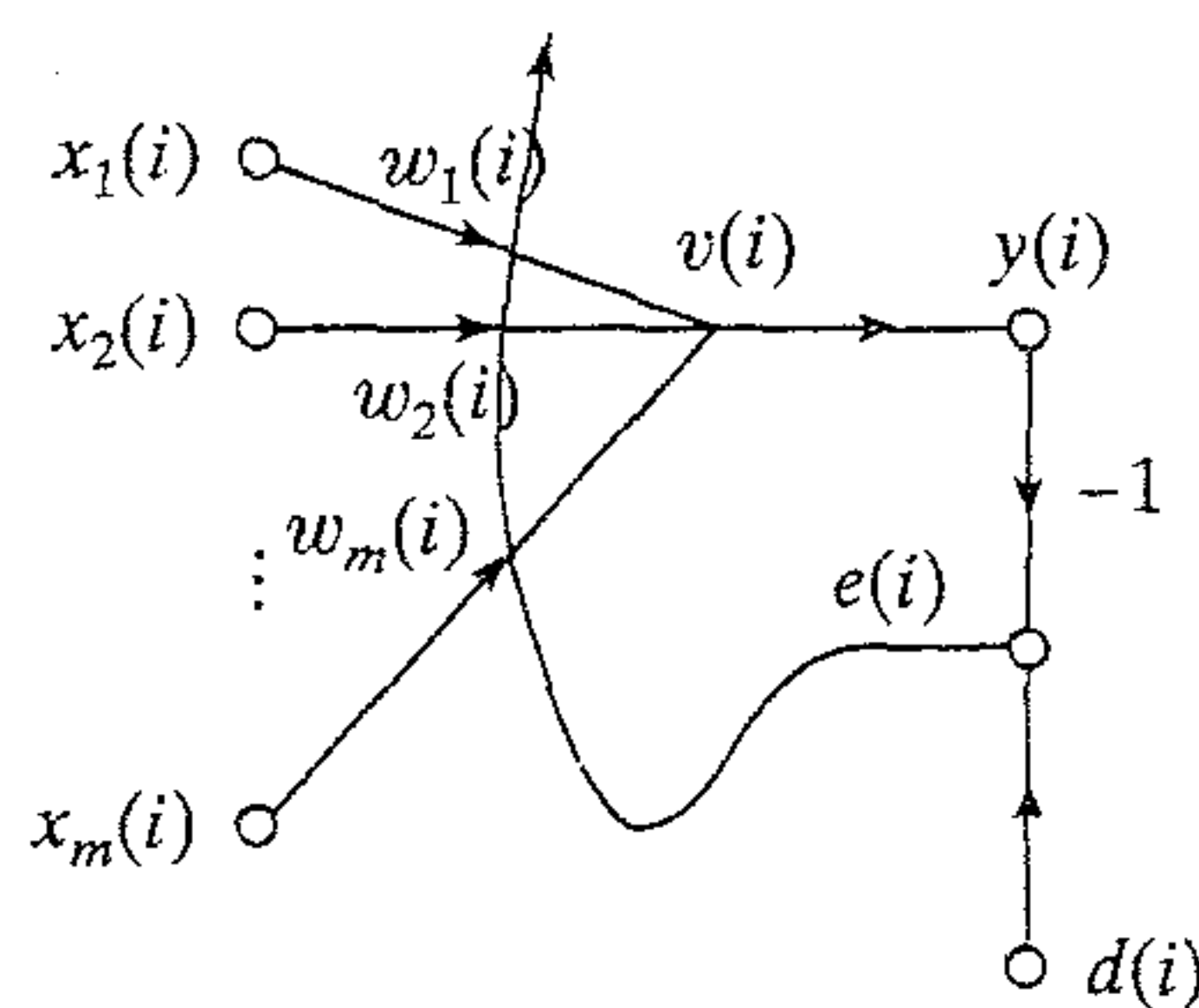
The chapter concludes with summary and discussion in Section 3.11.

## 3.2 ADAPTIVE FILTERING PROBLEM

Consider a *dynamical system*, the mathematical characterization of which is *unknown*. All that we have available on the system is a set of labeled input-output data generated by the system at discrete instants of time at some uniform rate. Specifically, when an  $m$ -dimensional stimulus  $\mathbf{x}(i)$  is applied across  $m$  input nodes of



(a)



(b)

**FIGURE 3.1** (a) Unknown dynamical system. (b) Signal-flow graph of adaptive model for the system.

the system, the system responds by producing a scalar output  $d(i)$ , where  $i = 1, 2, \dots, n, \dots$  as depicted in Fig. 3.1a. Thus, the external behavior of the system is described by the data set

$$\mathcal{T}: \{\mathbf{x}(i), d(i); i = 1, 2, \dots, n, \dots\} \quad (3.1)$$

where

$$\mathbf{x}(i) = [x_1(i), x_2(i), \dots, x_m(i)]^T$$

The samples comprising  $\mathcal{T}$  are identically distributed according to an unknown probability law. The dimension  $m$  pertaining to the input vector  $\mathbf{x}(i)$  is referred to as the *dimensionality of the input space* or simply as *dimensionality*.

The stimulus  $\mathbf{x}(i)$  can arise in one of two fundamentally different ways, one spatial and the other temporal:

- The  $m$  elements of  $\mathbf{x}(i)$  originate at different points in space; in this case we speak of  $\mathbf{x}(i)$  as a *snapshot* of data.
- The  $m$  elements of  $\mathbf{x}(i)$  represent the set of present and  $(m - 1)$  past values of some excitation that are *uniformly spaced in time*.

The problem we address is how to design a multiple input-single output *model* of the unknown dynamical system by building it around a single linear neuron. The neuronal model operates under the influence of an algorithm that *controls* necessary adjustments to the synaptic weights of the neuron, with the following points in mind:

- The algorithm starts from an *arbitrary setting* of the neuron's synaptic weights.
- Adjustments to the synaptic weights, in response to statistical variations in the system's behavior, are made on a *continuous* basis (i.e., time is incorporated into the constitution of the algorithm).

- Computations of adjustments to the synaptic weights are completed inside a time interval that is one sampling period long.

The neuronal model described is referred to as an *adaptive filter*. Although the description is presented in the context of a task clearly recognized as one of *system identification*, the characterization of the adaptive filter is general enough to have wide application.

Figure 3.1b shows a signal-flow graph of the adaptive filter. Its operation consists of two continuous processes:

1. *Filtering process*, which involves the computation of two signals:
  - An output, denoted by  $y(i)$ , that is produced in response to the  $m$  elements of the stimulus vector  $\mathbf{x}(i)$ , namely,  $x_1(i), x_2(i), \dots, x_m(i)$ .
  - An error signal, denoted by  $e(i)$ , that is obtained by comparing the output  $y(i)$  to the corresponding output  $d(i)$  produced by the unknown system. In effect,  $d(i)$  acts as a *desired response* or *target signal*.
2. *Adaptive process*, which involves the automatic adjustment of the synaptic weights of the neuron in accordance with the error signal  $e(i)$ .

Thus, the combination of these two processes working together constitutes a *feedback loop* acting around the neuron.

Since the neuron is linear, the output  $y(i)$  is exactly the same as the induced local field  $v(i)$ ; that is,

$$y(i) = v(i) = \sum_{k=1}^m w_k(i)x_k(i) \quad (3.2)$$

where  $w_1(i), w_2(i), \dots, w_m(i)$  are the  $m$  synaptic weights of the neuron, measured at time  $i$ . In matrix form we may express  $y(i)$  as an inner product of the vectors  $\mathbf{x}(i)$  and  $\mathbf{w}(i)$  as follows:

$$y(i) = \mathbf{x}^T(i)\mathbf{w}(i) \quad (3.3)$$

where

$$\mathbf{w}(i) = [w_1(i), w_2(i), \dots, w_m(i)]^T$$

Note that the notation for a synaptic weight has been simplified here by *not* including an additional subscript to identify the neuron, since we only have a single neuron to deal with. This practice is followed throughout the chapter. The neuron's output  $y(i)$  is compared to the corresponding output  $d(i)$  received from the unknown system at time  $i$ . Typically,  $y(i)$  is different from  $d(i)$ ; hence, their comparison results in the error signal:

$$e(i) = d(i) - y(i) \quad (3.4)$$

The manner in which the error signal  $e(i)$  is used to control the adjustments to the neuron's synaptic weights is determined by the cost function used to derive the adaptive filtering algorithm of interest. This issue is closely related to that of optimization. It is therefore appropriate to present a review of unconstrained optimization methods. The material is applicable not only to linear adaptive filters but also to neural networks in general.



$$\begin{aligned}
\mathbf{r}_{xd} &= E[\mathbf{x}(i)d(i)] \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}(i)d(i) \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^T(n)\mathbf{d}(n)
\end{aligned} \tag{3.31}$$

where  $E$  denotes the statistical expectation operator. Accordingly, we may reformulate the linear least-squares solution of Eq. (3.27) as follows:

$$\begin{aligned}
\mathbf{w}_o &= \lim_{n \rightarrow \infty} \mathbf{w}(n+1) \\
&= \lim_{n \rightarrow \infty} (\mathbf{X}^T(n)\mathbf{X}(n))^{-1} \mathbf{X}^T(n)\mathbf{d}(n) \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} (\mathbf{X}^T(n)\mathbf{X}(n))^{-1} \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^T(n)\mathbf{d}(n) \\
&= \mathbf{R}_x^{-1} \mathbf{r}_{xd}
\end{aligned} \tag{3.32}$$

where  $\mathbf{R}_x^{-1}$  is the inverse of the correlation matrix  $\mathbf{R}_x$ . The weight vector  $\mathbf{w}_o$  is called the *Wiener solution* to the linear optimum filtering problem in recognition of the contributions of Norbert Wiener to this problem (Widrow and Stearns, 1985; Haykin, 1996). Accordingly, we may make the following statement:

*For an ergodic process, the linear least-squares filter asymptotically approaches the Wiener filter as the number of observations approaches infinity.*

Designing the Wiener filter requires knowledge of the second-order statistics: the correlation matrix  $\mathbf{R}_x$  of the input vector  $\mathbf{x}(n)$  and the cross-correlation vector  $\mathbf{r}_{xd}$  between  $\mathbf{x}(n)$  and the desired response  $d(n)$ . However, this information is not available in many important situations encountered in practice. We may deal with an unknown environment by using a *linear adaptive filter*, adaptive in the sense that the filter is able to adjust its free parameters in response to statistical variations in the environment. A highly popular algorithm for doing this kind of adjustment on a continuing basis is the least-mean-square algorithm, which is intimately related to the Wiener filter.

### 3.5 LEAST-MEAN-SQUARE ALGORITHM

The *least-mean-square (LMS) algorithm* is based on the use of *instantaneous values* for the cost function, namely,

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2} e^2(n) \tag{3.33}$$

where  $e(n)$  is the error signal measured at time  $n$ . Differentiating  $\mathcal{E}(\mathbf{w})$  with respect to the weight vector  $\mathbf{w}$  yields

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = e(n) \frac{\partial e(n)}{\partial \mathbf{w}} \tag{3.34}$$

As with the linear least-squares filter, the LMS algorithm operates with a linear neuron so we may express the error signal as

$$e(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n) \quad (3.35)$$

Hence,

$$\frac{\partial e(n)}{\partial \mathbf{w}(n)} = -\mathbf{x}(n)$$

and

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}(n)} = -\mathbf{x}(n)e(n)$$

Using this latter result as an *estimate* for the gradient vector, we may write

$$\hat{\mathbf{g}}(n) = -\mathbf{x}(n)e(n) \quad (3.36)$$

Finally, using Eq. (3.36) for the gradient vector in Eq. (3.12) for the method of steepest descent, we may formulate the LMS algorithm as follows:

$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \eta \mathbf{x}(n)e(n) \quad (3.37)$$

where  $\eta$  is the learning-rate parameter. The feedback loop around the weight vector  $\hat{\mathbf{w}}(n)$  in the LMS algorithm behaves like a *low-pass filter*, passing the low frequency components of the error signal and attenuating its high frequency components (Haykin, 1996). The average time constant of this filtering action is inversely proportional to the learning-rate parameter  $\eta$ . Hence, by assigning a small value to  $\eta$ , the adaptive process will progress slowly. More of the past data are then remembered by the LMS algorithm, resulting in a more accurate filtering action. In other words, the inverse of the learning-rate parameter  $\eta$  is a measure of the *memory* of the LMS algorithm.

In Eq. (3.37) we have used  $\hat{\mathbf{w}}(n)$  in place of  $\mathbf{w}(n)$  to emphasize the fact that the LMS algorithm produces an *estimate* of the weight vector that would result from the use of the method of steepest descent. As a consequence, in using the LMS algorithm we sacrifice a distinctive feature of the steepest descent algorithm. In the steepest descent algorithm the weight vector  $\mathbf{w}(n)$  follows a well-defined trajectory in weight space for a prescribed  $\eta$ . In contrast, in the LMS algorithm the weight vector  $\hat{\mathbf{w}}(n)$  traces a random trajectory. For this reason, the LMS algorithm is sometimes referred to as a “stochastic gradient algorithm.” As the number of iterations in the LMS algorithm approaches infinity,  $\hat{\mathbf{w}}(n)$  performs a random walk (Brownian motion) about the Wiener solution  $\mathbf{w}_o$ . The important point is the fact that, unlike the method of steepest descent, the LMS algorithm does *not* require knowledge of the statistics of the environment.

A summary of the LMS algorithm is presented in Table 3.1, which clearly illustrates the simplicity of the algorithm. As indicated in this table, for the *initialization* of the algorithm, it is customary to set the initial value of the weight vector in the algorithm equal to zero.

**TABLE 3.1** Summary of the LMS Algorithm

*Training Sample:* Input signal vector =  $\mathbf{x}(n)$   
Desired response =  $d(n)$

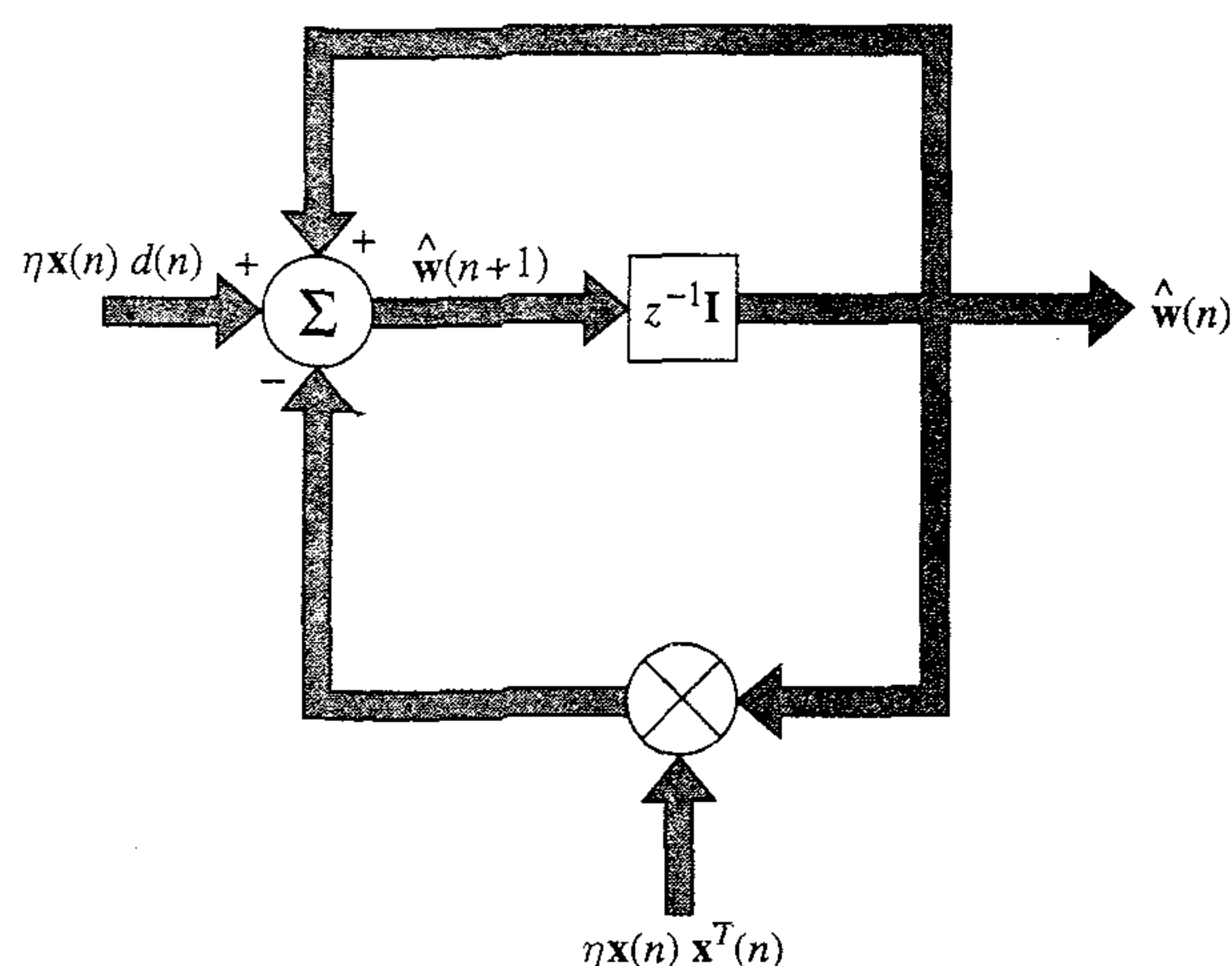
*User-selected parameter:*  $\eta$

*Initialization.* Set  $\hat{\mathbf{w}}(0) = \mathbf{0}$ .

*Computation.* For  $n = 1, 2, \dots$ , compute

$$e(n) = d(n) - \hat{\mathbf{w}}^T(n)\mathbf{x}(n)$$

$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \eta \mathbf{x}(n)e(n)$$



**FIGURE 3.3** Signal-flow graph representation of the LMS algorithm.

### Signal-Flow Graph Representation of the LMS Algorithm

By combining Eqs. (3.35) and (3.37) we may express the evolution of the weight vector in the LMS algorithm as follows:

$$\begin{aligned}\hat{\mathbf{w}}(n+1) &= \hat{\mathbf{w}}(n) + \eta \mathbf{x}(n)[d(n) - \mathbf{x}^T(n)\hat{\mathbf{w}}(n)] \\ &= [\mathbf{I} - \eta \mathbf{x}(n)\mathbf{x}^T(n)]\hat{\mathbf{w}}(n) + \eta \mathbf{x}(n)d(n)\end{aligned}\quad (3.38)$$

where  $\mathbf{I}$  is the identity matrix. In using the LMS algorithm, we recognize that

$$\hat{\mathbf{w}}(n) = z^{-1}[\hat{\mathbf{w}}(n+1)] \quad (3.39)$$

where  $z^{-1}$  is the *unit-delay operator*, implying storage. Using Eqs. (3.38) and (3.39), we may thus represent the LMS algorithm by the signal-flow graph depicted in Fig. 3.3. This signal-flow graph reveals that the LMS algorithm is an example of a *stochastic feedback system*. The presence of feedback has a profound impact on the convergence behavior of the LMS algorithm.

### Convergence Considerations of the LMS Algorithm

From control theory we know that the stability of a feedback system is determined by the parameters that constitute its feedback loop. From Fig. 3.3 we see that it is the lower feedback loop that adds variability to the behavior of the LMS algorithm. In par-



ticular, there are two distinct quantities, the learning-rate parameter  $\eta$  and the input vector  $\mathbf{x}(n)$ , that determine the transmittance of this feedback loop. We therefore deduce that the convergence behavior (i.e., stability) of the LMS algorithm is influenced by the statistical characteristics of the input vector  $\mathbf{x}(n)$  and the value assigned to the learning-rate parameter  $\eta$ . Casting this observation in a different way, we may state that for a specified environment that supplies the input vector  $\mathbf{x}(n)$ , we have to exercise care in the selection of the learning-rate parameter  $\eta$  for the LMS algorithm to be convergent.

The first criterion for convergence of the LMS algorithm is *convergence of the mean*, described by

$$E[\hat{\mathbf{w}}(n)] \rightarrow \mathbf{w}_o \quad \text{as } n \rightarrow \infty \quad (3.40)$$

where  $\mathbf{w}_o$  is the Wiener solution. Unfortunately, such a convergence criterion is of little practical value, since a sequence of zero-mean, but otherwise arbitrary, random vectors converges in this sense.

From a practical point of view, the convergence issue that really matters is *convergence in the mean square*, described by

$$E[e^2(n)] \rightarrow \text{constant as } n \rightarrow \infty \quad (3.41)$$

Unfortunately, a detailed convergence analysis of the LMS algorithm in the mean square is rather complicated. To make the analysis mathematically tractable, the following assumptions are usually made:

1. The successive input vectors  $\mathbf{x}(1), \mathbf{x}(2), \dots$  are statistically independent of each other.
2. At time step  $n$ , the input vector  $\mathbf{x}(n)$  is statistically independent of all previous samples of the desired response, namely  $d(1), d(2), \dots, d(n-1)$ .
3. At time step  $n$ , the desired response  $d(n)$  is dependent on  $\mathbf{x}(n)$ , but statistically independent of all previous values of the desired response.
4. The input vector  $\mathbf{x}(n)$  and desired response  $d(n)$  are drawn from Gaussian-distributed populations.

A statistical analysis of the LMS algorithm so based is called the *independence theory* (Widrow et al., 1976).

By invoking the elements of independence theory and assuming that the learning-rate parameter  $\eta$  is sufficiently small, it is shown in Haykin (1996) that the LMS is convergent in the mean square provided that  $\eta$  satisfies the condition

$$0 < \eta < \frac{2}{\lambda_{\max}} \quad (3.42)$$

where  $\lambda_{\max}$  is the *largest eigenvalue* of the correlation matrix  $\mathbf{R}_x$ . In typical applications of the LMS algorithm, however, knowledge of  $\lambda_{\max}$  is not available. To overcome this difficulty, the *trace* of  $\mathbf{R}_x$  may be taken as a conservative estimate for  $\lambda_{\max}$ , in which case the condition of Eq. (3.42) may be reformulated as

$$0 < \eta < \frac{2}{\text{tr}[\mathbf{R}_x]} \quad (3.43)$$

where  $\text{tr}[\mathbf{R}_x]$  denotes the trace of matrix  $\mathbf{R}_x$ . By definition, the trace of a square matrix is equal to the sum of its diagonal elements. Since each diagonal element of the correlation matrix  $\mathbf{R}_x$  equals the mean-square value of the corresponding sensor input, we may restate the condition for convergence of the LMS algorithm in the mean square as follows:

$$0 < \eta < \frac{2}{\text{sum of mean-square values of the sensor inputs}} \quad (3.44)$$

Provided the learning-rate parameter satisfies this condition, the LMS algorithm is also assured of convergence of the mean. That is, convergence in the mean square implies convergence of the mean, but the converse is not necessarily true.

### Virtues and Limitations of the LMS Algorithm

An important virtue of the LMS algorithm is its simplicity, as exemplified by the summary of the algorithm presented in Table 3.1. Moreover, the LMS algorithm is model independent and therefore *robust*, which means that small model uncertainty and small disturbances (i.e., disturbances with small energy) can only result in small estimation errors (error signals). In precise mathematical terms, the LMS algorithm is optimal in accordance with the  $H^\infty$  (or *minimax*) criterion (Hassibi et al., 1993, 1996). The basic philosophy of optimality in the  $H^\infty$  sense is to cater to the worst-case scenario<sup>4</sup>:

*If you do not know what you are up against, plan for the worst and optimize.*

For a long time the LMS algorithm was regarded as an instantaneous approximation to the gradient-descent algorithm. However, the  $H^\infty$  optimality of LMS provides this widely used algorithm with a rigorous footing. In particular, it explains its ability to work satisfactorily in a stationary as well as in a nonstationary environment. By a “nonstationary” environment we mean one where the statistics vary with time. In such an environment, the optimum Wiener solution takes on a time-varying form, and the LMS algorithm now has the additional task of *tracking* variations in the parameters of the Wiener filter.

The primary limitations of the LMS algorithm are its slow rate of convergence and sensitivity to variations in the eigenstructure of the input (Haykin, 1996). The LMS algorithm typically requires a number of iterations equal to about 10 times the dimensionality of the input space for it to reach a steady-state condition. The slow rate of convergence becomes particularly serious when the dimensionality of the input space becomes high. As for sensitivity to changes in environmental conditions, the LMS algorithm is particularly sensitive to variations in the *condition number* or *eigenvalue spread* of the correlation matrix  $\mathbf{R}_x$  of the input vector  $\mathbf{x}$ . The condition number of  $\mathbf{R}_x$ , denoted by  $\chi(\mathbf{R}_x)$ , is defined by

$$\chi(\mathbf{R}_x) = \frac{\lambda_{\max}}{\lambda_{\min}} \quad (3.45)$$

where  $\lambda_{\max}$  and  $\lambda_{\min}$  are the maximum and minimum eigenvalues of the matrix  $\mathbf{R}_x$ , respectively. The sensitivity of the LMS algorithm to variations in the condition number  $\chi(\mathbf{R}_x)$  becomes particularly acute when the training sample to which the input vector  $\mathbf{x}(n)$  belongs is *ill conditioned*, that is, when the condition number  $\chi(\mathbf{R}_x)$  is high.<sup>5</sup>



Note that in the LMS algorithm the *Hessian matrix*, defined as the second derivative of the cost function  $\mathcal{E}(\mathbf{w})$  with respect to  $\mathbf{w}$ , is equal to the correlation matrix  $\mathbf{R}_x$ ; see Problem 3.8. Thus, in the discussion presented here, we could have just as well spoken in terms of the Hessian as the correlation matrix  $\mathbf{R}_x$ .

### 3.6 LEARNING CURVES

An informative way of examining the convergence behavior of the LMS algorithm, or an adaptive filter in general, is to plot the *learning curve* of the filter under varying environmental conditions. The learning curve is a *plot of the mean-square value of the estimation error,  $\mathcal{E}_{av}(n)$ , versus the number of iterations,  $n$ .*

Imagine an experiment involving an *ensemble* of adaptive filters, with each filter operating under the control of a specific algorithm. It is assumed that the details of the algorithm, including initialization, are the same for all the filters. The differences between the filters arise from the *random* manner in which the input vector  $\mathbf{x}(n)$  and the desired response  $d(n)$  are drawn from the available training sample. For each filter we plot the squared value of the estimation error (i.e., the difference between the desired response and the actual filter output) versus the number of iterations. A *sample* learning curve so obtained consists of *noisy* exponentials, the noise being due to the inherently stochastic nature of the adaptive filter. To compute the *ensemble-averaged learning curve* (i.e., plot of  $\mathcal{E}_{av}(n)$  versus  $n$ ), we take the average of these sample learning curves over the ensemble of adaptive filters used in the experiment, thereby smoothing out the effects of noise.

Assuming that the adaptive filter is stable, we find that the ensemble-averaged learning curve starts from a large value  $\mathcal{E}_{av}(0)$  determined by the initial conditions, then decreases at some rate depending on the type of filter used, and finally converges to a steady-state value  $\mathcal{E}_{av}(\infty)$ , as illustrated in Fig. 3.4. On the basis of this learning curve we may define the *rate of convergence* of the adaptive filter as the number of iterations,  $n$ , required for  $\mathcal{E}_{av}(n)$  to decrease to some arbitrarily chosen value, such as 10 percent of the initial value  $\mathcal{E}_{av}(0)$ .

Another useful characteristic of an adaptive filter that is deduced from the ensemble-averaged learning curve is the *misadjustment*, denoted by  $\mathcal{M}$ . Let  $\mathcal{E}_{\min}$  denote the minimum mean-square error produced by the Wiener filter, designed on the basis of known values of the correlation matrix  $\mathbf{R}_x$  and cross-correlation vector  $\mathbf{r}_{xd}$ . We may define the *misadjustment* for the adaptive filter as follows (Widrow and Stearns, 1985; Haykin, 1996):

$$\begin{aligned}\mathcal{M} &= \frac{\mathcal{E}(\infty) - \mathcal{E}_{\min}}{\mathcal{E}_{\min}} \\ &= \frac{\mathcal{E}(\infty)}{\mathcal{E}_{\min}} - 1\end{aligned}\tag{3.46}$$

The misadjustment  $\mathcal{M}$  is a dimensionless quantity, providing a measure of how close the adaptive filter is to optimality in the mean-square error sense. The smaller  $\mathcal{M}$  is compared to unity, the more *accurate* is the adaptive filtering action of the algorithm. It is customary to express the misadjustment  $\mathcal{M}$  as a percentage. Thus, for example, a misadjustment of 10 percent means that the adaptive filter produces a mean-square error

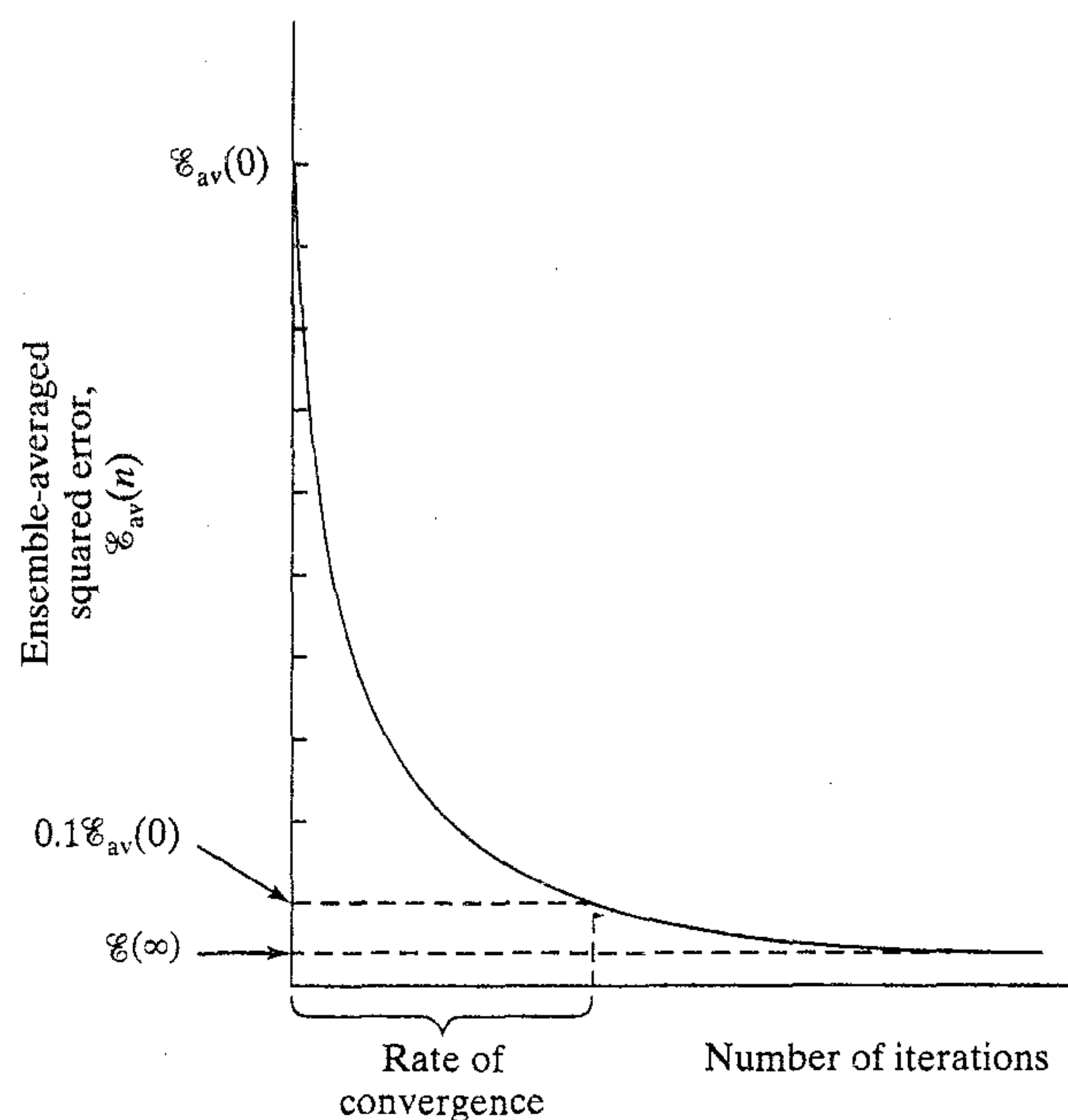


FIGURE 3.4 Idealized learning curve of the LMS algorithm.

(after adaptation is completed) that is 10 percent greater than the minimum mean-square error  $\mathcal{E}_{\min}$  produced by the corresponding Wiener filter. Such performance is ordinarily considered to be satisfactory in practice.

Another important characteristic of the LMS algorithm is the *settling time*. However, there is no unique definition for the settling time. We may, for example, approximate the learning curve by a single exponential with *average time constant*  $\tau_{av}$ , and so use  $\tau_{av}$  as a rough measure of the settling time. The smaller the value of  $\tau_{av}$  is, the faster the settling time will be (i.e., the faster the LMS algorithm will converge to a “steady-state” condition).

To a good degree of approximation, the misadjustment  $\mathcal{M}$  of the LMS algorithm is directly proportional to the learning-rate parameter  $\eta$ , whereas the average time constant  $\tau_{av}$  is inversely proportional to the learning-rate parameter  $\eta$  (Widrow and Stearns, 1985; Haykin, 1996). We therefore have conflicting results in the sense that if the learning-rate parameter is reduced so as to reduce the misadjustment, then the settling time of the LMS algorithm is increased. Conversely, if the learning-rate parameter is increased so as to accelerate the learning process, then the misadjustment is increased. Careful attention must be given to the choice of the learning parameter  $\eta$  in the design of the LMS algorithm in order to produce a satisfactory overall performance.

### 3.7 LEARNING-RATE ANNEALING SCHEDULES

The difficulties encountered with the LMS algorithm may be attributed to the fact that the learning-rate parameter is maintained constant throughout the computation, as shown by

$$\eta(n) = \eta_0 \quad \text{for all } n \quad (3.47)$$

This is the simplest possible form the learning-rate parameter can assume. In contrast, in *stochastic approximation*, which goes back to the classic paper by Robbins and Monro (1951), the learning-rate parameter is time-varying. The particular time-varying form most commonly used in stochastic approximation literature is described by

$$\eta(n) = \frac{c}{n} \quad (3.48)$$

where  $c$  is a constant. Such a choice is indeed sufficient to guarantee convergence of the stochastic approximation algorithm (Ljung, 1977; Kushner and Clark, 1978). However, when the constant  $c$  is large, there is a danger of parameter blowup for small  $n$ .

As an alternative to Eqs. (3.47) and (3.48), we may use the *search-then-converge schedule*, defined by Darken and Moody (1992)

$$\eta(n) = \frac{\eta_0}{1 + (n/\tau)} \quad (3.49)$$

where  $\eta_0$  and  $\tau$  are user-selected constants. In the early stages of adaptation involving a number of iterations  $n$ , small compared to the *search time constant*  $\tau$ , the learning-rate parameter  $\eta(n)$  is approximately equal to  $\eta_0$ , and the algorithm operates essentially as the “standard” LMS algorithm, as indicated in Fig. 3.5. Hence, by choosing a high value for  $\eta_0$  within the permissible range, we hope that the adjustable weights of the filter will find and hover about a “good” set of values. Then, for a number of iterations  $n$  large compared to the search time constant  $\tau$ , the learning-rate parameter  $\eta(n)$  approximates as  $c/n$ , where  $c = \tau\eta_0$ , as illustrated in Fig. 3.5. The algorithm now operates as a traditional stochastic approximation algorithm, and the weights converge to their optimum values. Thus the search-then-converge schedule has the potential to combine the desirable features of the standard LMS with traditional stochastic approximation theory.

### 3.8 PERCEPTRON

We now come to the second part of the chapter that deals with Rosenblatt’s perceptron, henceforth referred to simply as the *perceptron*. Whereas the LMS algorithm described in the preceding sections is built around a linear neuron, the perceptron is built around a nonlinear neuron, namely, the *McCulloch–Pitts model* of a neuron. From Chapter 1 we recall that such a neuronal model consists of a linear combiner followed by a hard limiter (performing the signum function), as depicted in Fig. 3.6. The summing node of the neuronal model computes a linear combination of the inputs applied to its synapses, and also incorporates an externally applied bias. The resulting sum, that is, the induced local field, is applied to a hard limiter. Accordingly, the neuron produces an output equal to  $+1$  if the hard limiter input is positive, and  $-1$  if it is negative.

In the signal-flow graph model of Fig. 3.6, the synaptic weights of the perceptron are denoted by  $w_1, w_2, \dots, w_m$ . Correspondingly, the inputs applied to the perceptron



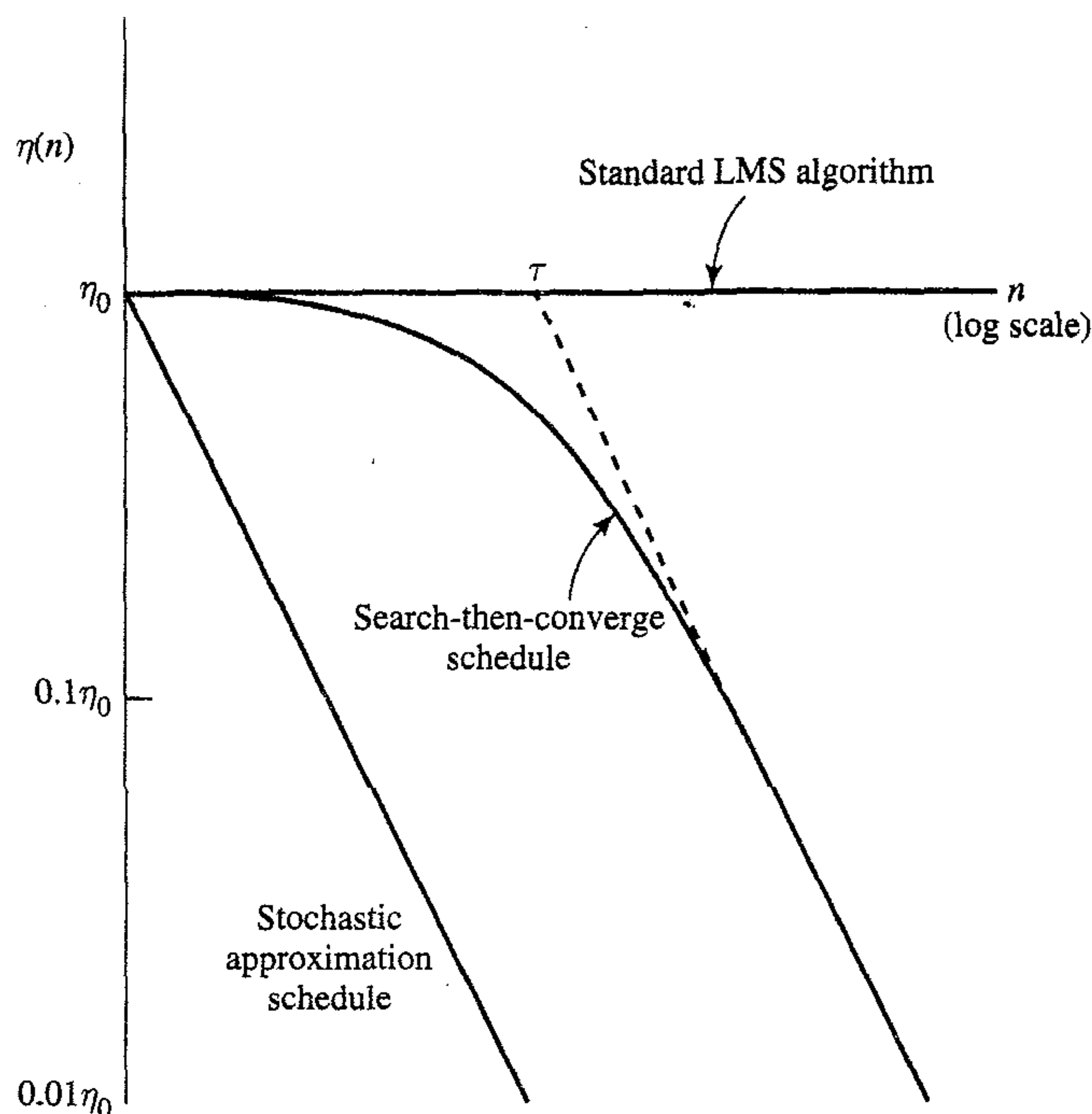


FIGURE 3.5 Learning-rate annealing schedules.

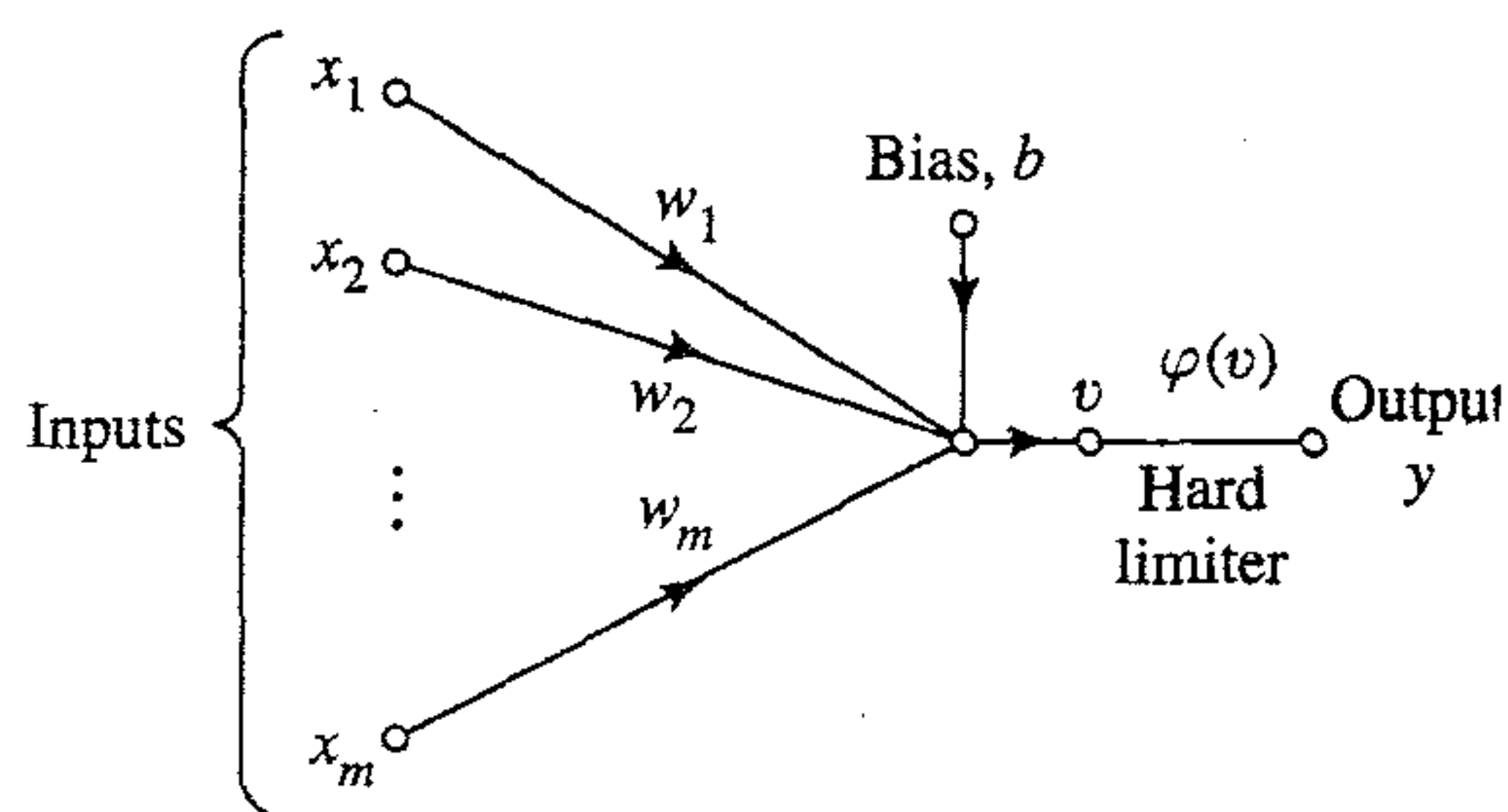


FIGURE 3.6 Signal-flow graph of the perceptron.

are denoted by  $x_1, x_2, \dots, x_m$ . The externally applied bias is denoted by  $b$ . From the model we find that the hard limiter input or induced local field of the neuron is

$$v = \sum_{i=1}^m w_i x_i + b \quad (3.50)$$

The goal of the perceptron is to correctly classify the set of externally applied stimuli  $x_1, x_2, \dots, x_m$  into one of two classes,  $\mathcal{C}_1$  or  $\mathcal{C}_2$ . The decision rule for the classification is to assign the point represented by the inputs  $x_1, x_2, \dots, x_m$  to class  $\mathcal{C}_1$  if the perceptron output  $y$  is  $+1$  and to class  $\mathcal{C}_2$  if it is  $-1$ .

To develop insight into the behavior of a pattern classifier, it is customary to plot a map of the decision regions in the  $m$ -dimensional signal space spanned by the  $m$

input variables  $x_1, x_2, x_m$ . In the simplest form of the perceptron there are two decision regions separated by a *hyperplane* defined by

$$\sum_{i=1}^m w_i x_i + b = 0 \quad (3.51)$$

This is illustrated in Fig. 3.7 for the case of two input variables  $x_1$  and  $x_2$ , for which the decision boundary takes the form of a straight line. A point  $(x_1, x_2)$  that lies above the boundary line is assigned to class  $\mathcal{C}_1$  and a point  $(x_1, x_2)$  that lies below the boundary line is assigned to class  $\mathcal{C}_2$ . Note also that the effect of the bias  $b$  is merely to shift the decision boundary away from the origin.

The synaptic weights  $w_1, w_2, \dots, w_m$  of the perceptron can be adapted on an iteration-by-iteration basis. For the adaptation we may use an error-correction rule known as the perceptron convergence algorithm.

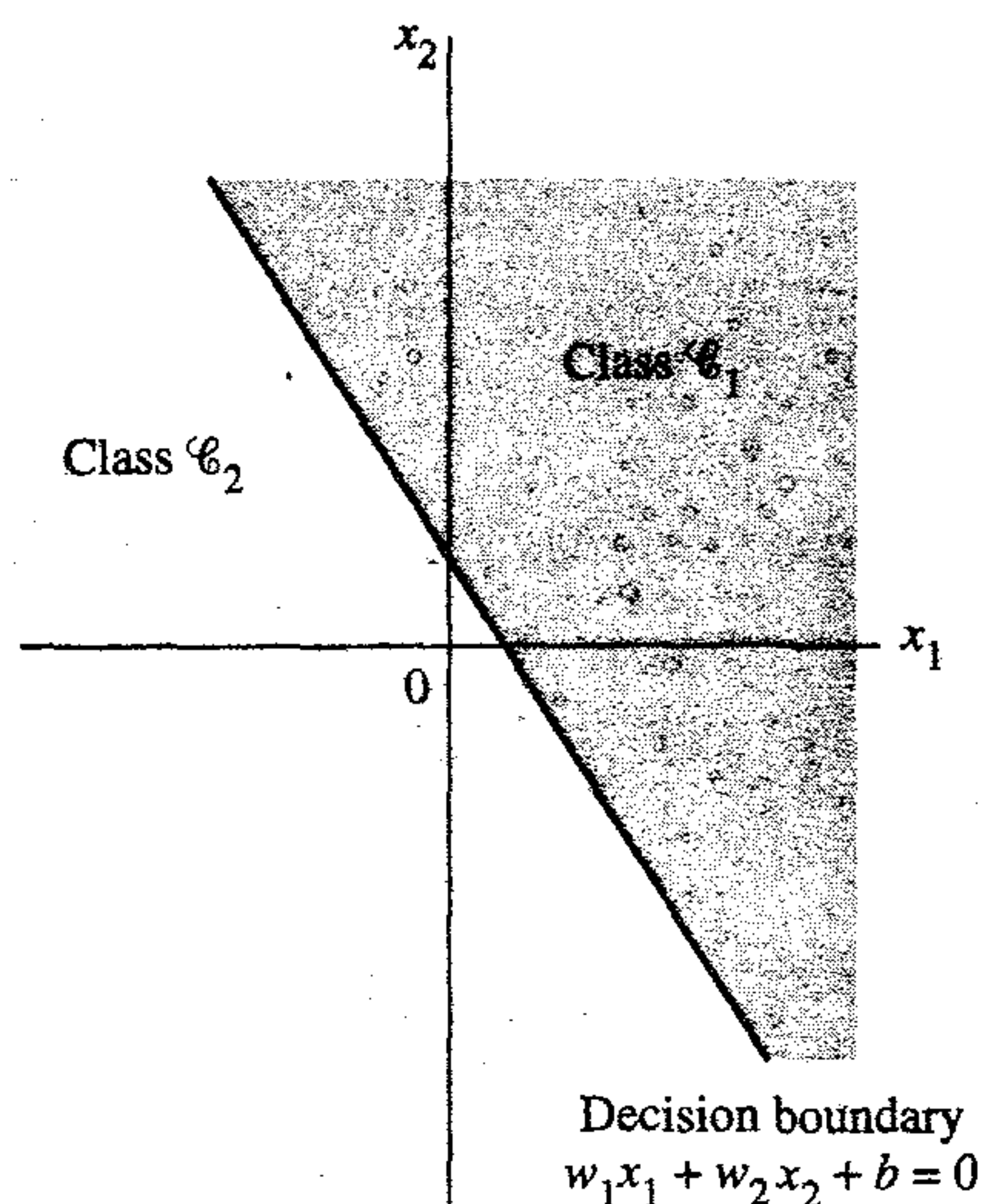
### 3.9 PERCEPTRON CONVERGENCE THEOREM

To derive the error-correction learning algorithm for the perceptron, we find it more convenient to work with the modified signal-flow graph model in Fig. 3.8. In this second model, which is equivalent to that of Fig. 3.6, the bias  $b(n)$  is treated as a synaptic weight driven by a fixed input equal to  $+1$ . We may thus define the  $(m+1)$ -by-1 input vector

$$\mathbf{x}(n) = [+1, x_1(n), x_2(n), \dots, x_m(n)]^T$$

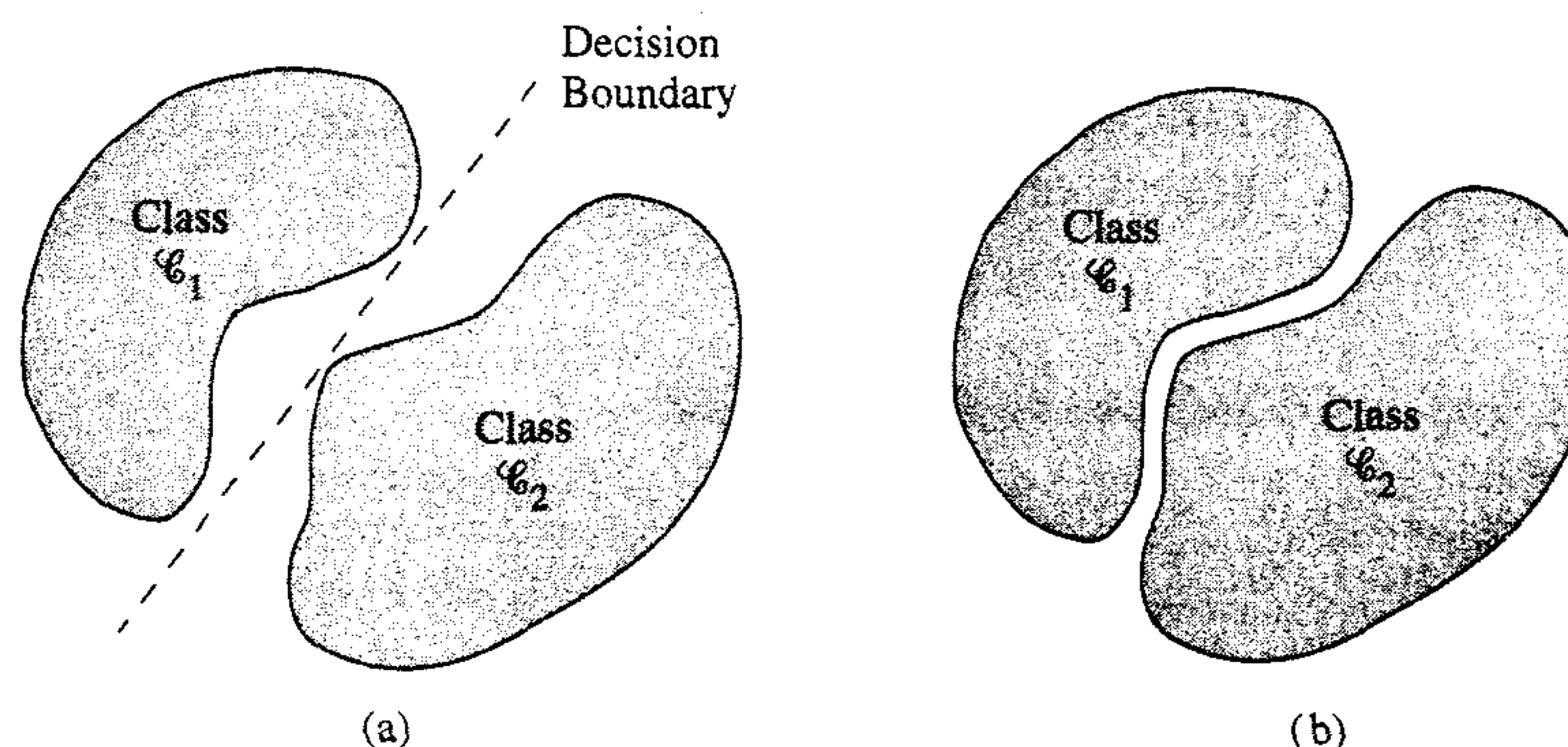
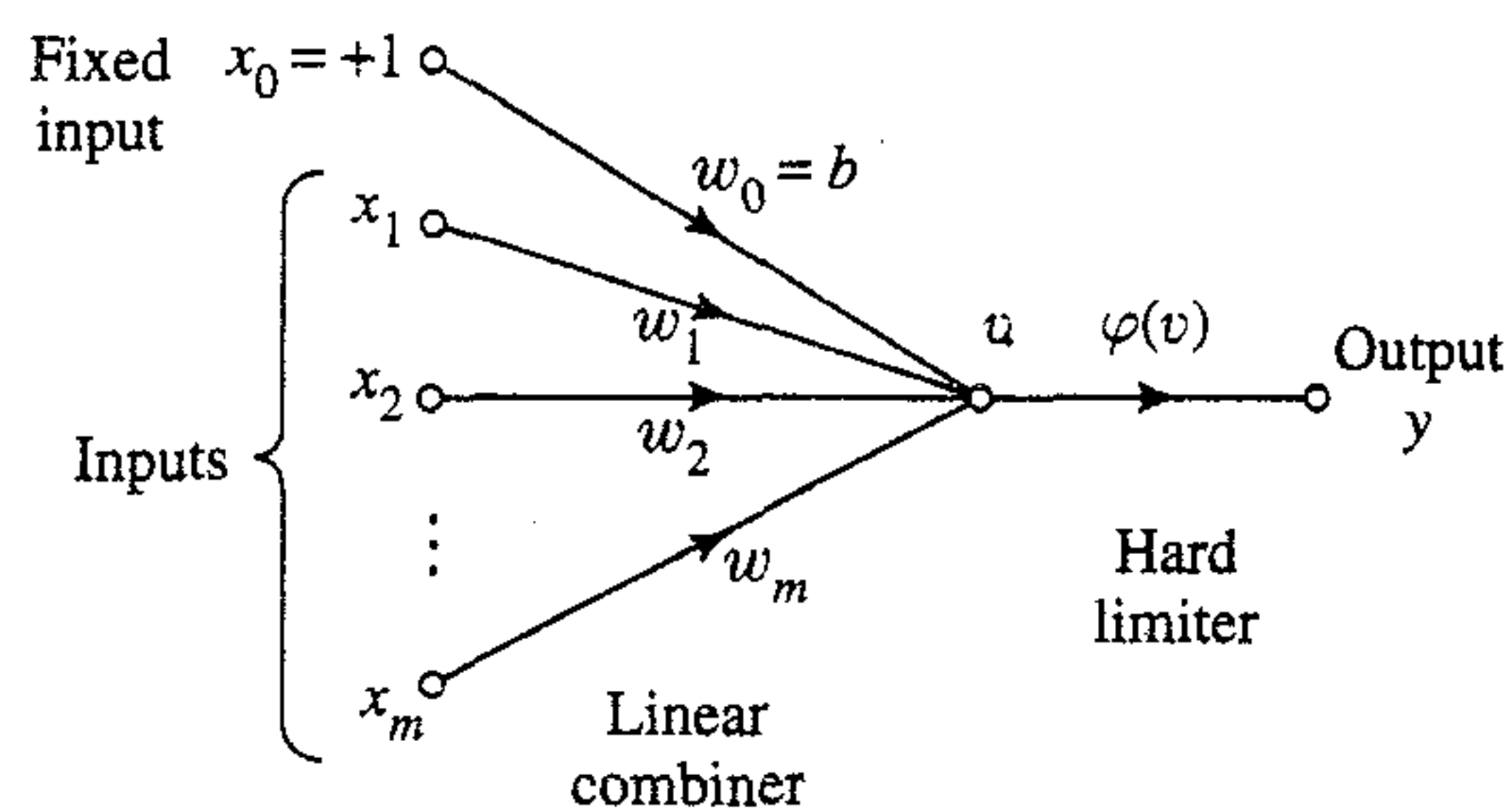
where  $n$  denotes the iteration step in applying the algorithm. Correspondingly we define the  $(m+1)$ -by-1 weight vector as

$$\mathbf{w}(n) = [b(n), w_1(n), w_2(n), \dots, w_m(n)]^T$$



**FIGURE 3.7** Illustration of the hyperplane (in this example, a straight line) as decision boundary for a two-dimensional, two-class pattern-classification problem.

**FIGURE 3.8** Equivalent signal-flow graph of the perceptron; dependence on time has been omitted for clarity.



**FIGURE 3.9** (a) A pair of linearly separable patterns. (b) A pair of non-linearly separable patterns.

Accordingly, the linear combiner output is written in the compact form

$$\begin{aligned} v(n) &= \sum_{i=0}^m w_i(n)x_i(n) \\ &= \mathbf{w}^T(n)\mathbf{x}(n) \end{aligned} \quad (3.52)$$

where  $w_0(n)$  represents the bias  $b(n)$ . For fixed  $n$ , the equation  $\mathbf{w}^T \mathbf{x} = 0$ , plotted in an  $m$ -dimensional space (plotted for some prescribed bias) with coordinates  $x_1, x_2, \dots, x_m$ , defines a hyperplane as the decision surface between two different classes of inputs.

For the perceptron to function properly, the two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  must be *linearly separable*. This, in turn, means that the patterns to be classified must be sufficiently separated from each other to ensure that the decision surface consists of a hyperplane. This requirement is illustrated in Fig. 3.9 for the case of a two-dimensional perceptron. In Fig. 3.9a, the two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are sufficiently separated from each other for us to draw a hyperplane (in this case a straight line) as the decision boundary. If, however, the two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are allowed to move too close to each other, as in Fig. 3.9b, they become nonlinearly separable, a situation that is beyond the computing capability of the perceptron.

Suppose then that the input variables of the perceptron originate from two linearly separable classes. Let  $\mathcal{X}_1$  be the subset of training vectors  $\mathbf{x}_1(1), \mathbf{x}_1(2), \dots$  that belong to class  $\mathcal{C}_1$ , and let  $\mathcal{X}_2$  be the subset of training vectors  $\mathbf{x}_2(1), \mathbf{x}_2(2), \dots$  that belong to class  $\mathcal{C}_2$ . The union of  $\mathcal{X}_1$  and  $\mathcal{X}_2$  is the complete training set  $\mathcal{X}$ . Given the sets



of vectors  $\mathcal{X}_1$  and  $\mathcal{X}_2$  to train the classifier, the training process involves the adjustment of the weight vector  $\mathbf{w}$  in such a way that the two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are linearly separable. That is, there exists a weight vector  $\mathbf{w}$  such that we may state

$$\begin{aligned} \mathbf{w}^T \mathbf{x} &> 0 \text{ for every input vector } \mathbf{x} \text{ belonging to class } \mathcal{C}_1 \\ \mathbf{w}^T \mathbf{x} &\leq 0 \text{ for every input vector } \mathbf{x} \text{ belonging to class } \mathcal{C}_2 \end{aligned} \quad (3.53)$$

In the second line of Eq. (3.53) we have arbitrarily chosen to say that the input vector  $\mathbf{x}$  belongs to class  $\mathcal{C}_2$  if  $\mathbf{w}^T \mathbf{x} = 0$ . Given the subsets of training vectors  $\mathcal{X}_1$  and  $\mathcal{X}_2$ , the training problem for the elementary perceptron is then to find a weight vector  $\mathbf{w}$  such that the two inequalities of Eq. (3.53) are satisfied.

The algorithm for adapting the weight vector of the elementary perceptron may now be formulated as follows:

1. If the  $n$ th member of the training set,  $\mathbf{x}(n)$ , is correctly classified by the weight vector  $\mathbf{w}(n)$  computed at the  $n$ th iteration of the algorithm, no correction is made to the weight vector of the perceptron in accordance with the rule:

$$\begin{aligned} \mathbf{w}(n+1) &= \mathbf{w}(n) && \text{if } \mathbf{w}^T \mathbf{x}(n) > 0 \text{ and } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_1 \\ \mathbf{w}(n+1) &= \mathbf{w}(n) && \text{if } \mathbf{w}^T \mathbf{x}(n) \leq 0 \text{ and } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_2 \end{aligned} \quad (3.54)$$

2. Otherwise, the weight vector of the perceptron is updated in accordance with the rule

$$\begin{aligned} \mathbf{w}(n+1) &= \mathbf{w}(n) - \eta(n)\mathbf{x}(n) && \text{if } \mathbf{w}^T(n)\mathbf{x}(n) > 0 \text{ and } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_2 \\ \mathbf{w}(n+1) &= \mathbf{w}(n) + \eta(n)\mathbf{x}(n) && \text{if } \mathbf{w}^T(n)\mathbf{x}(n) \leq 0 \text{ and } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_1 \end{aligned} \quad (3.55)$$

where the *learning-rate parameter*  $\eta(n)$  controls the adjustment applied to the weight vector at iteration  $n$ .

If  $\eta(n) = \eta > 0$ , where  $\eta$  is a constant independent of the iteration number  $n$ , we have a *fixed increment adaptation rule* for the perceptron.

In the sequel we first prove the convergence of a fixed increment adaptation rule for which  $\eta = 1$ . Clearly the value of  $\eta$  is unimportant, so long as it is positive. A value of  $\eta \neq 1$  merely scales the pattern vectors without affecting their separability. The case of a variable  $\eta(n)$  is considered later.

The proof is presented for the initial condition  $\mathbf{w}(0) = \mathbf{0}$ . Suppose that  $\mathbf{w}^T(n)\mathbf{x}(n) < 0$  for  $n = 1, 2, \dots$ , and the input vector  $\mathbf{x}(n)$  belongs to the subset  $\mathcal{X}_1$ . That is, the perceptron incorrectly classifies the vectors  $\mathbf{x}(1), \mathbf{x}(2), \dots$ , since the second condition of Eq. (3.53) is violated. Then, with the constant  $\eta(n) = 1$ , we may use the second line of Eq. (3.55) to write

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mathbf{x}(n) \quad \text{for } \mathbf{x}(n) \text{ belonging to class } \mathcal{C}_1 \quad (3.56)$$

Given the initial condition  $\mathbf{w}(0) = \mathbf{0}$ , we may iteratively solve this equation for  $\mathbf{w}(n+1)$ , obtaining the result

$$\mathbf{w}(n+1) = \mathbf{x}(1) + \mathbf{x}(2) + \dots + \mathbf{x}(n) \quad (3.57)$$

Since the classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are assumed to be linearly separable, there exists a solution  $\mathbf{w}_0$  for which  $\mathbf{w}_0^T \mathbf{x}(n) > 0$  for the vectors  $\mathbf{x}(1), \dots, \mathbf{x}(n)$  belonging to the subset  $\mathcal{X}_1$ . For a fixed solution  $\mathbf{w}_0$ , we may then define a positive number  $\alpha$  as

$$\alpha = \min_{\mathbf{x}(n) \in \mathcal{X}_1} \mathbf{w}_0^T \mathbf{x}(n) \quad (3.58)$$

Hence, multiplying both sides of Eq. (3.57) by the row vector  $\mathbf{w}_0^T$ , we get

$$\mathbf{w}_0^T \mathbf{w}(n+1) = \mathbf{w}_0^T \mathbf{x}(1) + \mathbf{w}_0^T \mathbf{x}(2) + \dots + \mathbf{w}_0^T \mathbf{x}(n)$$

Accordingly, in light of the definition given in Eq. (3.58), we have

$$\mathbf{w}_0^T \mathbf{w}(n+1) \geq n\alpha \quad (3.59)$$

Next we make use of an inequality known as the Cauchy–Schwarz inequality. Given two vectors  $\mathbf{w}_0$  and  $\mathbf{w}(n+1)$ , the *Cauchy–Schwarz inequality* states that

$$\|\mathbf{w}_0\|^2 \|\mathbf{w}(n+1)\|^2 \geq [\mathbf{w}_0^T \mathbf{w}(n+1)]^2 \quad (3.60)$$

where  $\|\cdot\|$  denotes the Euclidean norm of the enclosed argument vector, and the inner product  $\mathbf{w}_0^T \mathbf{w}(n+1)$  is a scalar quantity. We now note from Eq. (3.59) that  $[\mathbf{w}_0^T \mathbf{w}(n+1)]^2$  is equal to or greater than  $n^2 \alpha^2$ . From Eq. (3.60) we note that  $\|\mathbf{w}_0\|^2 \|\mathbf{w}(n+1)\|^2$  is equal to or greater than  $[\mathbf{w}_0^T \mathbf{w}(n+1)]^2$ . It follows therefore that

$$\|\mathbf{w}_0\|^2 \|\mathbf{w}(n+1)\|^2 \geq n^2 \alpha^2$$

or equivalently,

$$\|\mathbf{w}(n+1)\|^2 \geq \frac{n^2 \alpha^2}{\|\mathbf{w}_0\|^2} \quad (3.61)$$

We next follow another development route. In particular, we rewrite Eq. (3.56) in the form

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{x}(k) \quad \text{for } k = 1, \dots, n \quad \text{and} \quad \mathbf{x}(k) \in \mathcal{X}_1 \quad (3.62)$$

By taking the squared Euclidean norm of both sides of Eq. (3.62), we obtain

$$\|\mathbf{w}(k+1)\|^2 = \|\mathbf{w}(k)\|^2 + \|\mathbf{x}(k)\|^2 + 2\mathbf{w}^T(k)\mathbf{x}(k) \quad (3.63)$$

But, under the assumption that the perceptron incorrectly classifies an input vector  $\mathbf{x}(k)$  belonging to the subset  $\mathcal{X}_1$ , we have  $\mathbf{w}^T(k)\mathbf{x}(k) < 0$ . We therefore deduce from Eq. (3.63) that

$$\|\mathbf{w}(k+1)\|^2 \leq \|\mathbf{w}(k)\|^2 + \|\mathbf{x}(k)\|^2$$

or equivalently,

$$\|\mathbf{w}(k+1)\|^2 - \|\mathbf{w}(k)\|^2 \leq \|\mathbf{x}(k)\|^2, \quad k = 1, \dots, n \quad (3.64)$$

Adding these inequalities for  $k = 1, \dots, n$ , and invoking the assumed initial condition  $\mathbf{w}(0) = \mathbf{0}$ , we get the following inequality:

$$\begin{aligned}\|\mathbf{w}(n+1)\|^2 &\leq \sum_{k=1}^n \|\mathbf{x}(k)\|^2 \\ &\leq n\beta\end{aligned}\tag{3.65}$$

where  $\beta$  is a positive number defined by

$$\beta = \max_{\mathbf{x}(k) \in \mathcal{X}_1} \|\mathbf{x}(k)\|^2\tag{3.66}$$

Equation (3.65) states that the squared Euclidean norm of the weight vector  $\mathbf{w}(n+1)$  grows at most linearly with the number of iterations  $n$ .

The second result of Eq. (3.65) is clearly in conflict with the earlier result of Eq. (3.61) for sufficiently large values of  $n$ . Indeed, we can state that  $n$  cannot be larger than some value  $n_{\max}$  for which Eqs. (3.61) and (3.65) are both satisfied with the equality sign. That is,  $n_{\max}$  is the solution of the equation

$$\frac{n_{\max}^2 \alpha^2}{\|\mathbf{w}_0\|^2} = n_{\max} \beta$$

Solving for  $n_{\max}$ , given a solution vector  $\mathbf{w}_0$ , we find that

$$n_{\max} = \frac{\beta \|\mathbf{w}_0\|^2}{\alpha^2}\tag{3.67}$$

We have thus proved that for  $\eta(n) = 1$  for all  $n$ , and  $\mathbf{w}(0) = \mathbf{0}$ , and given that a solution vector  $\mathbf{w}_0$  exists, the rule for adapting the synaptic weights of the perceptron must terminate after at most  $n_{\max}$  iterations. Note also from Eqs. (3.58), (3.66), and (3.67) that there is *no* unique solution for  $\mathbf{w}_0$  or  $n_{\max}$ .

We may now state the *fixed-increment convergence theorem* for the perceptron as follows (Rosenblatt, 1962):

Let the subsets of training vectors  $\mathcal{X}_1$  and  $\mathcal{X}_2$  be linearly separable. Let the inputs presented to the perceptron originate from these two subsets. The perceptron converges after some  $n_0$  iterations, in the sense that

$$\mathbf{w}(n_0) = \mathbf{w}(n_0 + 1) = \mathbf{w}(n_0 + 2) = \dots$$

is a solution vector for  $n_0 \leq n_{\max}$ .

Consider next the *absolute error-correction procedure* for the adaptation of a single-layer perceptron, for which  $\eta(n)$  is variable. In particular, let  $\eta(n)$  be the smallest integer for which

$$\eta(n) \mathbf{x}^T(n) \mathbf{x}(n) > |\mathbf{w}^T(n) \mathbf{x}(n)|$$

With this procedure we find that if the inner product  $\mathbf{w}^T(n) \mathbf{x}(n)$  at iteration  $n$  has an incorrect sign, then  $\mathbf{w}^T(n+1) \mathbf{x}(n)$  at iteration  $n+1$  would have the correct sign. This suggests that if  $\mathbf{w}^T(n) \mathbf{x}(n)$  has an incorrect sign, we may modify the training sequence at iteration  $n+1$  by setting  $\mathbf{x}(n+1) = \mathbf{x}(n)$ . In other words, each pattern is presented repeatedly to the perceptron until that pattern is classified correctly.

Note also that the use of an initial value  $\mathbf{w}(0)$  different from the null condition merely results in a decrease or increase in the number of iterations required to converge,



depending on how  $\mathbf{w}(0)$  relates to the solution  $\mathbf{w}_0$ . Regardless of the value assigned to  $\mathbf{w}(0)$ , the perceptron is assured of convergence.

In Table 3.2 we present a summary of the *perceptron convergence algorithm* (Lippmann, 1987). The symbol  $\text{sgn}(\cdot)$ , used in step 3 of the table for computing the actual response of the perceptron, stands for the *signum function*:

$$\text{sgn}(v) = \begin{cases} +1 & \text{if } v > 0 \\ -1 & \text{if } v < 0 \end{cases} \quad (3.68)$$

We may thus express the *quantized response*  $y(n)$  of the perceptron in the compact form

$$y(n) = \text{sgn}(\mathbf{w}^T(n)\mathbf{x}(n)) \quad (3.69)$$

---

**TABLE 3.2** Summary of the Perceptron Convergence Algorithm

---

*Variables and Parameters:*

$\mathbf{x}(n)$  =  $(m+1)$ -by-1 input vector

$$= [+1, x_1(n), x_2(n), \dots, x_m(n)]^T$$

$\mathbf{w}(n)$  =  $(m+1)$ -by-1 weight vector

$$= [b(n), w_1(n), w_2(n), \dots, w_m(n)]^T$$

$b(n)$  = bias

$y(n)$  = actual response (quantized)

$d(n)$  = desired response

$\eta$  = learning-rate parameter, a positive constant less than unity

1. *Initialization.* Set  $\mathbf{w}(0) = \mathbf{0}$ . Then perform the following computations for time step  $n = 1, 2, \dots$

2. *Activation.* At time step  $n$ , activate the perceptron by applying continuous-valued input vector  $\mathbf{x}(n)$  and desired response  $d(n)$ .

3. *Computation of Actual Response.* Compute the actual response of the perceptron:

$$y(n) = \text{sgn}[\mathbf{w}^T(n)\mathbf{x}(n)]$$

where  $\text{sgn}(\cdot)$  is the signum function.

4. *Adaptation of Weight Vector.* Update the weight vector of the perceptron:

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \eta[d(n) - y(n)]\mathbf{x}(n)$$

where

$$d(n) = \begin{cases} +1 & \text{if } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_1 \\ -1 & \text{if } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_2 \end{cases}$$

5. *Continuation.* Increment time step  $n$  by one and go back to step 2.

---

Notice that the input vector  $\mathbf{x}(n)$  is an  $(m + 1)$ -by-1 vector whose first element is fixed at +1 throughout the computation. Correspondingly, the weight vector  $\mathbf{w}(n)$  is an  $(m + 1)$ -by-1 vector whose first element equals the bias  $b(n)$ . One other important point in Table 3.2: We have introduced a *quantized desired response*  $d(n)$ , defined by

$$d(n) = \begin{cases} +1 & \text{if } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_1 \\ -1 & \text{if } \mathbf{x}(n) \text{ belongs to class } \mathcal{C}_2 \end{cases} \quad (3.70)$$

Thus, the adaptation of the weight vector  $\mathbf{w}(n)$  is summed up nicely in the form of the *error-correction learning rule*:

$$\mathbf{w}(n + 1) = \mathbf{w}(n) + \eta[d(n) - y(n)]\mathbf{x}(n) \quad (3.71)$$

where  $\eta$  is the *learning-rate parameter*, and the difference  $d(n) - y(n)$  plays the role of an *error signal*. The learning-rate parameter is a positive constant limited to the range  $0 < \eta \leq 1$ . When assigning a value to it inside this range, we must keep in mind two conflicting requirements (Lippmann, 1987):

- *Averaging* of past inputs to provide stable weight estimates, which requires a small  $\eta$
- *Fast adaptation* with respect to real changes in the underlying distributions of the process responsible for the generation of the input vector  $\mathbf{x}$ , which requires a large  $\eta$

### 3.10 RELATION BETWEEN THE PERCEPTRON AND BAYES CLASSIFIER FOR A GAUSSIAN ENVIRONMENT

The perceptron bears a certain relationship to a classical pattern classifier known as the Bayes classifier. When the environment is Gaussian, the Bayes classifier reduces to a linear classifier. This is the same form taken by the perceptron. However, the linear nature of the perceptron is *not* contingent on the assumption of Gaussianity. In this section we study this relationship, and thereby develop further insight into the operation of the perceptron. We begin the discussion with a brief review of the Bayes classifier.

#### Bayes Classifier

In the *Bayes classifier* or *Bayes hypothesis testing procedure*, we minimize the *average risk*, denoted by  $\mathcal{R}$ . For a two-class problem, represented by classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , the average risk is defined by Van Trees (1968):

$$\begin{aligned} \mathcal{R} = & c_{11}p_1 \int_{\mathcal{X}_1} f_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} + c_{22}p_2 \int_{\mathcal{X}_2} f_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} \\ & + c_{21}p_1 \int_{\mathcal{X}_2} f_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} + c_{12}p_2 \int_{\mathcal{X}_1} f_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} \end{aligned} \quad (3.72)$$

Equations (1) and (2) provide two useful rules for the differentiation of a real-valued function with respect to a vector.

### 3. Positive definite matrix

An  $m$ -by- $m$  matrix  $\mathbf{R}$  is said to be nonnegative definite if it satisfies the condition

$$\mathbf{a}^T \mathbf{R} \mathbf{a} \geq 0 \quad \text{for any vector } \mathbf{a} \in \mathbb{R}^m$$

If this condition is satisfied with the inequality sign, the matrix  $\mathbf{R}$  is said to be positive definite.

An important property of a positive definite matrix  $\mathbf{R}$  is that it is *nonsingular*, that is, the inverse matrix  $\mathbf{R}^{-1}$  exists.

Another important property of a positive definite matrix  $\mathbf{R}$  is that its eigenvalues, or roots of the characteristic equation

$$\det(\mathbf{R}) = 0$$

are all positive.

### 4. Robustness

The  $H^\infty$  criterion is due to Zames (1981), and it is developed in Zames and Francis (1983). The criterion is discussed in Doyle et al. (1989), Green and Limebeer (1995), and Hassibi et al. (1998).

5. To overcome the limitations of the LMS algorithm, namely, slow rate of convergence and sensitivity to variations in the condition number of the correlation matrix  $\mathbf{R}_x$ , we may use the *recursive least-squares (RLS) algorithm*, which follows from a recursive implementation of the linear least-squares filter described in Section 3.4. The RLS algorithm is a special case of the Kalman filter, which is known to be the optimum linear filter for a nonstationary environment. Most importantly, the Kalman filter exploits all past data extending up to and including the time instant at which the computations are made. For more details about the RLS algorithm and its relationship to the Kalman filter, see Haykin (1996). The Kalman filter is discussed in Chapter 15.

## PROBLEMS

### Unconstrained optimization

- 3.1 Explore the method of steepest descent involving a single weight  $w$  by considering the following cost function:

$$\mathcal{E}(w) = \frac{1}{2} \sigma^2 - r_{xd} w + \frac{1}{2} r_x w^2$$

where  $\sigma^2$ ,  $r_{xd}$ , and  $r_x$  are constants.

- 3.2 Consider the cost function

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2} \sigma^2 - \mathbf{r}_{xd}^T \mathbf{w} + \frac{1}{2} \mathbf{w}^T \mathbf{R}_x \mathbf{w}$$

where  $\sigma^2$  is some constant, and

$$\mathbf{r}_{xd} = \begin{bmatrix} 0.8182 \\ 0.354 \end{bmatrix}$$

$$\mathbf{R}_x = \begin{bmatrix} 1 & 0.8182 \\ 0.8182 & 1 \end{bmatrix}$$

- (a) Find the optimum value  $\mathbf{w}^*$  for which  $\mathcal{E}(\mathbf{w})$  reaches its minimum value.



(b) Use the method of steepest descent to compute  $\mathbf{w}^*$  for the following two values of learning-rate parameter:

(i)  $\eta = 0.3$

(ii)  $\eta = 1.0$

For each case, plot the trajectory traced by the evolution of the weight vector  $\mathbf{w}(n)$  in the  $\mathbf{W}$ -plane.

*Note:* The trajectories obtained for cases (i) and (ii) of part (b) should correspond to the pictures displayed in Fig. 3.2.

3.3 Consider the cost function of Eq. (3.24) that represents a modified form of the sum of error squares defined in Eq. (3.17). Show that the application of the Gauss–Newton method to Eq. (3.24) yields the weight-update described in Eq. (3.23).

### LMS Algorithm

3.4 The correlation matrix  $\mathbf{R}_x$  of the input vector  $\mathbf{x}(n)$  in the LMS algorithm is defined by

$$\mathbf{R}_x = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

Define the range of values for the learning-rate parameter  $\eta$  of the LMS algorithm for it to be convergent in the mean square.

3.5 The *normalized LMS algorithm* is described by the following recursion for the weight vector:

$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \frac{\eta}{\|\mathbf{x}(n)\|^2} e(n)\mathbf{x}(n)$$

where  $\eta$  is a positive constant and  $\|\mathbf{x}(n)\|$  is the Euclidean norm of the input vector  $\mathbf{x}(n)$ . The error signal  $e(n)$  is defined by

$$e(n) = d(n) - \hat{\mathbf{w}}^T(n)\mathbf{x}(n)$$

where  $d(n)$  is the desired response. For the normalized LMS algorithm to be convergent in the mean square, show that

$$0 < \eta < 2$$

3.6 The LMS algorithm is used to implement the generalized sidelobe canceler shown in Fig. 2.16. Set up the equations that define the operation of this system, assuming the use of a single neuron for the neural network.

3.7 Consider a linear predictor with its input vector made up of the samples  $x(n-1)$ ,  $x(n-2)$ , ...,  $x(n-m)$ , where  $m$  is the prediction order. The requirement is to use the LMS algorithm to make a prediction  $\hat{x}(n)$  of the input sample  $x(n)$ . Set up the recursions that may be used to compute the tap weight  $w_1, w_2, \dots, w_m$  of the predictor.

3.8 The ensemble-averaged counterpart to the sum of error squares viewed as a cost function is the mean-square value of the error signal:

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} E[e^2(n)] \\ &= \frac{1}{2} E[(d(n) - \mathbf{x}^T(n)\mathbf{w})^2] \end{aligned}$$

- (a) Assuming that the input vector  $\mathbf{x}(n)$  and desired response  $d(n)$  are drawn from a stationary environment, show that

$$J(\mathbf{w}) = \frac{1}{2} \sigma_d^2 - \mathbf{r}_{xd}^T \mathbf{w} + \frac{1}{2} \mathbf{w}^T \mathbf{R}_x \mathbf{w}$$

where

$$\sigma_d^2 = E[d^2(n)]$$

$$\mathbf{r}_{xd} = E[\mathbf{x}(n) d(n)]$$

$$\mathbf{R}_x = E[\mathbf{x}(n) \mathbf{x}^T(n)]$$

- (b) For this cost function, show that the gradient vector and Hessian matrix of  $J(\mathbf{w})$  are as follows, respectively:

$$\mathbf{g} = -\mathbf{r}_{xd} + \mathbf{R}_x \mathbf{w}$$

$$\mathbf{H} = \mathbf{R}_x$$

- (c) In the *LMS/Newton algorithm*, the gradient vector  $\mathbf{g}$  is replaced by its instantaneous value (Widrow and Stearns, 1985). Show that this algorithm, incorporating a learning-rate parameter  $\eta$ , is described by

$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \eta \mathbf{R}_x^{-1} \mathbf{x}(n) (d(n) - \mathbf{x}^T(n) \mathbf{w}(n))$$

The inverse of the correlation matrix  $\mathbf{R}_x$ , assumed to be positive definite, is calculated ahead of time.

- 3.9** In this problem we revisit the correlation matrix memory discussed in Section 2.11. A shortcoming of this memory is that when a key pattern  $\mathbf{x}_j$  is presented to it, the actual response  $\mathbf{y}$  produced by the memory may not be close enough (in a Euclidean sense) to the desired response (memorized pattern)  $\mathbf{y}_j$  for the memory to associate perfectly. This shortcoming is inherited from the use of Hebbian learning that has no provision for feedback from the output to the input. As a remedy for this shortcoming, we may incorporate an error-correction mechanism into the design of the memory, forcing it to associate properly (Anderson, 1983).

Let  $\hat{\mathbf{M}}(n)$  denote the memory matrix learned at iteration  $n$  of the error-correction learning process. The memory matrix  $\hat{\mathbf{M}}(n)$  learns the information represented by the associations:

$$\mathbf{x}_k \rightarrow \mathbf{y}_k, \quad k = 1, 2, \dots, q$$

- (a) Adapting the LMS algorithm for the problem at hand, show that the updated value of the memory matrix is defined by

$$\hat{\mathbf{M}}(n+1) = \hat{\mathbf{M}}(n) + \eta [\mathbf{y}_k - \hat{\mathbf{M}}(n) \mathbf{x}_k] \mathbf{x}_k^T$$

where  $\eta$  is the learning-rate parameter.

- (b) For autoassociation,  $\mathbf{y}_k = \mathbf{x}_k$ . For this special case, show that as the number of iterations,  $n$ , approaches infinity, the memory autoassociates perfectly, as shown by

$$\mathbf{M}(\infty) \mathbf{x}_k = \mathbf{x}_k, \quad k = 1, 2, \dots, q$$

- (c) The result described in part (b) may be viewed as an *eigenvalue problem*. In that context,  $\mathbf{x}_k$  represents an eigenvector of  $\mathbf{M}(\infty)$ . What are the eigenvalues of  $\mathbf{M}(\infty)$ ?

- 3.10** In this problem we investigate the effect of bias on the condition number of a correlation matrix, and therefore the performance of the LMS algorithm.

Consider a random vector  $\mathbf{X}$  with covariance matrix

$$\mathbf{C} = \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

and mean vector

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

- (a) Calculate the condition number of the covariance matrix  $\mathbf{C}$ .
  - (b) Calculate the condition number of the correlation matrix  $\mathbf{R}$ .
- Comment on the effect of the bias  $\boldsymbol{\mu}$  on the performance of the LMS algorithm.

### Rosenblatt's Perceptron

- 3.11** In this problem, we consider another method for deriving the update equation for Rosenblatt's perceptron. Define the *perceptron criterion function* (Duda and Hart, 1973):

$$J_p(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}(\mathbf{w})} (-\mathbf{w}^T \mathbf{x})$$

where  $\mathcal{X}(\mathbf{w})$  is the set of samples misclassified by the choice of weight vector  $\mathbf{w}$ . Note that  $J_p(\mathbf{w})$  is defined as zero if there are no misclassified samples, and the output is misclassified if  $\mathbf{w}_x^T \leq 0$ .

- (a) Demonstrate geometrically that  $J_p(\mathbf{w})$  is proportional to the sum of Euclidean distances from the misclassified samples to the decision boundary.
- (b) Determine the gradient of  $J_p(\mathbf{w})$  with respect to the weight vector  $\mathbf{w}$ .
- (c) Using the result obtained in part (b), show that the weight-update for the perceptron is:

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \eta(n) \sum_{\mathbf{x} \in \mathcal{X}(\mathbf{w}(n))} \mathbf{x}$$

where  $\mathcal{X}(\mathbf{w}(n))$  is the set of samples misclassified by the use of weight vector  $\mathbf{w}(n)$ , and  $\eta(n)$  is the learning-rate parameter. Show that this result, for the case of a single-sample correction, is basically the same as that described by Eqs. (3.54) and (3.55).

- 3.12** Verify that Eqs. (3.68)–(3.71), summarizing the perceptron convergence algorithm, are consistent with Eqs. (3.54) and (3.55).
- 3.13** Consider two one-dimensional, Gaussian-distributed classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  that have a common variance equal to 1. Their mean values are

$$\mu_1 = -10$$

$$\mu_2 = +10$$

These two classes are essentially linearly separable. Design a classifier that separates these two classes.

- 3.14** Suppose that in the signal-flow graph of the perceptron shown in Fig. 3.6 the hard limiter is replaced by the sigmoidal nonlinearity:

$$\varphi(v) = \tanh\left(\frac{v}{2}\right)$$



where  $v$  is the induced local field. The classification decisions made by the perceptron are defined as follows:

*Observation vector  $\mathbf{x}$  belongs to class  $\mathcal{C}_1$  if the output  $y > \theta$  where  $\theta$  is a threshold; otherwise,  $\mathbf{x}$  belongs to class  $\mathcal{C}_2$ .*

Show that the decision boundary so constructed is a hyperplane.

- 3.15** (a) The perceptron may be used to perform numerous logic functions. Demonstrate the implementation of the binary logic functions AND, OR, and COMPLEMENT.  
 (b) A basic limitation of the perceptron is that it cannot implement the EXCLUSIVE OR function. Explain the reason for this limitation.
- 3.16** Equations (3.86) and (3.87) define the weight vector and bias of the Bayes classifier for a Gaussian environment. Determine the composition of this classifier for the case when the covariance matrix  $\mathbf{C}$  is defined by

$$\mathbf{C} = \sigma^2 \mathbf{I}$$

where  $\sigma^2$  is a constant.